

Convergence is not guaranteed, however, if $g''(x)$ is bounded and does not change sign on $[a, b]$, then we can rescale nonconvergent problems by choosing $\alpha \neq 0$ and letting

$$\phi(x) = \alpha g'(x) + x.$$

This works since $\alpha g'(x) = 0$ iff $g'(x) = 0$.

There are ways to carefully calculate α , however it is often easier to just try a few values. (Ex pg 26).

Sec 2.2 Multivariate Problems.

In a multivariate opt. problem we seek to max/min a real valued function g of a p -dim vector $x = (x_1, \dots, x_p)^T$. At iteration t , $x^{(t)} = (x_1^{(t)}, \dots, x_p^{(t)})^T$

Many of the general principles still apply.

- iterative algorithms
- often take steps based on linearization of g' from Taylor series, secant approx, etc.
- convergence criteria are in the same spirit

Convergence criteria

Need $D(u, v)$, a distance measure for p -dim vectors.
 Ex $D(u, v) = \sum_{i=1}^p |u_i - v_i|$ or $D(u, v) = \sqrt{\sum_{i=1}^p (u_i - v_i)^2}$
 Then we form abs, rel convergence from.

$$D(x^{(t+1)}, x^{(t)}) < \varepsilon \quad \text{or} \quad \frac{D(x^{(t+1)}, x^{(t)})}{D(x^{(t)}, 0)} < \varepsilon$$

L_p norms

Sec. 2.2.1 Newton's Method & Fisher Scoring.

For Newton's method we approx $g(x^*)$ by the Taylor series

$$g(x^*) = g(x^{(t)}) + (x^* - x^{(t)})^T \vec{g}'(x^{(t)}) + (x^* - x^{(t)})^T \vec{g}''(x^{(t)}) (x^* - x^{(t)}) / 2$$

and we max this function by setting the gradient of the r.h.s. equal to zero.

(Recall that the gradient of f at x is $f'(x) = (\frac{df(x)}{dx_1}, \dots, \frac{df(x)}{dx_p})$
 & we have

$$0 = g'(x^{(t)}) + \cancel{g''(x^{(t)})} g''(x^{(t)}) (x^* - x^{(t)})$$

and the algorithm becomes

$$x^{(t+1)} = x^{(t)} - g''(x^{(t)})^{-1} g'(x^{(t)}) \quad \text{Newton's}$$

or similarly. for an MLE we have

$$\theta^{(t+1)} = \theta^{(t)} + I(\theta^{(t)})^{-1} l'(\theta^{(t)}) \quad \text{Fisher Scoring.}$$

Again the Methods are asymptotically equivalent (Ex. pg 33).
 and have similar problems to the univariate case.

Sec 2.2.2. Newton-like methods.

Computation of the Hessian, $g''(x^{(t)}) = \left[\frac{d^2 f(x)}{dx_i dx_j} \right]_{ij}$
 can be expensive & so many methods rely upon
 eg. of the form

$$x^{(t+1)} = x^{(t)} - (M^{(t)})^{-1} g'(x^{(t)})$$

where $M^{(t)}$ is a $p \times p$ approx. (Ex Fisher Scoring).

22.2.1 Ascent Algorithms

With Newton's method the steps are not necess. uphill, i.e. $g(x^{(t+1)}) > g(x^{(t)})$. If we force this, then it is called an ascent algorithm.

What does this remind you of?

Method of steepest ascent: w/ $M^{(t)} = -I$.
 $x^{(t+1)} = x^{(t)} + g'(x^{(t)})$ is taking a step in the steepest direction uphill (indicated by the gradient).

We can also use ~~small~~ scaled steps
 $x^{(t+1)} = x^{(t)} + \alpha^{(t)} g'(x^{(t)})$ for $\alpha^{(t)} > 0$
 to control convergence. or more generally

$$x^{(t+1)} = x^{(t)} + \alpha^{(t)} (M^{(t)})^{-1} g'(x^{(t)}).$$

Often w/ scaling $\alpha^{(t)} > 0$ is chosen to be a contraction or step length parameter whose value can shrink to ensure ascent at each step.

(See ex 2.7 pg 33).

In other words if a step turns out to be downhill, we adjust to ensure uphill.

Ex/ Backtracking:

- Start each step w/ $\alpha^{(t)} = 1$.
- if step is downhill ($g(x^{(t+1)}) < g(x^{(t)})$)
 let $\alpha^{(t)} = \frac{1}{2} \alpha^{(t)}$ & try again
- Repeat until step is uphill.

(See example pg 28)

Backtracking

Will converge under formal conditions, can be slow.

pg 39

Fixed-Point Methods.

If $M^{(t)} = M \forall t$ we have the fixed pt. method

$$x^{(t+1)} = x^{(t)} - M^{-1} g'(x^{(t)})$$

A reasonable choice is $M = g''(x^{(0)})$.
 If M is diagonal, then this is eq. to applying the univariate-scaled fixed-point algorithm to each component.

Secant-Like Methods.

We can replace $g''(x^{(t)})$ w/ a matrix $M^{(t)}$ of finite discrete difference quotients.

Ex/ let $g'_i(x) = dg(x)/dx_i$ (ith element of $g'(x)$)
 and $e_j = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$ ← jth position

Then

$$M_{ij}^{(t)} = \frac{g'_i(x^{(t)} + h_{ij}^{(t)} e_j) - g'_i(x^{(t)})}{h_{ij}^{(t)}}$$

for some constants $h_{ij}^{(t)}$.

If $h_{ij}^{(t)} = h$ we get convergence order 1.

If $h_{ij}^{(t)} = x_j^{(t)} - x_j^{(t-1)} \forall i$ we get convergence order similar to secant method in univariate case.

There are quite a few other methods discussed in the text that we do not have time to cover.

Pg 89 Chapter 4 EM (Expectation Maximization) Method

The EM algorithm is an iterative opt. strategy motivated by a notion of missingness & by consideration of the cond. dist of what is missing given what is observed.

Assume: we have observed data from r.v. X along with missing (unobserved or latent) data from r.v. Z . We wish to envision complete data from $Y = (X, Z)$.

Given observed data x we want to maximize a likelihood $L(\theta|x)$, but we want to do so w/o calc $L(\theta|x)$ directly but rather working with $Y|\theta \propto Z|(x, \theta)$.
 $L(\theta|Y)$ which may be more tractable.

Notation.

X observed data
 Y complete data
 Z missing data
 $f_X(x|\theta)$ density of observed data
 $f_Y(y|\theta)$ density of complete data
 M be the many to fewer mapping $X = M(Y)$

Intuitively EM:

- ① Fills in z based on $x \propto \theta$
- ② Restimates θ based on $Y = (x, z)$

Then the missing data amounts to a marginalization model in which we observe X having density

$$f_X(x|\theta) = \int_{Z: M(Y)=x} f_Y(y|\theta)$$

And the conditional density ~~of~~ of the missing z given observed x is

$$f_{z|x}(z|x, \theta) = \frac{f_Y(y|\theta)}{f_X(x|\theta)}$$

And similarly we will view our likelihood $L(\theta|x)$ as a marginalization of the complete data Likelihood $L(\theta|y) = L(\theta|x, z)$.

Sec 4.2. The EM Algorithm

Let $\theta^{(t)}$ be our est at iteration $t = 0, 1, \dots$.
Define $Q(\theta|\theta^{(t)})$ to be the expectation for the joint log likelihood of Y conditioned on $X=x$

$$\begin{aligned} Q(\theta|\theta^{(t)}) &= E[\log L(\theta|Y) | x, \theta^{(t)}] \\ &= E[\log f_Y(y|\theta) | x, \theta^{(t)}] \\ &= \int (\log f_Y(y|\theta)) f_{z|x}(z|x, \theta^{(t)}) dz \end{aligned}$$

Recall Z is the only random part of Y once $X=x$.

Then the EM algorithm is (starting w/ $\theta^{(0)}$)

- ① E step: Compute $Q(\theta|\theta^{(t)})$
- ② M step: Maximize $Q(\theta|\theta^{(t)})$ w.r.t. θ and set $\theta^{(t+1)}$ equal to this maximizer.
- ③ Return to E step unless stopping criteria has been met.

Step ② \rightarrow use earlier methods.

Stopping criteria similar to before: built upon $|Q(\theta^{(t+1)}|\theta^{(t)}) - Q(\theta^{(t)}|\theta^{(t)})|$ or $(\theta^{(t+1)} - \theta^{(t)})^T (\theta^{(t+1)} - \theta^{(t)})$.

Ex/4.1 Very Simple

$y_1, y_2 \sim \text{i.i.d. Exp}(\theta)$ so $f(y_i) = \theta e^{-\theta y_i}$ $E[y_i] = \frac{1}{\theta}$

Suppose $y_1 = 5$ is observed (x) and y_2 is missing (z) $\therefore Y = (y_1, y_2)$.

- Write down expression for complete density

$$f_Y(Y|\theta) = \theta e^{-\theta y_1} \theta e^{-\theta y_2}$$

also useful.

$$\begin{aligned} f_{z|x}(z|x, \theta^{(t)}) \\ = f_{z|y}(z|y, \theta^{(t)}) \\ = \theta e^{-\theta z} \end{aligned}$$

- Write down the log likelihood function of the complete Y
 $\log L(\theta|y) = \log f_Y(Y|\theta) = 2 \log \theta - \theta y_1 - \theta y_2$

- Find $Q(\theta|\theta^{(t)}) = E[\log L(\theta|y)] \big| y_1, \theta^{(t)}$

$$= E[2 \log \theta - \theta y_1 - \theta y_2 \big| y_1, \theta^{(t)}]$$

$$= 2 \log \theta - \theta \cdot 5 - \theta E[y_2 | \theta^{(t)}]$$

$$= 2 \log \theta - 5\theta - \theta / \theta^{(t)}$$

Now we need to maximize $Q(\theta|\theta^{(t)})$
 by solving for the root of $Q'(\theta|\theta^{(t)})$
 $= 2/\theta - 5 - 1/\theta^{(t)} = 0.$

$$\text{Hence } \theta^{(t+1)} = \frac{2(\theta^{(t)})}{5(\theta^{(t)} + 1)}.$$

and repeat.

Further comments:

- One of the most appealing & central results is that the sequence $\{\theta^{(k)}\}$ converges, at least to a local ~~minimum~~ maximum. proof pg 95. and for well behaved problems θ is a global max.

- We move uphill at each step. - The increase in the observed data log likelihood function $\ell(\theta|x)$ at each step is one of its most attractive features v/s Newton's Method.
- Rate of convergence: The rate is only linear (Newton quadratic) \therefore is criticized as being slow. convergence rate linked to proportion of data missing - more missing \Rightarrow slower convergence
However there are many techniques to speed up EM. (Sec 4.3).
- Starting Points: A main drawback is that its limiting position is often sensitive to initial guesses.

Ex 4.2 more complicated example.

C carbonia CC CI CT
I insularia II IT
T typica TT

$$p_C, p_I, p_T \Rightarrow \begin{matrix} p_C^2 & 2p_C p_I & 2p_C p_T & p_I^2 & 2p_I p_T & p_T^2 \\ C & CI & CT & II & IT & TT \end{matrix}$$

$$\text{Also } p_C + p_I + p_T = 1$$

$$n = n_C + n_I + n_T \quad \Sigma \text{ observed} = (n_C, n_I, n_T)$$

$$Y = \text{complete data} = (n_{CC}, n_{CI}, n_{CT}, n_{II}, n_{IT}, n_{TT})$$

$$\omega / X = M(Y) = (n_{CC} + n_{CI} + n_{CT}, n_{II} + n_{IT}, n_{TT})$$

Since $P_T = 1 - P_C - P_I$

$$P = (P_C, P_I).$$

1) Multinomial R.V. • complete log likelihood function.

$$f_y(y|P) = \frac{n!}{n_C! n_I! n_T!} (P_C^2)^{n_C} (2P_C P_I)^{n_I} (2P_C P_T)^{n_T} \cdot (P_I^2)^{n_{II}} (2P_I P_T)^{n_{IT}} (P_T^2)^{n_{TT}}$$

∴

$$\begin{aligned} \log f_y(y|P) &= n_C \log P_C^2 + n_I \log 2P_C P_I + n_T \log 2P_C P_T \\ &\quad + n_{II} \log P_I^2 + n_{IT} \log 2P_I P_T + n_{TT} \log P_T^2 \\ &\quad + \log(n_C! n_I! n_T!) \end{aligned}$$

The complete data is

$$Y = (N_C, N_I, N_T, N_{II}, N_{IT}, N_{TT})$$

and only

$$N_{II} + N_{IT} + N_{TT} = n_T \text{ is observed.}$$

$$\text{To find } Q(P|P^{(t)}) = E[\log f_y(y|P) | x, P^{(t)}]$$

we need

$$E[N_{??} | n_C, n_I, n_T, P^{(t)}] \text{ for each ?? pair}$$

$$E[N_C | x, P^{(t)}] = n_C^{(t)} = \frac{n_C (P_C^{(t)})^2}{(P_C^{(t)})^2 + 2P_C^{(t)} P_I^{(t)} + 2P_C^{(t)} P_T^{(t)}}$$

$$E[N_I | x, P^{(t)}] = n_I^{(t)} = \frac{n_C (2P_C^{(t)} P_I^{(t)})}{(P_C^{(t)})^2 + 2P_C^{(t)} P_I^{(t)} + 2P_C^{(t)} P_T^{(t)}}$$

$$E[N_T | x, P^{(t)}] = n_T^{(t)} = \frac{n_C (2P_C^{(t)} P_T^{(t)})}{(P_C^{(t)})^2 + 2P_C^{(t)} P_I^{(t)} + 2P_C^{(t)} P_T^{(t)}}$$

Each phenotype has its own dist.
Multinomial!
 $E[x_i] = n p_i$

$$E[N_{II} | X, p^{(t)}] = n_{II}^{(t)} = \frac{n_I (p_I^{(t)})^2}{(p_I^{(t)})^2 + 2p_I^{(t)} p_T^{(t)}}$$

$$E[N_{IT} | X, p^{(t)}] = n_{IT}^{(t)} = \frac{n_I 2(p_I^{(t)} p_T^{(t)})}{(\Delta)}$$

$$E[N_{TT} | X, p^{(t)}] = n_{TT} = n_T$$

$$\text{and } E\left[\binom{n}{n_{cc} \dots n_{pp}}\right] = k(n_{cc}, n_{II}, n_{IT}, n_{TT}; p_T)$$

So

does not depend on p .

$$Q(p | p^{(t)}) = n_{cc}^{(t)} \log p_c^2 + n_{II}^{(t)} \log 2p_I p_c + n_{IT}^{(t)} \log 2p_I p_T \\ + n_{II}^{(t)} \log p_I^2 + n_{IT}^{(t)} \log p_I p_T + n_{TT} \log p_T^2$$

and so we need to take partials w.r.t p_c & p_I

$$\frac{dQ(p | p^{(t)})}{dp_c} = \frac{2n_{cc}^{(t)} + n_{II}^{(t)} + n_{IT}^{(t)}}{1 - p_c - p_I} \\ - \frac{n_{IT}^{(t)} + n_{II}^{(t)} + 2n_{TT}^{(t)}}{1 - p_c - p_I}$$

$$\frac{dQ(p | p^{(t)})}{dp_I} = \frac{2n_{II}^{(t)} + n_{IT}^{(t)} + n_{IT}^{(t)}}{1 - p_c - p_I} \\ - \frac{2n_{IT}^{(t)} + n_{IT}^{(t)} + n_{IT}^{(t)}}{1 - p_c - p_I}$$

We need to set these to 0 & solve.

$$p_c^{(t+1)} = \frac{2n_{cc}^{(t)} + n_{II}^{(t)} + n_{IT}^{(t)}}{2n}$$

$$p_I^{(t+1)} = \frac{2n_{II}^{(t)} + n_{IT}^{(t)} + n_{IT}^{(t)}}{2n}$$

$$p_T^{(t+1)} = \frac{2n_{TT}^{(t)} + n_{IT}^{(t)} + n_{IT}^{(t)}}{2n}$$

Don't actually
need missing info
ex algorithm
run on
pg 94.

all the exp
are numbers

Final EM Example

Bayesian posterior mode w/ EM.

Review: Bayesian Inference Sec 1.5 pg 11

Consider a Bayesian problem with

- likelihood $L(\theta|x)$
- prior $f(\theta)$ - density assigned to θ before observing the data
- missing data or parameters Z

We wish to find the posterior density

$$f(\theta|x) = f(\theta) f(x|\theta) \cdot K = f(\theta) L(\theta|x) \cdot K$$

where K is a normalizing constant.

To use the EM method, the E-step requires

$$\begin{aligned} Q(\theta | \theta^{(t)}) &= E[\log\{L(\theta|Y)f(\theta) \cdot K(Y) | x, \theta^{(t)}\}] \\ &= E[\log L(\theta|Y) | x, \theta^{(t)}] + \log f(\theta) \\ &\quad + E[\log K(Y) | x, \theta^{(t)}] \end{aligned}$$

where the last term can be ignored when maximizing w.r.t. θ .

So, we have our standard MLE Q with the addition of the log prior.