

(1)

Homework 3 - Optimization & EM will be shortly emailed.

1. For this Q, code and calculations are attached.

2. For $i = 1:n$, consider the hierarchical model $\begin{cases} Y_i | \tau_i \sim N_p(\mu, \frac{1}{\tau_i} \Psi) \\ \tau_i \sim \frac{1}{\nu} \chi_\nu^2 \end{cases}$ iid.

Then we have an AA model for EM with

$Y_0 := Y = (Y_1, \dots, Y_n)$, $Y_m := \tau = (\tau_1, \dots, \tau_n)$, $\theta := (\mu, \Psi)$,
assuming ν is provided. Note $\mu \in \mathbb{R}^p$, $\Psi \in \mathbb{R}^{p \times p}$.

We are asked to give the updates $\theta^{(t)} \mapsto \theta^{(t+1)}$.

First, we give $\log p(Y_0, Y_m | \theta) = \log p(Y | \tau, \mu, \Psi) + \log p(\tau)$,
suppressing those terms independent of μ, Ψ :

$$p(Y | \tau, \mu, \Psi) = \prod_{i=1}^n \left(\frac{\tau_i}{2\pi} \right)^{p/2} |\Psi|^{-1/2} \exp \left\{ -\frac{\tau_i}{2} (y_i - \mu)^T \Psi^{-1} (y_i - \mu) \right\}$$

$$p(\tau) = \prod_{i=1}^n \nu^{-1} 2^{-\frac{\nu}{2}} (\Gamma(\frac{\nu}{2}))^{-1} \tau_i^{\frac{\nu}{2}-1} \exp \left\{ -\frac{\tau_i}{2} \right\}$$

Note in particular that $p(\tau)$ will not affect the maximization of $Q(\theta | \theta^t)$.

Taking the log, the relevant terms are:

$$\sum_{i=1}^n \left\{ -\frac{1}{2} \log |\Psi| - \frac{\tau_i}{2} (y_i - \mu)^T \Psi^{-1} (y_i - \mu) \right\}$$

$$= -\frac{1}{2} \left\{ n \log |\Psi| + \sum_{i=1}^n \tau_i (y_i - \mu)^T \Psi^{-1} (y_i - \mu) \right\}.$$

Now, we examine the second term. A quadratic form can be



$$* \text{tr}(cA) = \sum_i (cA)_{ii} = \sum_i cA_{ii} = c \text{tr} A$$

(2)

seen as the trace of a certain matrix product:

$$\text{tr}(\underbrace{xx^T A}_{\text{data matrix}}) = \sum_i (xx^T)_{ii} A_{ii} = \sum_i x_i A_{ii} x_i = \sum_i (x^T A)_i x_i = x^T A x$$

Hence, we may rewrite this term as

$$\begin{aligned} \sum_{i=1}^n \gamma_i \text{tr}((y_i - u)(y_i - u)^T \Psi^{-1}) &= \sum_{i=1}^n \text{tr}(\gamma_i (y_i - u)(y_i - u)^T \Psi^{-1}) \\ &= \text{tr} \left\{ \underbrace{\left(\sum_{i=1}^n \gamma_i (y_i - u)(y_i - u)^T \right)}_{\text{call this matrix } B} \Psi^{-1} \right\} \end{aligned}$$

by additivity
of trace

Then our log-likelihood looks like

$$-\frac{1}{2} \left\{ n \log |\Psi| + \text{tr}(B \Psi^{-1}) \right\}. \text{ Now we consider}$$

$Q(\theta | \theta') = \mathbb{E}_{\tau | Y, \theta'} [\log p(Y, \tau | \theta) | Y, \theta']$, where θ' is a particular value of u, Ψ . Considering only the same terms, τ appears only linearly, in the expression for B :

$$\mathbb{E}[B | Y, \theta'] = \sum_{i=1}^n \mathbb{E}[\tau_i | Y, \theta'] (y_i - u)(y_i - u)^T.$$

For convenience we denote $\tau'_i := \mathbb{E}[\tau_i | Y, \theta']$, and we are given the fact that

$$\tau'_i = \frac{\nu + \rho}{\nu + (y_i - u')^T (\Psi')^{-1} (y_i - u')}$$

Denote the corresponding matrix B' . Then (up to relevant terms), $Q(\theta | \theta') = -\frac{1}{2} \{ n \log |\Psi| + \text{tr}(B' \Psi^{-1}) \}$.



(3)

HW 3 cont.

2. cont.

Now, given $\theta^{(t)}$, the E-step is computing $\tilde{\gamma}'(\theta^{(t)}) = \langle \tilde{\gamma}_1' \dots \tilde{\gamma}_n' \rangle$, using $\theta^{(t)}$ as the value of θ' .

The M-step is finding the argmax wrt each parameter.

To find $\mu^{(t+1)}$, we need to find the maximum of the (negative-semidefinite) quadratic form $-\frac{1}{2} \text{tr}(\beta' \Psi^{-1}) = -\frac{1}{2} \sum_{i=1}^n \tilde{\gamma}_i' (\mathbf{y}_i - \mu)' \Psi^{-1} (\mathbf{y}_i - \mu)$.

Fortunately such a form always has a unique maximizer, found by observing the stationary point. In general, applying $\frac{\partial}{\partial x_k} := \partial_k$ to $x^T A x$ yields

$$\begin{aligned} \partial_k(x^T A x) &= \partial_k \sum_{ij} A_{ij} x_i x_j = \sum_{ij} A_{ij} [x_i \delta_{jk} + x_j \delta_{ik}] \\ &= \sum_j [A_{jk} + A_{kj}] x_j = (A + A^T)_{k \circ} x \end{aligned}$$

where $k \circ$ subscript denotes the k^{th} row. Applying this in our case, concatenating the p row-equations, and setting to 0 yields

$$\sum_{i=1}^n \tilde{\gamma}_i' (\Psi^{-1})' (\mathbf{y}_i - \mu) = 0. \quad \text{As a covariance matrix is symmetric, pos.-definite, this simplifies to}$$

$$\sum_{i=1}^n \tilde{\gamma}_i' \Psi^{-1} \mathbf{y}_i = \sum_{i=1}^n \tilde{\gamma}_i' \Psi^{-1} \mu, \quad \text{and so}$$

$$\mu^{(t+1)} = \underbrace{\frac{\sum_{i=1}^n \tilde{\gamma}_i' \mathbf{y}_i}{\sum_{i=1}^n \tilde{\gamma}_i'}}$$

$$\text{For } \Psi, \text{ we are supplied that } \underset{A}{\operatorname{argmax}} \left\{ n \log |A^{-1}| - \text{tr}(A^{-1} B) \right\} = \frac{1}{n} B,$$

Clearly our $Q(\theta|\theta')$ can easily be rearranged into this form.

(3)

(4).

$$|\mathbf{4}^{-1}| = \frac{1}{|\mathbf{4}|} \text{ since } |\mathbf{4}||\mathbf{4}^{-1}| = |\mathbf{4}\mathbf{4}^{-1}| = 1$$

Doing so yields $\frac{1}{2} \{ n \log |\mathbf{4}^{-1}| - \text{tr}(\mathbf{4}^{-1} \mathbf{B}') \},$

so $\boxed{\mathbf{4}^{(t+1)} = \frac{1}{2n} \mathbf{B}'}$, (where \mathbf{B}' has γ' evaluated at $\theta^{(t)}$). \square

3. For $i=1:n$, consider the hierarchical model $\begin{cases} Y_i | \lambda_i \sim \text{Pois}(\lambda_i) \\ \lambda_i | \beta \sim \text{Gamma}(1, \beta) \end{cases}$

This is an SA model $Y_m = \lambda, Y_0 = Y, \theta = \beta$.

for reference, $Y \sim \text{Pois}(\lambda), p(y) = \frac{1}{y!} e^{-\lambda} \lambda^y$
 $X \sim \text{Gamma}(\alpha, \beta), p(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp\{-\beta x\}$

- (a) Derive the EM model for β in the SA framework.

To compute $Q(\beta | \beta')$, we need the β -having terms of $\log p(Y, \lambda | \beta)$.

These are: $\sum_i \{ \log \beta - \lambda_i \beta \} = n \log \beta - \beta \sum_i \lambda_i$

Taking the conditional expectation, we have

$$n \log \beta - \beta \sum_i \mathbb{E}[\lambda_i | Y, \beta'], \text{ so we need Bayes rule.}$$

$$\begin{aligned} p(\lambda_i | Y, \beta) &\propto p(Y_i | \lambda_i) p(\lambda_i | \beta) \text{ as a function of } \lambda_i; \text{ that is} \\ p(\lambda_i | Y, \beta) &\propto e^{-\lambda_i} \lambda_i^{Y_i} e^{-\beta \lambda_i} = \lambda_i^{Y_i} e^{-(\beta+1)\lambda_i} \sim \text{Gamma}(Y_i + 1, \beta + 1) \end{aligned}$$

If $X \sim \text{Gamma}(\alpha, \beta)$, $\mathbb{E}[x] = \frac{\alpha}{\beta}$. Hence,

$$Q(\beta | \beta') = n \log \beta - \beta \sum_i \frac{Y_i + 1}{\beta' + 1} \text{ for fixed } \beta'.$$

Taking $\beta' = \beta^{(t)}$, our iteration is found by the arg max:



(5)

$$\text{let } \bar{y} := \frac{1}{n} \sum y_i$$

HW 3 cont.

3@cont.

$$\text{do } Q(\beta | \beta^{(t)}) = \frac{n}{\beta} - \sum_i \frac{y_i + 1}{\beta^{(t)} + 1} \Rightarrow 0; \text{ thus, } \beta^{(t+1)} = \frac{\beta^{(t)} + 1}{\bar{y} + 1}.$$

(Here the E-step was performed analytically.)

(b)

Construct an AA model which preserves $p(Y|\beta)$ in the above model.

Here we need to trivialize the distribution of $\lambda|\beta$; because it is part of an exponential family, consider $\tilde{\lambda}_i = \beta \lambda_i$ as a transformation.

$$\text{Then } p(\tilde{\lambda}_i \leq x) = p(\beta \lambda_i \leq x) = p(\lambda_i \leq \frac{x}{\beta}) = \int_0^{\frac{x}{\beta}} \beta e^{-\beta \lambda} d\lambda \\ = 1 - e^{-x}.$$

Hence $p(\tilde{\lambda}_i) = e^{-\tilde{\lambda}_i}$, or $\tilde{\lambda}_i \sim \text{Gamma}(1, 1)$.

Consequently, our model becomes $\begin{cases} Y_i | \tilde{\lambda}_i, \beta \sim \text{Pois}\left(\frac{1}{\beta} \tilde{\lambda}_i\right) \\ \tilde{\lambda}_i \sim \text{Gamma}(1, 1) \end{cases}$

Now, our log-likelihood (with only β -relevant terms) is

$$\sum_i \left\{ -\frac{1}{\beta} \tilde{\lambda}_i + y_i \log \frac{1}{\beta} \right\} = -\frac{1}{\beta} \sum_i \tilde{\lambda}_i - n \bar{y} \log \beta$$

Taking the conditional expectation, we have

$$-\frac{1}{\beta} \sum_i \mathbb{E}[\tilde{\lambda}_i | Y, \beta] - n \bar{y} \log \beta, \text{ so we need Bayes rule.}$$

$$p(\tilde{\lambda}_i | Y, \beta) \propto p(Y_i | \tilde{\lambda}_i, \beta) p(\tilde{\lambda}_i) \text{ as a function of } \tilde{\lambda}_i; \text{ that is,} \\ p(\tilde{\lambda}_i | Y, \beta) \propto e^{-\frac{1}{\beta} \tilde{\lambda}_i} \left(\frac{1}{\beta} \tilde{\lambda}_i\right)^{y_i} e^{-\tilde{\lambda}_i} \sim \text{Gamma}(y_i + 1, \frac{1}{\beta} + 1).$$

$$\text{So, } \mathbb{E}[\tilde{\lambda}_i | Y, \beta] = \frac{y_i + 1}{\frac{1}{\beta} + 1} = \underbrace{\frac{\beta}{\beta + 1} (y_i + 1)}_{}, \text{ Putting this together,}$$

$$Q(\beta | \beta') = -\frac{1}{\beta} \left(\frac{\beta'}{\beta + 1}\right) n(\bar{y} + 1) - n \bar{y} \log \beta.$$

So to find the argmax,



(6)

$$\frac{\partial}{\partial \beta} Q(\beta | \beta^{(t)}) = \frac{1}{\beta^2} \left(\frac{\beta^{(t)}}{\beta^{(t)} + 1} \right) n(\bar{y} + 1) - \frac{1}{\beta} n\bar{y} \Rightarrow 0.$$

Thus $\boxed{\beta^{(t+1)} = \frac{\beta^{(t)}}{\beta^{(t)} + 1} \left(1 + \frac{1}{\bar{y}} \right)}$

- (c) Using the SA and AA from (a), derive the corresponding Intertwoven EM algorithm,

Let's take SA as A1, AA as A2. The IEM Q-function is

$$Q_I(\theta | \theta') = E_{A2} \left[E_{A1} \left[\log P_{A1}(Y_0, Y_m | \theta) \mid Y_0, \tilde{Y}_m, G_{A2}(\theta') \right] \mid Y_0, \theta' \right],$$

where E_{A2} integrates $\tilde{Y}_m \mid Y_0, \theta'$, E_{A1} integrates $Y_m \mid Y_0, \tilde{Y}_m, G_{A2}(\theta')$, and G_{A2} maps θ to its EM update under A2.

$\log P_{A1}$ up to β -terms is $n \log \beta - \beta \sum_i \lambda_i$, so taking E_{A1} yields $n \log \beta - \beta \sum_i E_{A1}[\lambda_i \mid Y, \tilde{\lambda}, G_{A2}(\beta')]$

The relationship between $\lambda_i, \tilde{\lambda}$ is deterministic, so we have the same A1 expectation as before: $\frac{Y_i + 1}{\beta + 1}$. Conditioning on β , we get $E_{A1}[\lambda_i \mid Y, \tilde{\lambda}, G_{A2}(\beta')] = (Y_i + 1) / \left[\left(\frac{\beta'}{\beta' + 1} \right) \left(1 + \frac{1}{\tilde{\lambda}} \right) + 1 \right]$

Because $\tilde{\lambda}$ doesn't appear, E_{A2} has no effect, so:

$$Q_I(\beta | \beta') = n \left[\log \beta - \beta \frac{1}{\left(\frac{\beta'}{\beta' + 1} \right) \left(1 + \frac{1}{\tilde{\lambda}} \right) + 1} (\bar{y} + 1) \right]$$

$$= n \left[\log \beta - \beta \frac{1}{\left(\frac{\beta'}{\beta' + 1} \right) \frac{1}{\bar{y}} + \frac{1}{1 + \bar{y}}} \right].$$

up to β -terms.



3③ cont.

HW 3 cont.

To find the argmax, $\partial_B Q(\beta|\beta^{(+)}) = n \left[\frac{1}{\beta} - \frac{1}{(\frac{\beta^{(+)}}{\beta^{(+)}}+1)} \frac{1}{\bar{y}} + \frac{1}{1+\bar{y}} \right] \Rightarrow 0$,

so $\beta^{(t+1)} = \frac{\beta^{(+)}}{\beta^{(+)}} + \frac{1}{\bar{y}} + \frac{1}{1+\bar{y}}$.

(d) Derive the observed-data log-likelihood and compute the MLE.
(we use the SA model)

We know $p(y|\beta) = \int p(y|\lambda, \beta) p(\lambda|\beta) d\lambda$
= $\prod_i \int p(y_i|\lambda_i) p(\lambda_i|\beta) d\lambda_i$,

so that $\log p(y|\beta) = \sum_i \log \int p(y_i|\lambda_i) p(\lambda_i|\beta) d\lambda_i$.

$$\begin{aligned} \int p(y_i|\lambda_i) p(\lambda_i|\beta) d\lambda_i &= \frac{\beta}{y_i!} \underbrace{\int e^{-\lambda_i} \lambda_i^{y_i} e^{-\beta \lambda_i} d\lambda_i}_{\sim \text{Gamma}(y_i+1, \beta+1)} \\ &= \frac{\beta}{y_i!} \frac{y_i!}{(\beta+1)^{y_i+1}} = \frac{\beta}{\beta+1} (\beta+1)^{-y_i} \end{aligned}$$

Hence $\log p(y|\beta) = \sum_i \log \beta - (y_i+1) \log(\beta+1)$
= $n [\log \beta - (\bar{y}+1) \log(\beta+1)]$

To find the MLE, $\partial_B \log p(y|\beta) = n \left[\frac{1}{\beta} - (\bar{y}+1) \frac{1}{\beta+1} \right] \Rightarrow 0$,

then $\hat{\beta} = \frac{1}{\bar{y}}$

(e) Compare the linear rate of convergence of each EM algorithm,
and discuss when each will perform well.

(8)

The SA update is $\beta \mapsto \frac{\beta+1}{\bar{\gamma}+1}$. Its fixed points satisfy $\beta^* = \frac{\beta^*+1}{\bar{\gamma}+1}$, so that $\beta^* = \frac{1}{\bar{\gamma}}$ as desired. The linear convergence rate is found by

$$\lim_{\beta \rightarrow \beta^*} \frac{f(\beta) - \beta^*}{\beta - \beta^*} = \lim_{\beta \rightarrow \frac{1}{\bar{\gamma}}} \frac{\frac{\beta+1}{\bar{\gamma}+1} - \frac{1}{\bar{\gamma}}}{\beta - \frac{1}{\bar{\gamma}}} = \frac{\bar{\gamma}\beta + \bar{\gamma} - \bar{\gamma} - 1}{(\bar{\gamma}\beta - 1)(1 + \bar{\gamma})} \rightarrow \frac{1}{1 + \bar{\gamma}}$$

The AA update is $\beta \mapsto \frac{\beta}{\beta+1} \left(1 + \frac{1}{\bar{\gamma}}\right)$. Its f.p. satisfy $\beta^* = \frac{\beta^*}{\beta^*+1} \left(1 + \frac{1}{\bar{\gamma}}\right)$, so that $\beta^* = 0$ or $\frac{1}{\bar{\gamma}}$. The derivative of the map is positive, and $1 + \frac{1}{\bar{\gamma}} > 1$, so that 0 is unstable and $\frac{1}{\bar{\gamma}}$ is stable. The linear rate is

$$(f'(0) = 1 + \frac{1}{\bar{\gamma}} > 1) \nearrow \quad (f'(\frac{1}{\bar{\gamma}}) = \frac{1}{1 + \frac{1}{\bar{\gamma}}} < 1) \nearrow$$

$$\lim_{\beta \rightarrow \frac{1}{\bar{\gamma}}} \frac{\frac{\beta}{\beta+1} \left(1 + \frac{1}{\bar{\gamma}}\right) - \frac{1}{\bar{\gamma}}}{\beta - \frac{1}{\bar{\gamma}}} = \frac{\beta(\bar{\gamma}+1) - (\beta+1)}{(\beta+1)(\bar{\gamma}\beta - 1)} = \frac{1}{\beta+1} \rightarrow \frac{\bar{\gamma}}{1 + \bar{\gamma}}$$

The IEM update is $\beta \mapsto \frac{\beta}{\beta+1} \frac{1}{\bar{\gamma}} + \frac{1}{1 + \bar{\gamma}}$. Its f.p. satisfy $\beta^* = \frac{\beta^*}{\beta^*+1} \frac{1}{\bar{\gamma}} + \frac{1}{1 + \bar{\gamma}}$, so that β^* satisfies a quadratic eqn. with solns. $\beta^* = \frac{1}{\bar{\gamma}}, -\frac{\bar{\gamma}}{1 + \bar{\gamma}}$. Because negative β are inadmissible, and the f.p. at $\frac{1}{\bar{\gamma}}$ is stable, we have

$$(f'(\frac{1}{\bar{\gamma}}) = \frac{\bar{\gamma}}{(1 + \bar{\gamma})^2} < 1) \nearrow$$

the desired result. The linear rate is

$$\lim_{\beta \rightarrow \frac{1}{\bar{\gamma}}} \frac{\frac{\beta}{\beta+1} \frac{1}{\bar{\gamma}} + \frac{1}{1 + \bar{\gamma}} - \frac{1}{\bar{\gamma}}}{\beta - \frac{1}{\bar{\gamma}}} = \frac{\beta\bar{\gamma} + \bar{\gamma} - \bar{\gamma} - 1}{(\beta+1)(\bar{\gamma}+1)(\beta\bar{\gamma} - 1)} = \frac{1}{(\beta+1)(\bar{\gamma}+1)} \rightarrow \frac{\bar{\gamma}}{(1 + \bar{\gamma})^2}$$

Notice that this rate is the product of the individual rates, and so guaranteed to outperform either. When $\bar{\gamma} \ll 1$, the AA method beats the SA, and is on par with IEM. When $\bar{\gamma} \gg 1$, the SA method beats the AA, and IEM $\sim \frac{1}{2(1+\bar{\gamma})}$ is about twice as fast as SA. When $\bar{\gamma} \approx 1$, SA $\sim \frac{1}{2}$, AA $\sim \frac{1}{2}$, and IEM $\sim \frac{1}{4}$; again about twice as fast.