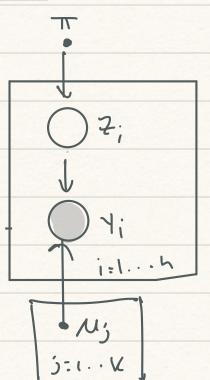
# The expectation-maximization (EM) algorithm

## Gaussian mixture model:

#### **Model parameters**



"Complete data likelihood" or "joint probability"

"Likelihood" or "marginal likelihood / evidence"

"Maximum likelihood estimation" or "type-II MLE"

Note: There are two parallel interpretations. In some models, we will interpret Z as "missing data", in other models we will think of Z as "latent variables". When Z are "missing data", then P(Z, Y) is the "complete data likelihood" and P(Y) is just the "likelihood". When Z are "latent variables" then P(Y, Z) is the "joint", P(Y | Z) is the "likelihood", and P(Y) is the "marginal likelihood", and maximizing P(Y) is doing "type II MLE" (i.e., maximum marginal likelihood estimation). But this is all just semantics / interpretation; the math is what's important. We will call P(Y) the marginal likelihood or evidence, but some literature calls it the likelihood.

### Evidence lower bound (ELBO)

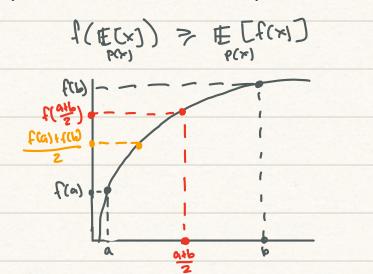
For any distribution q(Z) that is absolutely continuous wrt P(Y, Z):

= 
$$\log \mathbb{E} \left[ \frac{P_{\Theta}(Y|Z)}{q(Z)} \right]$$
 If  $p(Z=z, Y) = 0$  then  $q(Z=z) = 0$ , for all z.

Switching the expectation and log, we then have by Jensen's inequality:

This is a function of q and the parameters, which we will call the ELBO:

Jensen's inequality: For any concave function f(x), the expectation of f(x) lower bounds f(expectation of x).

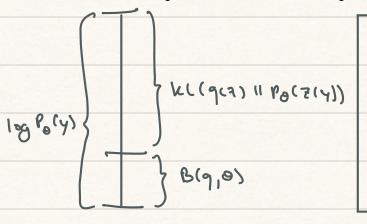


Let's do some ELBO surgery:

p(y) does not depend on z

this is a (negative) KL-divergence

Since the KL divergence must be non-negative, the ELBO must lower bound the evidence.



The Kullback-Leibler divergence gives a notion of similarity between two distributions q and p over the same support. It is not a metric (not symmetric) and is 0 only if p=q.

- · KL(911P) 70
- · KL(9(1p) =0 IFF 9=P
- $\text{KL}(9||p) \neq \text{KL}(p||9)$ "not necessarily"

### The EM algorithm

Choal: 
$$\hat{\theta}$$
 to argument  $p_{\theta}(y)$ 

initialize  $\theta_{\theta}$ 

for  $m=1,2,...$  until convergence in  $\theta_{m}$ :

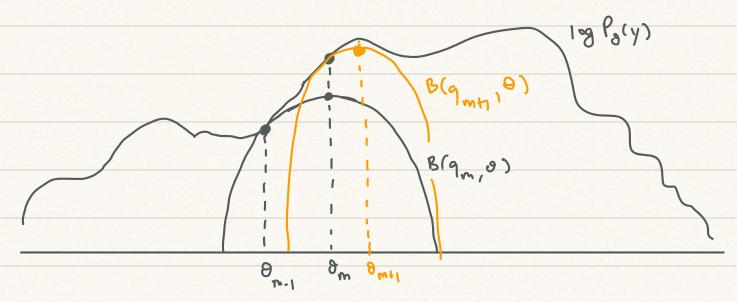
(E-step)  $q_{m} = \underset{q}{\operatorname{argmax}} B(q, \theta_{m-1})$ 

(M-step)  $\theta_{m} = \underset{q}{\operatorname{argmax}} B(q_{m}, \theta)$ 

return  $\theta_{m}$ 

- This will converge to a local mode of p(y).
- Random restarts are a good idea (to help find better local modes).
- The ELBO is tight after the E-step and always increases.

### EM as a minorize-maximize algorithm



"minorize": to construct a surrogate function that touches the function-tomaximize at the current point, and lies below it everywhere else.

"maximize": to find a new point that maximizes the surrogate function.

$$E - step$$

$$q_{m} = avg_{max} B(q_{1} \sigma_{m-1})$$

$$= avg_{mex} - kel(q_{1} 11 l_{0}(2_{1} y)) + locat$$

$$= P_{0m-1}(2_{1} y)$$

$$M - slep$$

$$g_{m} = avg_{max} B(q_{m-1} \theta)$$

$$\frac{\partial}{\partial m} = \alpha v y n \alpha x \quad \mathcal{B}(q_m, \theta)$$

$$= \alpha v y n \alpha x \quad \mathcal{E}\left[\log P_{\theta}(1, 2)\right] - \mathcal{E}\left[\log P_{\theta}(1, 2)\right]$$

$$= \alpha v y m \alpha x \quad \mathcal{E}\left[\log P_{\theta}(1, 2)\right]$$

$$= \alpha v y m \alpha x \quad \mathcal{E}\left[\log P_{\theta}(1, 2)\right]$$

$$= \alpha v y m \alpha x \quad \mathcal{E}\left[\log P_{\theta}(1, 2)\right]$$

Plugging in the E-step, a more compact way to write the algorithm is as follows:

Where now we have something resembling the transition operator in MCMC (i.e., an objective function that depends on the current iterate.

### EM for the Gaussian mixture model

E. slep: 
$$q(z) = P(z|y) = TT P(z; |y)$$

$$P_{s}(z_{i-j}|y) = r_{ij} = \frac{\mathcal{N}(y_{i}; \mu_{i}, 1) \pi_{j}}{\sum \mathcal{N}(y_{i}; \mu_{i}', 1) \pi_{j}'}$$

M-step:

This breaks into independent optimization problems:

The optimization problem for the means further breaks into K independent ones:

This is just a weighted least-squares problem:

$$M_{j}^{m} = \frac{\sum_{i} r_{ij}^{m} Y_{i}}{\sum_{i} r_{ij}^{m}} \quad A_{j}$$

Meanwhile, this is just a crossentropy loss problem:

## EM for the expfam mixture model

Now say the class-conditional likelihood is any exponential family...

Again, the M-step breaks into separate optimization problems, and the problem for pi is the same as before, so we only focus on estimating the natural parameters of the likelihood:

These again separate into independent optimization. Each one looks like this:

Taking the gradient of the objective and setting it to zero...

$$\frac{\nabla \gamma_{ij}}{\nabla \gamma_{ij}} = \frac{\nabla \gamma_{ij} + (\gamma_{ij})}{\nabla \gamma_{ij}} = \frac{\nabla \gamma_{ij} + (\gamma_{ij})}{\nabla \gamma_{ij}}$$

Note that for any (minimal) exponential family, we have the fact that the mapping from the natural parameter to the mean parameter is exactly the gradient of the log normalizer:

Therefore the maximizer is the weighted average of the sufficient statistics passed through the inverse of the mean map:

ms M; = M-1 ( - Zrist(Yi))

### EM for MAP estimation

Let's now put priors on the parameters and do maximum-a-posteriori (MAP) estimation:

The ELBO now becomes:

Note that the q distribution is still only over the latent variables Z:

The M-step becomes:

For example, let's add the following priors to the Gaussian mixture model:

The optimization problem for one of the means is just the same as before, but now adding the log prior to the objective function:

This is still just a weighted least squares problem:

Now consider the optimization problem for Pi; the objective function is the same as the one before but now adding the log Dirichlet prior:

TI 
$$\leftarrow$$
 augmax  $\sum \sum V_{ij} \log T_i + \log D \cdot \nu (T, d)$ 

TI  $\leftarrow$  augmax.  $\sum \sum V_{ij} \log T_j + \sum (d_{j-1}) \log T_j$ 

You can solve this using Lagrange multipliers (a good exercise). The solution is:

T; 
$$\leftarrow \frac{\alpha_{j} + \sum_{i \neq j-1}}{\sum_{j} (\alpha_{j} + \sum_{i \neq j-1})}$$

Last time, we derived the complete conditionals for both of these parameters with the same priors. The complete conditionals were:

$$P(M; |-) = M(M; m; x; )$$
 $m_{j} = \lambda_{j}^{-1}(M_{0} + \sum_{j=1}^{j} y_{j})$ 
 $\lambda_{j} = \lambda_{0} + \sum_{j=1}^{j} y_{j}$ 
 $P(\pi |-) = Dr(\pi; \lambda_{n})$ 
 $\lambda_{n,j} = \alpha_{j} + \sum_{j=1}^{j} y_{j}$ 

The M-step update for the Gaussian looks exactly like the mean of its complete conditional, except that the indicators have been replaced with expectations of the indicators (under q(Z)).

Similarly, the M-step update for Pi is exactly the mode of the Pirichlet complete conditional, except again the indicators have been replaced with expectations of indicators.

This is not a coincidence!

## EM and exponential families: a match made in heaven

An alternative way to go about deriving the M-step update is as follows. First note that:

And further note that:

Therefore:

In our model (and in many) the parameters are conditionally independent in their posterior, so this breaks apart into separate optimization problems (as before). If we have all conditionally conjugate priors, then the complete conditional for every parameter can be written as:

where f(...) is an exponential family, and lambda\_n is its posterior natural parameter, which is a function of Z and Y.

And so the optimal update to eta is the mode of an exponential family f(...), which is the same as the complete conditional except replacing the natural parameter with its expectation under q(Z)!

