

## EM Algorithm

---

An iterative optimization strategy motivated by a notion of missingness and by consideration of the conditional distribution of what is missing given what is observed.

Can be very simple to implement. Can reliably find an optimum through stable, uphill steps.

Difficult likelihoods often arise when data are missing. EM simplifies such problems. In fact, the ‘missing data’ may not truly be missing: they may be only a conceptual ploy to exploit the EM simplification!

## Notation

**X** : Observed variables.

**Z** : Missing or latent variables.

**Y** : Complete data  $Y = (\mathbf{X}, \mathbf{Z})$ .

In Bayesian settings, **X**, **Z**, and **Y** often refer to sets of parameters, rather than data.

Suppose we seek to maximize  $L(\boldsymbol{\theta}|\mathbf{x})$  with respect to  $\boldsymbol{\theta}$ .

Define  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$  to be the expectation of the joint log likelihood for the complete data, conditional on the observed data  $\mathbf{X} = \mathbf{x}$ . Namely

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = \text{E} \left\{ \log L(\boldsymbol{\theta}|\mathbf{Y}) \mid \mathbf{x}, \boldsymbol{\theta}^{(t)} \right\} \quad (49)$$

$$= \text{E} \left\{ \log f_{\mathbf{Y}}(\mathbf{y}|\boldsymbol{\theta}) \mid \mathbf{x}, \boldsymbol{\theta}^{(t)} \right\} \quad (50)$$

$$= \int [\log f_{\mathbf{Y}}(\mathbf{y}|\boldsymbol{\theta})] f_{\mathbf{Z}|\mathbf{X}}(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^{(t)}) d\mathbf{z} \quad (51)$$

where (51) emphasizes that  $\mathbf{Z}$  is the only random part of  $\mathbf{Y}$  once we are given  $\mathbf{X} = \mathbf{x}$ .

# The EM Algorithm

---

Start with  $\theta^{(0)}$ . Then

1. E step: Compute  $Q(\theta|\theta^{(t)})$ .
2. M step: Maximize  $Q(\theta|\theta^{(t)})$  with respect to  $\theta$ . Set  $\theta^{(t+1)}$  equal to the maximizer of  $Q$ .
3. Return to the E step unless a stopping criterion has been met.

### Trivial example:

$Y_1, Y_2 \sim \text{i.i.d. Exp}(\theta)$  with  $y_1 = 5$  observed but  $y_2$  missing.

The complete data log likelihood function is

$$\log\{L(\theta|\mathbf{y})\} = \log\{f_{\mathbf{Y}}(\mathbf{y}|\theta)\} = 2 \log\{\theta\} - \theta y_1 - \theta y_2. \quad (52)$$

Thus

$$Q(\theta|\theta^{(t)}) = 2 \log\{\theta\} - 5\theta - \theta/\theta^{(t)} \quad (53)$$

since  $E\{Y_2|y_1, \theta^{(t)}\} = E\{Y_2|\theta^{(t)}\} = 1/\theta^{(t)}$  follows from independence.

The maximizer of  $Q(\theta|\theta^{(t)})$  is the root of  $2/\theta - 5 - 1/\theta^{(t)} = 0$ . Thus  $\theta^{(t+1)} = \frac{2\theta^{(t)}}{5\theta^{(t)}+1}$ .

Converges quickly to  $\hat{\theta} = 0.2$ .

The E step and M step do not need to be re-derived at each iteration

This example is not realistic. Easy analytic solution. Taking the required expectation is trickier in real applications because one needs to know the conditional distribution of the complete data given the missing data.

## Example: Peppered moths

---

Wing color determined by a single gene with three possible alleles which we denote C, I, and T. C is dominant to I, and T is recessive to I. Thus genotypes CC, CI, and CT result in the *carbonaria* phenotype having solid black coloring. The TT genotype results in the *typica* phenotype with light-colored wings. The II and IT genotypes produce an intermediate *insularia* phenotype with mottled wings.

There are six possible genotypes, but only three phenotypes are measurable in field work.

In UK and USA, *carbonaria* nearly replaced the paler phenotypes in areas affected by coal-fired industries. This change in allele frequencies in the population is cited as an instance when we may observe micro-evolution occurring on a human time scale. The theory (supported by experiments) is that “differential predation by birds on moths that are variously conspicuous against backgrounds of different reflectance” induces selectivity that favors *carbonaria* in times/regions where sooty, polluted conditions reduce the reflectance of the surface of tree bark on which the moths rest. When improved environmental standards reduced pollution, prevalence of the lighter-colored phenotypes increased and *carbonaria* prevalence plummeted.

## How do we estimate allele frequencies from phenotype counts?

Under Hardy-Weinberg, if the allele frequencies in the population are  $p_C$ ,  $p_I$ , and  $p_T$ , then the genotype frequencies should be  $p_C^2$ ,  $2p_Cp_I$ ,  $2p_Cp_T$ ,  $p_I^2$ ,  $2p_Ip_T$ , and  $p_T^2$ , for genotypes CC, CI, CT, II, IT, and TT, respectively.

Capture  $n$  moths, of which there are  $n_C$ ,  $n_I$ , and  $n_T$  of the *carbonaria*, *insularia*, and *typica* phenotypes, respectively. Thus  $n = n_C + n_I + n_T$ . If we knew the genotype of each moth rather than merely its phenotype, we could generate genotype counts  $n_{CC}$ ,  $n_{CI}$ ,  $n_{CT}$ ,  $n_{II}$ ,  $n_{IT}$ , and  $n_{TT}$ . Allele frequencies could then be easily tabulated: each moth with genotype CI contributes one C allele and one I allele, whereas a II moth contributes two I alleles.

**Observed data:**  $\mathbf{x} = (n_C, n_I, n_T)$

**Complete data:**  $\mathbf{y} = (n_{CC}, n_{CI}, n_{CT}, n_{II}, n_{IT}, n_{TT})$

**Complete data log likelihood is multinomial:**

$$\begin{aligned} \log\{f_{\mathbf{Y}}(\mathbf{y}|\mathbf{p})\} &= n_{CC} \log\{p_C^2\} + n_{CI} \log\{2p_C p_I\} + n_{CT} \log\{2p_C p_T\} \\ &\quad + n_{II} \log\{p_I^2\} + n_{IT} \log\{2p_I p_T\} + n_{TT} \log\{p_T^2\} \\ &\quad + \log \binom{n}{n_{CC} \ n_{CI} \ n_{CT} \ n_{II} \ n_{IT} \ n_{TT}}. \end{aligned} \quad (54)$$

**Many-to-one mapping from  $\mathbf{Y}$  to  $\mathbf{X}$ :**  $(n_C, n_I, n_T) = (n_{CC} + n_{CI} + n_{CT}, n_{II} + n_{IT}, n_{TT})$ .



### E-Step:

Conditional on  $n_C$  and a parameter vector  $\mathbf{p}^{(t)} = (p_C^{(t)}, p_I^{(t)})$ , the latent counts for the three *carbonaria* genotypes have a three-cell multinomial distribution with count parameter  $n_C$  and cell probabilities proportional to  $(p_C^{(t)})^2$ ,  $2p_C^{(t)}p_I^{(t)}$ , and  $2p_C^{(t)}p_T^{(t)}$ .

Same reasoning for other genotypes.

Thus

$$E\{N_{CC}|n_C, n_I, n_T, \mathbf{p}^{(t)}\} = n_{CC}^{(t)} = \frac{n_C(p_C^{(t)})^2}{(p_C^{(t)})^2 + 2p_C^{(t)}p_I^{(t)} + 2p_C^{(t)}p_T^{(t)}} \quad (55)$$

and so forth.

Thus

$$Q(\mathbf{p}|\mathbf{p}^{(t)}) = n_{CC}^{(t)} \log\{p_C^2\} + n_{CI}^{(t)} \log\{2p_C p_I\} + n_{CT}^{(t)} \log\{2p_C p_T\} \\ + n_{II}^{(t)} \log\{p_I^2\} + n_{IT}^{(t)} \log\{2p_I p_T\} + n_{TT} \log\{p_T^2\} + k(n_C, n_I, n_T, \mathbf{p}^{(t)}). \quad (56)$$

## M-Step:

Differentiating with respect to  $p_c$  and  $p_i$  yields

$$\frac{dQ(\mathbf{p}|\mathbf{p}^{(t)})}{dp_c} = \frac{2n_{cc}^{(t)} + n_{ci}^{(t)} + n_{ct}^{(t)}}{p_c} - \frac{2n_{tt}^{(t)} + n_{ct}^{(t)} + n_{it}^{(t)}}{1 - p_c - p_i} \quad (57)$$

$$\frac{dQ(\mathbf{p}|\mathbf{p}^{(t)})}{dp_i} = \frac{2n_{ii}^{(t)} + n_{it}^{(t)} + n_{ci}^{(t)}}{p_i} - \frac{2n_{tt}^{(t)} + n_{ct}^{(t)} + n_{it}^{(t)}}{1 - p_c - p_i}. \quad (58)$$

Setting these derivatives equal to zero and solving for  $p_c$  and  $p_i$  completes the M step, yielding

$$p_c^{(t+1)} = \frac{2n_{cc}^{(t)} + n_{ci}^{(t)} + n_{ct}^{(t)}}{2n} \quad (59)$$

$$p_i^{(t+1)} = \frac{2n_{ii}^{(t)} + n_{it}^{(t)} + n_{ci}^{(t)}}{2n} \quad (60)$$

$$p_t^{(t+1)} = \frac{2n_{tt}^{(t)} + n_{ct}^{(t)} + n_{it}^{(t)}}{2n}, \quad (61)$$

**Observed data:**  $n_C = 85$ ,  $n_I = 196$ , and  $n_T = 341$

Table 5: EM results for peppered moth example.  $R^{(t)}$  is the relative convergence criterion;  $D_C^{(t)}$ , and  $D_I^{(t)}$  are ratios of consecutive errors.

$t$	$p_C^{(t)}$	$p_I^{(t)}$	$R^{(t)}$	$D_C^{(t)}$	$D_I^{(t)}$
0	0.333333	0.333333			
1	0.081994	0.237406	$5.7 \times 10^{-1}$	0.0425	0.337
2	0.071249	0.197870	$1.6 \times 10^{-1}$	0.0369	0.188
3	0.070852	0.190360	$3.6 \times 10^{-2}$	0.0367	0.178
4	0.070837	0.189023	$6.6 \times 10^{-3}$	0.0367	0.176
5	0.070837	0.188787	$1.2 \times 10^{-3}$	0.0367	0.176
6	0.070837	0.188745	$2.1 \times 10^{-4}$	0.0367	0.176
7	0.070837	0.188738	$3.6 \times 10^{-5}$	0.0367	0.176
8	0.070837	0.188737	$6.4 \times 10^{-6}$	0.0367	0.176

## The nature of EM

**Ascent:** Each M-step increases the log likelihood.

**Convergence:** is linear (slow!). Rate is inversely related to the proportion of missing data.

**Optimization transfer:**

$$l(\boldsymbol{\theta}|\mathbf{x}) \geq Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) + l(\boldsymbol{\theta}^{(t)}|\mathbf{x}) - Q(\boldsymbol{\theta}^{(t)}|\boldsymbol{\theta}^{(t)}) = G(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}). \quad (62)$$

The last two terms in  $G(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$  are constant with respect to  $\boldsymbol{\theta}$ , so  $Q$  and  $G$  are maximized at the same  $\boldsymbol{\theta}$ . Further,  $G$  is tangent to  $l$  at  $\boldsymbol{\theta}^{(t)}$ , and lies everywhere below  $l$ . We say that  $G$  is a *minorizing function* for  $l$ . EM transfers optimization from  $l$  to the surrogate function  $G$ , which is more convenient to maximize.

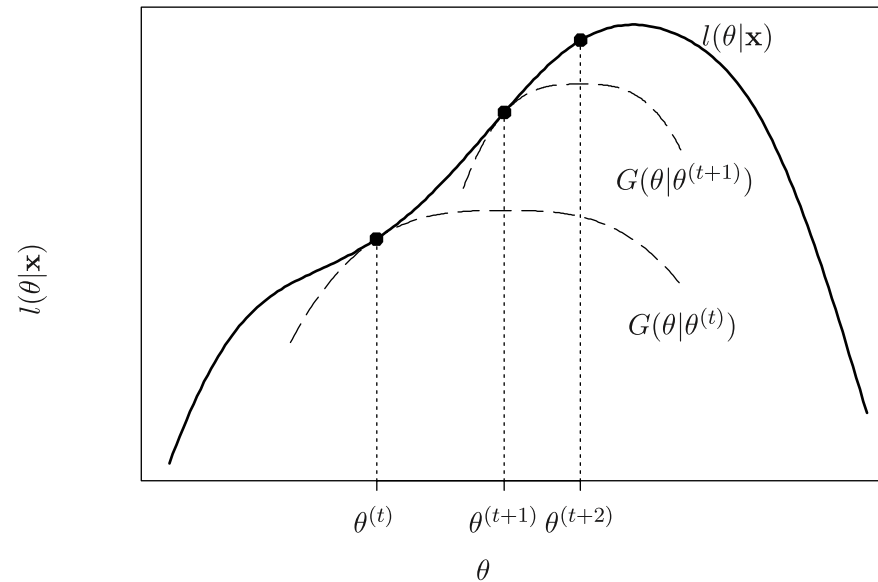


Figure 21: One-dimensional illustration of EM algorithm as a minorization or optimization transfer strategy.

Each E step forms a minorizing function  $G$ , and each M step maximizes it to provide an uphill step.

## EM in exponential families

1. **E step:** Compute the expected values of the sufficient statistics for the complete data, given the observed data and using the current parameter guesses,  $\boldsymbol{\theta}^{(t)}$ . Let  $\mathbf{s}^{(t)} = E\{\mathbf{s}(\mathbf{Y})|\mathbf{x}, \boldsymbol{\theta}^{(t)}\} = \int \mathbf{s}(\mathbf{y})f_{\mathbf{Z}|\mathbf{X}}(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^{(t)}) d\mathbf{z}$ .
2. **M step:** Set  $\boldsymbol{\theta}^{(t+1)}$  to the value which makes the unconditional expectation of the sufficient statistics for the complete data equal to  $\mathbf{s}^{(t)}$ . In other words,  $\boldsymbol{\theta}^{(t+1)}$  solves  $E\{\mathbf{s}(\mathbf{Y})|\boldsymbol{\theta}\} = \mathbf{s}^{(t)}$ .
3. Return to the E step unless a convergence criterion has been met.

## The missing information principle

Can show

$$\log f_{\mathbf{X}}(\mathbf{x}|\boldsymbol{\theta}) = Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) - H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \quad (63)$$

where

$$H(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = \text{E} \left\{ \log f_{\mathbf{Z}|\mathbf{X}}(\mathbf{Z}|\mathbf{x}, \boldsymbol{\theta}) \mid \mathbf{x}, \boldsymbol{\theta}^{(t)} \right\}. \quad (64)$$

Taking second order partial derivatives of (63) and negating both sides yields

$$-l''(\boldsymbol{\theta}|\mathbf{x}) = -\mathbf{Q}''(\boldsymbol{\theta}|\boldsymbol{\omega})|_{\boldsymbol{\omega}=\boldsymbol{\theta}} + \mathbf{H}''(\boldsymbol{\theta}|\boldsymbol{\omega})|_{\boldsymbol{\omega}=\boldsymbol{\theta}} \quad (65)$$

where the primes on  $\mathbf{Q}''$  and  $\mathbf{H}''$  denote derivatives with respect to the first argument,  $\boldsymbol{\theta}$ .

Rewrite (65) as

$$\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta}) = \hat{\mathbf{i}}_{\mathbf{Y}}(\boldsymbol{\theta}) - \hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{X}}(\boldsymbol{\theta}) \quad (66)$$

where  $\hat{\mathbf{i}}_{\mathbf{X}}(\boldsymbol{\theta}) = -l''(\boldsymbol{\theta}|\mathbf{x})$  is the observed information, and  $\hat{\mathbf{i}}_{\mathbf{Y}}(\boldsymbol{\theta})$  and  $\hat{\mathbf{i}}_{\mathbf{Z}|\mathbf{X}}(\boldsymbol{\theta})$  will be called the complete information and the missing information, respectively.

The observed information equals the complete information minus the missing information.

## Obtaining variance estimates

---

1. Bootstrap
2. SEM algorithm (next)

Less common methods include Louis's method (using missing information principle), empirical information, and numerical differentiation.



## Supplemented EM (SEM) algorithm

---

Let  $\Psi$  denotes the EM mapping, having fixed point  $\hat{\theta}$  and Jacobian matrix  $\Psi'(\theta)$  with  $(i, j)$ th element equaling  $\frac{d\Psi_i(\theta)}{d\theta_j}$ . It can be shown that

$$\Psi'(\hat{\theta})^T = \hat{\mathbf{i}}_{Z|X}(\hat{\theta})\hat{\mathbf{i}}_Y(\hat{\theta})^{-1} \quad (67)$$

Further use of the missing information principle leads to

$$\widehat{\text{var}}\{\hat{\theta}\} = \hat{\mathbf{i}}_Y(\hat{\theta})^{-1} \left( \mathbf{I} + \Psi'(\hat{\theta})^T (\mathbf{I} - \Psi'(\hat{\theta})^T)^{-1} \right). \quad (68)$$

SEM is numerically stable and requires little extra work.

## SEM algorithm:

1. Run the EM algorithm to convergence, finding  $\hat{\theta}$ .
2. Restart the algorithm from some  $\theta^{(0)}$  near to  $\hat{\theta}$ . For  $t = 0, 1, 2, \dots$ 
  - (a) Take a standard E step and M step to produce  $\theta^{(t+1)}$  from  $\theta^{(t)}$ .
  - (b) For  $j = 1, \dots, p$ , define  $\theta^{(t)}(j) = (\hat{\theta}_1, \dots, \hat{\theta}_{j-1}, \theta_j^{(t)}, \hat{\theta}_{j+1}, \dots, \hat{\theta}_p)$  and

$$r_{ij}^{(t)} = \frac{\Psi_i(\theta^{(t)}(j)) - \hat{\theta}_i}{\theta_j^{(t)} - \hat{\theta}_j} \quad (69)$$

for  $i = 1, \dots, p$ . (Recall that  $\Psi(\hat{\theta}) = \hat{\theta}$ .)

(c) Stop when all  $r_{ij}^{(t)}$  have converged

3. The  $(i, j)$ th element of  $\Psi'(\hat{\theta})$  equals  $\lim_{t \rightarrow \infty} r_{ij}^{(t)}$ . Plug the final estimate of  $\Psi'(\hat{\theta})$  into (68) to get the variance.

## Facilitating the E-step

---

### Monte Carlo EM (MCEM)

Replace the  $t$ th E step with

1. Draw missing datasets  $\mathbf{Z}_1^{(t)}, \dots, \mathbf{Z}_{m^{(t)}}^{(t)}$  i.i.d. from  $f_{\mathbf{Z}|\mathbf{X}}(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}^{(t)})$ . Each  $\mathbf{Z}_j^{(t)}$  is a vector of all the missing values needed to complete the observed dataset, so  $\mathbf{Y}_j = (\mathbf{x}, \mathbf{Z}_j)$  denotes a completed dataset where the missing values have been replaced by  $\mathbf{Z}_j$ .
2. Calculate  $\hat{Q}^{(t+1)}(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) = \frac{1}{m^{(t)}} \sum_{j=1}^{m^{(t)}} \log f_{\mathbf{Y}}(\mathbf{Y}_j^{(t)}|\boldsymbol{\theta})$ .

Then  $\hat{Q}^{(t+1)}(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$  is a Monte Carlo estimate of  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$ .

The M step is modified to maximize  $\hat{Q}^{(t+1)}(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$ .

Increase  $m^{(t)}$  as iterations progress to reduce the Monte Carlo variability of  $\hat{Q}$ . MCEM will not converge in the same sense as ordinary EM, rather values of  $\boldsymbol{\theta}^{(t)}$  will bounce around the true maximum, with a precision that depends on  $m^{(t)}$ .

## Facilitating the M-step

---

### ECM algorithm

Replaces the M step with a series of computationally simpler conditional maximization (CM) steps.

Call the collection of simpler CM steps after the  $t$ th E step a CM *cycle*. Thus, the  $t$ th iteration of ECM is comprised of the  $t$ th E step and the  $t$ th CM cycle.

Let  $S$  denote the total number of CM steps in each CM cycle.

For  $s = 1, \dots, S$ , the  $s$ th CM step in the  $t$ th cycle requires the maximization of  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})$  subject to (or conditional on) a constraint, say

$$\mathbf{g}_s(\boldsymbol{\theta}) = \mathbf{g}_s(\boldsymbol{\theta}^{(t+(s-1)/S)}) \quad (70)$$

where  $\boldsymbol{\theta}^{(t+(s-1)/S)}$  is the maximizer found in the  $(s - 1)$ th CM step of the current cycle.

When the entire cycle of  $S$  steps of CM has been completed, we set  $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t+S/S)}$  and proceed to the E step for the  $(t + 1)$ th iteration.

The art of constructing an effective ECM algorithm lies in choosing the constraints cleverly.

## Choice 1: Iterated Conditional Modes / Gauss-Seidel

Partition  $\boldsymbol{\theta}$  into  $S$  subvectors,  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_S)$ .

In the  $s$ th CM step, maximize  $Q$  with respect to  $\boldsymbol{\theta}_s$  while holding all other components of  $\boldsymbol{\theta}$  fixed.

This amounts to the constraint induced by the function

$$g_s(\boldsymbol{\theta}) = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{s-1}, \boldsymbol{\theta}_{s+1}, \dots, \boldsymbol{\theta}_S).$$

## Choice 2:

At the  $s$ th CM step, maximize  $Q$  with respect to all other components of  $\theta$  while holding  $\theta_s$  fixed.

Then  $g_s(\theta) = \theta_s$ .

Additional systems of constraints can be imagined, depending on the particular problem context.

A variant of ECM inserts an E step between each pair of CM steps, thereby updating  $Q$  at every stage of the CM cycle.

(Sketched) Example:

Multivariate regression (Meng and Rubin, 1993)

Let  $U_1, \dots, U_n$  be independent where

$$U_i \sim N_d(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}) \quad (71)$$

for  $\boldsymbol{\mu}_i = V_i \boldsymbol{\beta}$ , where the  $V_i$  are known  $d \times p$  design matrices, and  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$  are unknown parameters.

To estimate  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$  when some elements of some of the  $U_i$  are missing, we can partition each CM cycle into two stages: one for  $\boldsymbol{\beta}$  and one for  $\boldsymbol{\Sigma}$ .



The algebra is messy, but the idea turns out to be simple.

- **E-step:** Find the expected values of the complete data sufficient statistics conditional on the observed data and current guesses  $\beta^{(t)}, \Sigma^{(t)}$ . (The sufficient stats are  $\sum_{i=1}^n U_{ij}$  for  $j = 1, \dots, d$  and  $\sum_{i=1}^n U_{ij}U_{ik}$  for  $j, k = 1, \dots, d$ .)
- **CM 1:**  $\beta^{(t+1/2)}$  is the weighted least squares estimate given  $\Sigma = \Sigma^{(t)}$ , computed using the sufficient stats from the E-step.
- **CM 2:**  $\Sigma^{(t+2/2)}$  is the sample covariance matrix of the completed data, using  $\beta = \beta^{(t+1/2)}$ .
- **Return to the E-step**

## EM gradient algorithm

Replace the M step with a single step of Newton's method, thereby approximating the maximum without actually solving for it exactly.

Instead of maximizing, choose:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \mathbf{Q}''(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})^{-1} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \mathbf{Q}'(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \quad (72)$$

$$= \boldsymbol{\theta}^{(t)} - \mathbf{Q}''(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})^{-1} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} \mathbf{l}'(\boldsymbol{\theta}^{(t)}|\mathbf{x}) \quad (73)$$

Ascent is ensured for canonical parameters in exponential families. Backtracking can ensure ascent in other cases; inflating steps can speed convergence.

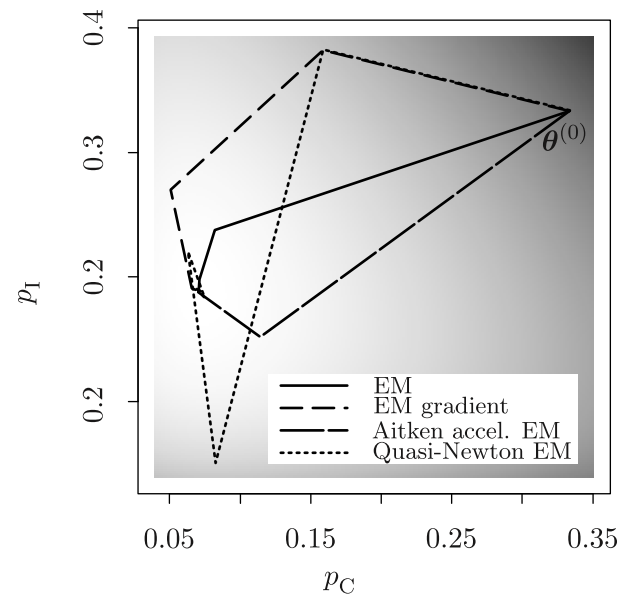


Figure 22: Steps taken by the EM gradient algorithm (long dashes). Ordinary EM steps are shown with the solid line. Steps from two methods from later sections (Aitken and quasi-Newton acceleration) are also shown, as indicated in the key. The observed data log likelihood is shown with the grey scale, with light shading corresponding to high likelihood. All algorithms were started from  $p_C = p_I = 1/3$ .

## EM Acceleration Methods

---

Standard EM is reliable but slow. Acceleration methods use the EM setup to motivate particular forms of Newton-like steps.

**Aitken acceleration:** Equivalent to applying the Newton-Raphson method to find a zero of  $\Psi(\boldsymbol{\theta}) - \boldsymbol{\theta}$  where  $\Psi$  is the mapping defined by the ordinary EM algorithm producing  $\boldsymbol{\theta}^{(t+1)} = \Psi(\boldsymbol{\theta}^{(t)})$ .

**Quasi-Newton acceleration:** Take

$$\mathbf{M}^{(t)} = \mathbf{Q}''(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t)})\Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(t)}} - \mathbf{B}^{(t)}$$

where  $\mathbf{B}^{(t)}$  is crafted to make  $\mathbf{M}^{(t)}$  approximate  $\mathbf{l}''(\boldsymbol{\theta}^{(t)}|\mathbf{x})$ . A potentially superior strategy approximates  $(\mathbf{l}'')^{-1}$  instead of  $\mathbf{l}''$ .

## A final example: the Baum-Welch algorithm

---

A *hidden Markov model* can be used to describe the joint probability of a sequence of unobserved (‘hidden’) discrete state variables,  $\mathbf{H} = (H_0, \dots, H_n)$ , and a sequence of corresponding observed variables  $\mathbf{O} = (O_0, \dots, O_n)$  for which  $O_i$  is dependent on  $H_i$  for each  $i$ .

We say that  $H_i$  emits  $O_i$ .

Let the (discrete) state spaces for elements of  $\mathbf{H}$  and  $\mathbf{O}$  be  $\mathcal{H}$  and  $\mathcal{E}$ , respectively.

Such models are extremely useful in statistical genetics, signal processing, and environmental time series. For example, the hidden states can be “weather patterns” and the observed variables can be measurable weather events like rain, snow, smog, etc.

**The Baum-Welch algorithm fits such models. It is an EM algorithm.**

The parameters of a HMM are  $\theta = (\pi, \mathbf{P}, \mathbf{E})$ , where  $\pi$  is a vector of initial state probabilities,  $\mathbf{P}$  is a matrix of transition probabilities,  $\mathbf{E}$  is a matrix of emission probabilities.

For an observed sequence  $\mathbf{o}$ , define the *forward variables* to be

$$\alpha(i, h) = P[\mathbf{O}_{\leq i} = \mathbf{o}_{\leq i}, H_i = h] \quad (74)$$

and the *backward variables* to be

$$\beta(i, h) = P[\mathbf{O}_{> i} = \mathbf{o}_{> i} | H_i = h] \quad (75)$$

for  $i = 1, \dots, n$  and each  $h \in \mathcal{H}$ . Dependence on  $\theta$  is suppressed.

There are efficient recursive formulas to calculate the  $\alpha(i, h)$  and  $\beta(i, h)$ .

**The E-step boils down to:**

Let  $N(h)$  denote the number of times  $H_0 = h$ , let  $N(h, h^*)$  denote the number of transitions from  $h$  to  $h^*$ , and let  $N(h, o)$  denote the number of emissions of  $o$  when the underlying state is  $h$ . Then:

$$E\{N(h)\} = \frac{\alpha(0, h)\beta(0, h)}{P[\mathbf{O} = \mathbf{o}|\boldsymbol{\theta}]} \quad (76)$$

$$E\{N(h, h^*)\} = \sum_{i=0}^{n-1} \frac{\alpha(i, h)p(h, h^*)e(h^*, o_{i+1})\beta(i+1, h^*)}{P[\mathbf{O} = \mathbf{o}|\boldsymbol{\theta}]} \quad (77)$$

$$E\{N(h, o)\} = \sum_{i: O_i=o} \frac{\alpha(i, h)\beta(i, h)}{P[\mathbf{O} = \mathbf{o}|\boldsymbol{\theta}]} \quad (78)$$

The M-step boils down to:

$$\pi(h)^{(t+1)} = \frac{\mathbb{E}\{N(h)|\boldsymbol{\theta}^{(t)}\}}{\sum_{h^* \in \mathcal{H}} \mathbb{E}\{N(h^*)|\boldsymbol{\theta}^{(t)}\}} \quad (79)$$

$$p(h, h^*)^{(t+1)} = \frac{\mathbb{E}\{N(h, h^*)|\boldsymbol{\theta}^{(t)}\}}{\sum_{h^{**} \in \mathcal{H}} \mathbb{E}\{N(h, h^{**})|\boldsymbol{\theta}^{(t)}\}} \quad (80)$$

$$e(h, o)^{(t+1)} = \frac{\mathbb{E}\{N(h, o)|\boldsymbol{\theta}^{(t)}\}}{\sum_{o^* \in \mathcal{E}} \mathbb{E}\{N(h, o^*)|\boldsymbol{\theta}^{(t)}\}}. \quad (81)$$



(Oversimplified) HMM data:

A region is hypothesized to cycle between two weather patterns: “tending-wet” and “tending-dry”. The probability of measurable precipitation is  $w$  for a day during a tending-wet cycle and  $d$  for a day in the tending-dry cycle. Cycles switch from dry to wet or vice versa with probability  $s$ .

**Observed Data:** was there measurable precipitation on each of 200 consecutive days?

**Latent Data:** which pattern was extant on each day.

**Goal:** Estimate  $w$ ,  $d$ , and  $s$ .

Dry pattern (unobservable) ■  
Wet pattern (unobservable) ■  
Rainy day |

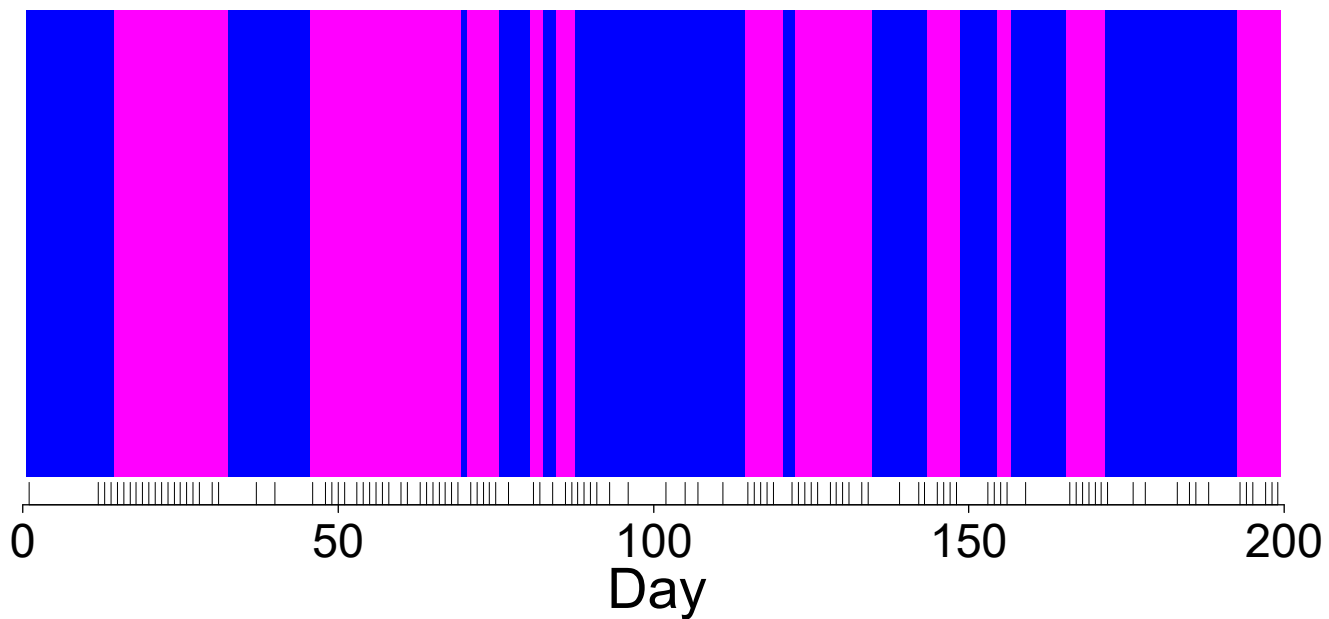


Figure 23: Simple data for a meteorological hidden Markov model.

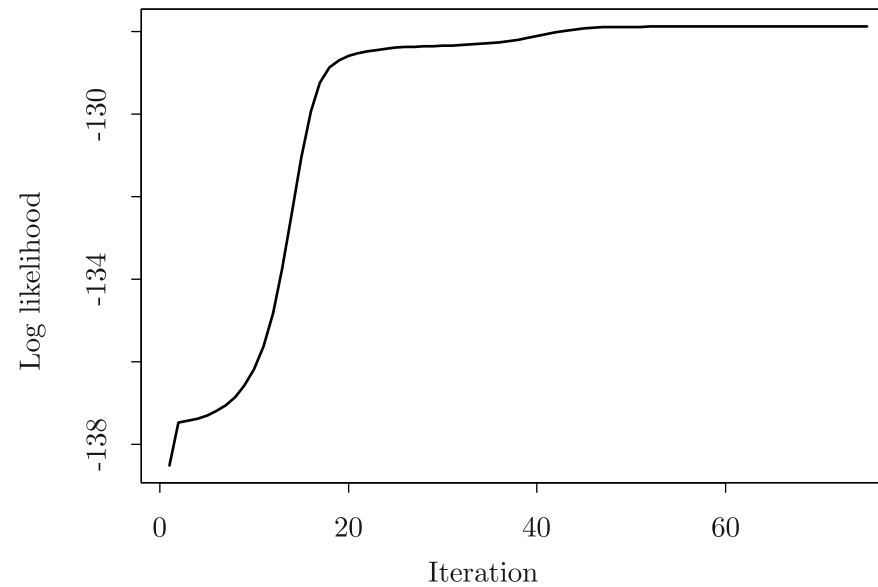


Figure 24: Log likelihood for parameters of hidden Markov model, against Baum-Welch iteration.

**As an EM algorithm, the Baum-Welch approach reliably ascends as iterations progress.**

Iteration, $t$	$d^{(t)}$	$w^{(t)}$	$s^{(t)}$
0	0.4500	0.5500	0.4000
1	0.4941	0.6042	0.3983
5	0.4632	0.6365	0.3872
10	0.3853	0.7127	0.3467
20	0.2217	0.8658	0.1750
30	0.2349	0.8528	0.1444
40	0.2475	0.8497	0.1314
50	0.2595	0.8468	0.1190
75	0.2608	0.8460	0.1174
100	0.2608	0.8460	0.1174

Table 6: Selected iterations during Baum-Welch search for MLEs for hidden Markov model for rain data.

The data were simulated with  $d = 0.25$ ,  $w = 0.85$ , and  $s = 0.10$ .