

Secant method

If  $f'$  is difficult to compute (almost always) ~~or~~ we are lazy (always) the secant method provides an approximation.

Newton step/secant method

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

$$x_{n+1} = x_n - \frac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})}$$

$$\Rightarrow f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$$

~~Pro~~ Note: Need 2 starting values.

Pro: Easy as pie

Con: Convergence only super linear.

Statistics application

$\ell(\theta): \mathbb{R}^K \rightarrow \mathbb{R}$ , log-likelihood,  $\theta = (\theta_1, \dots, \theta_K)$

$\ell'(\theta): \mathbb{R}^K \rightarrow \mathbb{R}^K$ , ~~score~~ gradient

$\ell''(\theta): \mathbb{R}^K \rightarrow \mathbb{R}^{K \times K}$ , hessian

$$x_n - x_{n-1} =$$

$$x_n - x_{n-1} + x_{n-1} - x_{n-2}$$

$$\varepsilon_n - \varepsilon_{n-1}$$

$$\varepsilon_{n+1} = \varepsilon_n -$$

$$\frac{f(x_n + \varepsilon_n)(x_n - x_{n-1})}{f(x_n + \varepsilon_n) - f(x_{n-1} + \varepsilon_{n-1})}$$

$$\frac{\varepsilon_n}{\left( \varepsilon_n f' + \frac{\varepsilon_n^2}{2} f'' \right) (x_n - x_{n-1})}$$

$$\frac{\varepsilon_n f' + \frac{\varepsilon_n^2}{2} f'' - \varepsilon_{n-1} f' - \frac{\varepsilon_{n-1}^2}{2} f''}{\varepsilon_n f' + \frac{\varepsilon_n^2}{2} f'' - \varepsilon_{n-1} f' - \frac{\varepsilon_{n-1}^2}{2} f''}$$

$$= \varepsilon_n -$$

$$\frac{(\varepsilon_n f' + \frac{\varepsilon_n^2}{2} f'')(x_n - x_{n-1})}{(\varepsilon_n - \varepsilon_{n-1}) f' - \left( \frac{\varepsilon_n^2}{2} - \frac{\varepsilon_{n-1}^2}{2} \right) f''}$$

## Fisher Scoring for GLM,

$$Y \sim \mathcal{G}(\mu) \quad , \quad p(Y|\mu) \quad \mathbb{E}[Y] = \mu$$

$$g(\mu) = X\beta \quad \text{Var}(Y) = V(\mu)$$

$$g(y) \approx g(\mu) + (y - \mu)g'(\mu)$$

↓

$z$  (working response)

$$\eta = g(\mu)$$

① Start with  $\hat{\mu}_i$

② Set  $z_i = \eta_i + (y_i - \hat{\mu}_i) \left( \frac{\partial \eta_i}{\partial \mu_i} \right)$

③ Weighted regression of  $z$  on  $X$ , i.e. solve

$$X^T W X \beta_n = X^T W z$$

$$\beta_n = (X^T W X)^{-1} X^T W z$$

where

$$W = \left[ \left( \frac{\partial \eta_i}{\partial \mu_i} \right)^2 V(\mu_i) \right]^{-1}$$

ex:  $y_i \sim \text{Poisson}(\mu_i)$

$$g(\mu_i) = \log \mu_i = \eta_i$$

$$\frac{\partial \eta_i}{\partial \mu_i} = \frac{1}{\mu_i} \quad V(\mu_i) = \mu_i$$

① set  $\hat{\mu}_i$

②  $z_i = \log \mu_i + (y_i - \mu_i) / \mu_i$

③ Regress  $z$  on  $X$  where

$$W = \left[ \frac{1}{\hat{\mu}_i^2} \hat{\mu}_i \right]^{-1} = \hat{\mu}_i$$

Poisson regression via Newton

$$L(\beta) \propto e^{-\mu} \mu^y$$

$$\begin{aligned} \ell(\beta) &= y \log \mu - \mu \\ &= y X\beta - \exp(X\beta) \end{aligned}$$

$$\ell'(\beta) = Xy - X \exp(X\beta)$$

$$\ell''(\beta) = -X \underbrace{\exp(X\beta)}_{\mu} X$$

$$\left\{ \begin{array}{l} y^T X\beta - \exp(X\beta) \\ X^T y - X^T \mu \\ X^T (y - \mu) \end{array} \right.$$

$$X^T y - X^T \mu$$

$$-X^T W X$$

$$\beta_{n+1} = \beta_n + (-X^T W X)^{-1} (-X^T (y - \mu))$$

$$= \beta_n + (X^T W X)^{-1} X^T (\mu z - \mu \eta)$$

$$= \beta_n + (X^T W X)^{-1} X^T W_n (z_n - \eta_n)$$

$$= (X^T W X)^{-1} X^T W z + \underbrace{\left( \beta_n - (X^T W X)^{-1} X^T W \eta \right)}_{=0}$$

$$= (X^T W X)^{-1} X^T W z$$

$$= 0$$

$$= \log \mu$$

$$= X\beta$$

$$z = \eta + (y - \mu) / \mu$$

$$\mu z = \mu \eta + y - \mu$$

$$y - \mu = \mu z - \mu \eta$$

$$\eta = g(\mu) = \log \mu$$

$$\eta = z - \frac{y - \mu}{\mu}$$

$$X\beta$$

$$\mu = \exp$$

$$\log =$$

## Summary of Minimization

$$\min_x f(x) \quad \text{for } f: \mathbb{R}^k \rightarrow \mathbb{R}, \quad x \in \mathbb{R}^k$$

### Line Search: ~~Steepest~~

#### ① Steepest descent

$$x_{n+1} = x_n + \alpha [-f'(x_n)]$$

$\swarrow$  direction of steepest descent  
 $\searrow$  scalar step length

$\Rightarrow$  linear convergence

#### ② Newton

$$x_{n+1} = x_n + 1 \underbrace{[f''(x_n)]^{-1} [-f'(x_n)]}_{\text{Newton direction}}$$

step length = 1

$\Rightarrow$  quadratic convergence

$\Rightarrow$  in stat, estimate of asymptotic covariance

#### ③ Quasi-Newton

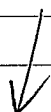
$$x_{n+1} = x_n + \alpha B_n [-f'(x_n)]$$

$$B_n = \arg \min_B \|B - B_{n-1}\|$$

$$B = B^T \quad \text{or} \quad f'(x_n) - f'(x_{n-1}) = B_n (x_n - x_{n-1})$$

# ④ Conjugate Gradient ("modified steepest descent")

$$x_{n+1} = x_n + \alpha p_n$$



$$p_{n+1} = -f'_n(x_{n+1}) + \frac{\|f'_n(x_{n+1})\|^2}{\|f'_n(x_n)\|^2} p_n$$

$$p_0 = -f'(x_0)$$

~~$$p_{n+1} = f_0'^T f_1 + \frac{\|f_1'\|^2}{\|f_0'\|^2} p_0^T f_0$$~~

## Coordinate Descent:

$$x_{n+1}^{(k)} = \arg \min_{x^{(k)}} f(x_n^{(1)}, x_n^{(2)}, \dots, x_n^{(k)}, \dots, x_n^{(p)})$$

$\Rightarrow$  has "descent property", non-increasing

Do this  
⇒

Method of steepest descent/ascent uses

$$\theta_{n+1} = \theta_n - \alpha l'(\theta_n)$$

Where  $\alpha$  is chosen so that  $\theta_{n+1}$  has a larger likelihood value than  $\theta_n$ .

Summary:

Knowledge of  $l''$

✓

① Newton: Fastest (quadratic convergence), requires calculating  $l''$ , gives asymptotic  $\text{Var}(\hat{\theta})$ , unstable if starting value too far.

② Score: Superlinear convergence, but equiv. to Newton if  $l''$  does not depend on  $y$  (true in many common cases), requires calculating  $E[l'']$ , can be unstable, but often quite stable in typical statistics apps.

③ Quasi-Newton: Superlinear convergence, does not require  $l''$  or  $E[l'']$ , more stable than Newton, No estimate of  $\text{Var}(\hat{\theta})$

④ ~~Steepest descent~~: ~~linear convergence~~, ~~very stable~~

④ Steepest descent: linear convergence, stable

X

## The EM algorithm

EM stands for Expectation-Maximization.

Originally outlined in DLR (1977) but ideas go much further back. DLR unified many different ideas and gave examples of the broad applicability of EM.

EM is not an "algorithm". It is an algorithm for creating other algorithms.

## Generalized Additive Models (Hastie + Tibshirani)

Usual linear model has:

~~$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_p X_{pi} + \epsilon_i$$~~

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_p X_{pi} + \epsilon_i$$

GAMS say:

$$Y_i = \alpha + S_1(X_{1i}) + S_2(X_{2i}) + \dots + S_p(X_{pi}) + \epsilon_i$$

$$\Rightarrow Y = \alpha + \sum_{j=1}^p S_j(X_j) + \epsilon \quad (*)$$

where  $S_j()$  are smooth.

$S_j()$  can be any kind of "smoother", even a mixture of different kinds

ex. Smoothers: parametric splines, smoothing splines, penalized splines, loess, running lines, running median



GAM algorithm:  $\rightarrow Y = (Y_1, \dots, Y_n)$

Given model  $Y = \alpha + \sum_{j=1}^p S_j(X_j) + \varepsilon$

① Initialize  $\alpha = \frac{1}{n} \sum_{i=1}^n Y_i$ ,  $s_j = s_j^0 = 0$

② For  $j=1, 2, \dots, p$

Let  $r_j = Y - \alpha - \sum_{k \neq j} s_k(X_k)$

where  ~~$s_k = (s_k(X_{k1}), s_k(X_{k2}), \dots, s_k(X_{kn}))$~~

where  $s_k = (s_k(X_{k1}), s_k(X_{k2}), \dots, s_k(X_{kn}))$

i.e.  $s_k()$  evaluated at data points for  $X_k$

~~$s_j$~~   $s_j = \text{Smooth}(r_j | X_j)$

$\hat{s}_j = s_j - \sum_{i=1}^n s_j(X_{ij})$

③ Evaluate  $\Delta = \sum_{j=1}^p \|s_j - s_j^0\|$

or  $\Delta = \frac{\sum_{j=1}^p \|s_j - s_j^0\|}{\sum \|s_j^0\|}$

If  $\Delta < \varepsilon$ , stop

Else set  $s_j^0 = s_j$  for  $j=1 \rightarrow p$ .

~~Goto~~ Goto ②

Local Scoring

Set adjusted response

$$z_i = \eta_i + (y_i - \mu_i) \frac{\partial \eta_i}{\partial \mu_i}$$

Fit an additive model for

$z \sim X_1 \dots X_p$  with observation weights

$$W_i = \left[ \frac{\partial \eta_i}{\partial \mu_i} \right]^{-2} V(\mu_i)^{-1}$$

GAM algorithm = "backfitting"

~ Alternating conditional expectation  
~ Cyclic coordinate descent

⇒ linear convergence algorithm

⇒ sidesteps curse of dimensionality by additivity constraint (big!)

## EM Algorithm

EM stands for Expectation-~~Parameter~~ Maximization.

Originally by DLR, 1977 but ideas go much further back. DLR united many different ideas and put them in a statistical framework.

EM is not strictly an "algorithm". It is an (abstract) algorithm for creating other algorithms.

The basic principle of EM is straightforward.

We observe some data  $Y$  ~~but~~ but there are some data that are unavailable or "missing". Call these data  $Z$ .

The observed data  $Y$  with the missing data  $Z$  are the complete data  $X = (Y, Z)$  <sup>complete data</sup>

- ① The complete data have a joint density  $g(Y, Z | \theta)$
- ② Because of missing  $Z$ , we cannot evaluate  $g$ . We observe  $Y$  with a joint density

$$f(Y | \theta) = \int g(Y, Z | \theta) dZ$$

$$\ell(\theta | Y) = \log f(Y | \theta)$$

- ③  $\ell(\theta | Y)$  is hard to evaluate! (because of  $\int$ )  
(although, maybe not!)

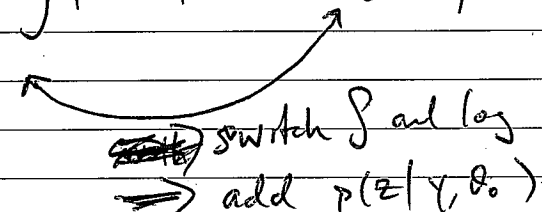
meta-algorithm?

Direct ML maximizes  $\ell(\theta|y)$ . This may be possible!

Heuristically, the EM algorithm is as follows:

① E-step. Given estimate  $\theta = \theta_0$

$$\begin{aligned} \text{Define } Q(\theta|\theta_0) &= E[\log g(y, z|\theta) | y, \theta_0] \\ &= \int p(z|y, \theta_0) \log g(y, z|\theta) dz \end{aligned}$$



② Maximize  $Q(\theta|\theta_0)$  wrt  $\theta$ , ~~and set~~

~~the set~~  $\hat{\theta}_1 = \arg \max_{\theta} Q(\theta|\theta_0)$   
~~Set~~  $\theta_0 = \hat{\theta}_1$  and go to ①

~~Create a Sequence~~

$$\text{Set } \theta_{n+1} = \arg \max_{\theta} Q(\theta|\theta_n)$$

Under broad assumptions,  $\theta_n \xrightarrow[n \rightarrow \infty]{} \hat{\theta}$ , MLE

(More later)

$$E(\log X) = \int \log x f(x) dx$$

Ex

$$Y_1 \rightarrow y_m, Z_{m+1} \rightarrow z_n \sim \text{Poisson}(\mu)$$

$$g(x|\mu) \propto \prod_{i=1}^n x_i e^{-\mu}$$

$$\log g(x|\mu) = \sum_{i=1}^n X_i \log \mu - \mu$$

$$= \log \mu \sum X_i - n \mu$$

$$\mathbb{E} \left[ \sum_{i=1}^n X_i \mid y \right] = \sum_{i=1}^m y_i + \sum_{i=m+1}^n \tilde{\mu}$$

$$\mathbb{E}_{\mu} \left[ \sum_{i=1}^n X_i \right] = \mu n$$

Ex: Censored exponential data

$Y_1, \dots, Y_n \sim \text{Exp}(\lambda)$  but some cases are censored on right

Let ~~complete~~ observed data be

$$(Y_1, \dots, Y_n) \rightarrow (\min(Y_i, c_i), \delta_i), \dots, (\min(Y_n, c_n), \delta_n)$$

where  $\delta_i = 1$  if  $Y_i \leq c_i$  and

$= 0$  if  $Y_i$  censored

$$g(x|\lambda) \propto \prod_{i=1}^n \frac{1}{\lambda} \exp(-x_i/\lambda)$$

$$\log g(x|\lambda) = -n \log \lambda - \frac{1}{\lambda} \sum x_i$$

$$= -n \log \lambda - \frac{1}{\lambda} \left\{ \sum_{\text{obs}} Y_i + \sum_{\text{censored}} Z_i \right\}$$

$$E[\log g(x|\lambda) | Y, \tilde{\lambda}]$$

$$= -n \log \lambda - \frac{1}{\lambda} \left[ \sum_{\text{obs}} Y_i + \sum_{\text{m.c.s}} c_i + \frac{1}{\lambda} \right]$$

$$E[Z_i | Y_i, c_i]$$

$$\Rightarrow \hat{\lambda} = \frac{1}{n} \left[ \sum_{\text{obs}} Y_i + \sum_{\text{m.c.s}} c_i + \frac{1}{\lambda} \right]$$

One sample

$$Y_1, \dots, Y_n \sim N(\mu, \sigma^2)$$

$$Y_1, Y_2, Y_3, Y_4, \dots, Y_n \sim N(\mu, \sigma^2)$$

$$L(\mu, \sigma^2) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (Y_i - \mu)^2\right)$$

$$\ell = \sum_{i=1}^n -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum (Y_i - \mu)^2$$

$$= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum (Y_i - \mu)^2$$

$$\frac{\partial \ell}{\partial \mu} = -\frac{1}{\sigma^2} \sum (Y_i - \mu) \quad \begin{matrix} Y_i^2 - 2\mu Y_i + \mu^2 \\ \sigma_0^2 + \mu_0^2 - 2\mu \mu_0 - \mu^2 \\ \sum Y_i^2 - 2\mu \sum Y_i + n\mu^2 \end{matrix}$$

$$E[Y_i | \mu_0, \sigma_0^2] = \mu_0$$

$$E[Y_i^2] = \sigma^2 + \mu^2$$

Bivariate

$$Y_i = \begin{pmatrix} Y_{i1} \\ Y_{i2} \end{pmatrix} \sim N \left( \underbrace{\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}}_{\mu}, \underbrace{\begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ 0 & \sigma_2^2 \end{pmatrix}}_V \right)$$

$$L \propto \prod_{i=1}^n |V|^{-1/2} \exp\left(-\frac{1}{2} (Y_i - \mu)^T V^{-1} (Y_i - \mu)\right)$$

$$\ell = \sum_{i=1}^n -\frac{1}{2} \log |V| - \frac{1}{2} (Y_i - \mu)^T V^{-1} (Y_i - \mu)$$

$$= -\frac{n}{2} \log |V| - \frac{1}{2} \text{tr} \left( (Y - \mu) V^{-1} (Y - \mu)^T \right)$$

$$= -\frac{1}{2} \text{tr} \left( (Y - \mu) (Y - \mu)^T V^{-1} \right)$$

$$= -\frac{1}{2} \text{tr} \left( (Y Y^T - Y \mu^T - \mu Y^T + \mu \mu^T) V^{-1} \right)$$

$$n \times 2 \quad 2 \times 2 \quad 2 \times n$$

$$n \times 2 \quad n \times 2 \quad 2 \times 2$$

$$n \times 2 \quad 2 \times n$$

$$n \times n$$

$$l = -\frac{n}{2} \log |V| - \frac{1}{2} \text{tr} \left( (Y - \mu)^T (Y - \mu) V^{-1} \right) \\ - \frac{1}{2} \text{tr} \left( (Y^T Y - Y^T \mu - \mu^T Y + \mu^T \mu) V^{-1} \right)$$

$$\begin{pmatrix} \sum_{i=1}^n y_{i1}^2 & \sum_{i=1}^n y_{i1} y_{i2} \\ \sum_{i=1}^n y_{i1} y_{i2} & \sum_{i=1}^n y_{i2}^2 \end{pmatrix}$$

$$y_{ij}^{*2} = \left[ \hat{\mu}_1 + \hat{\rho} (y_{i2} - \hat{\mu}_2) \right]^2 + (1 - \hat{\rho}^2) \hat{\sigma}_2^2$$

$$= (E y_{i1})^2 + \text{Var}(y_{i1}) \left| \hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}_1^2, \hat{\sigma}_2^2, \hat{\rho} \right|$$

$$E[y_{i1}^* y_{i2}^* | v] = E[y_{i1}^*] E[y_{i2}^*] + \hat{\rho} \hat{\sigma}_1 \hat{\sigma}_2$$



Ms class  
11/26

$$\log f(y|\theta) = \log \int g(y, z|\theta) dz$$

$$\log f(y|\theta) - \log f(y|\theta_0)$$

$$= \log \int g(y, z|\theta) dz - \log \int g(y, z|\theta_0) dz$$

$$= \log \frac{\int g(y, z|\theta) dz}{\int g(y, z|\theta_0) dz}$$

$$= \log \frac{\int g(y, z|\theta_0) \frac{g(y, z|\theta)}{g(y, z|\theta_0)} dz}{\int g(y, z|\theta_0) dz}$$

$$= \log \int h(z|y, \theta_0) \frac{g(y, z|\theta)}{g(y, z|\theta_0)} dz$$

$$\geq \int h(z|y, \theta_0) \log \frac{g(y, z|\theta)}{g(y, z|\theta_0)} dz$$

$$= \int h(z|y, \theta_0) \log g(y, z|\theta) dz$$

$$- \int h(z|y, \theta_0) \log g(y, z|\theta_0) dz$$

$$\log f(y|\theta) \geq \log f(y|\theta_0) + \mathbb{E}_h[\log g(y, z|\theta)] - \mathbb{E}_h[\log g(y, z|\theta_0)]$$

$f(y|\theta)$  obs.

$g(y, z|\theta)$  complete

Ascent property of EM algorithm

The sequence  $\theta_{n+1} = \arg \max_{\theta} Q(\theta|\theta_n)$  has the property that for each  $n$ ,

$$\ell(\theta_{n+1}|y) \geq \ell(\theta_n|y)$$

$$\Rightarrow \log f(y|\theta_{n+1}) \geq \log f(y|\theta_n)$$

w/ strict inequality when  $Q(\theta_{n+1}|\theta_n) > Q(\theta_n|\theta_n)$

~~Define  $D(f||g) = -\log E$~~

Define  $D(f||g) = E_f[\log f/g]$

$D(f||g) \geq 0 \Leftrightarrow f = g$

$D(f||g) \geq 0$  (information inequality)

Pf:

$$\begin{aligned} D(f||g) &= E_f[\log f/g] \\ &= E_f[-\log g/f] \end{aligned}$$

Jensen's inequality,  $\geq -\log E_f[g/f]$

( $-\log$  is a convex function)

$$= -\log \int g/f \cdot f$$

$$= -\log 1$$

$$= 0$$

$D(\cdot||\cdot)$  is like a distance (but no  $\Delta$  inequality)

~~Q~~

Recall

$$Q(\theta_{n+1}|\theta_n) \geq Q(\theta_n|\theta_n)$$

$$\begin{aligned} \ell(\theta_{n+1}^*) - \cancel{\ell(\theta_{n+1}^*)} &= \log f(y|\theta_{n+1}) - \cancel{\log f(y|\theta_{n+1})} \\ &= \log f(y|\theta_{n+1}) + \underbrace{Q(\theta_{n+1}|\theta_n) - Q(\theta_{n+1}|\theta_n)}_0 \end{aligned}$$

$$= Q(\theta_{n+1}|\theta_n) - [Q(\theta_{n+1}|\theta_n) - \log f(y|\theta_{n+1})]$$

$$= Q(\theta_{n+1}|\theta_n) - [\mathbb{E}_p[\log g(y, z|\theta_{n+1}) | y, \theta_n] - \log f(y|\theta_{n+1})]$$

$$= Q(\theta_{n+1}|\theta_n) - \mathbb{E}_p \left[ \log \frac{g(y, z|\theta_{n+1})}{f(y|\theta_{n+1})} \middle| y, \theta_n \right]$$

$$= Q(\theta_{n+1}|\theta_n) - \mathbb{E}_p [\log p(z|y, \theta_{n+1}) | y, \theta_n]$$

$\hookrightarrow p(z|y, \theta_n)$

$$\geq Q(\theta_n|\theta_n) - \mathbb{E}_p [\log p(z|y, \theta_{n+1}) | y, \theta_n]$$

$$\geq Q(\theta_n|\theta_n) - \mathbb{E}_p [\log p(z|y, \theta_n) | y, \theta_n]$$

$$= Q(\theta_n|\theta_n) - \mathbb{E}_p \left[ \log \frac{g(y, z|\theta_n)}{f(y|\theta_n)} \middle| y, \theta_n \right]$$

$$= Q(\theta_n|\theta_n) - \mathbb{E}_p [\log g(y, z|\theta_n) | y, \theta_n] + \log f(y|\theta_n)$$

$$= Q(\theta_n|\theta_n) - Q(\theta_n|\theta_n) + \log f(y|\theta_n)$$

$$= \log f(y|\theta_n)$$

$$= \ell(\theta_n)$$

assumption  
by information  
inequality

2

Roughly,

If  $l$  is ~~convex~~ <sup>concave</sup> then  $\{l(\theta_n)\}$  is  
a monotone increasing, bounded, seq. of numbers.

$\Rightarrow$  there is a limit.

~~But does  $\theta_n \rightarrow \hat{\theta}$~~

But does this mean  $\theta_n \rightarrow \hat{\theta}$ , MLE?

Not necessarily, but yes for exponential families.

For exp. fam., there is always a unique maximizer.

Note: We do not require that

$$\theta_{n+1} = \arg \max_{\theta} Q(\theta | \theta_n).$$

We only need  $Q(\theta_{n+1} | \theta_n) \geq Q(\theta_n | \theta_n)$ .

$\hookrightarrow$  This algorithm is Generalized EM (GEM).

Ex. ~~Exponential~~ Normal w/ missing data

$$\mathbf{X} = (X_1, X_2) \sim N(\mu, \Sigma)$$

$$\mu = (\mu_1, \mu_2), \quad \Sigma = \begin{bmatrix} \sigma^2 & v \\ v & \tau^2 \end{bmatrix}$$

$$\begin{bmatrix} X_1 & X_2 \\ 3 & 5 \\ 6.1 & 8.3 \\ ? & 7.4 \\ 4 & ? \end{bmatrix}$$

$$\theta = (\mu_1, \mu_2, \sigma^2, \tau^2, v)$$

Ex: One-way random effects

$$y_{ij} \stackrel{iid}{\sim} N(\mu_i, \sigma^2)$$

$$y_{ij} = \mu_i + \varepsilon_{ij} \sim N(0, \sigma^2)$$

$$\mu_i \stackrel{iid}{\sim} N(\alpha, \tau^2)$$

$$p(y_{i1}, y_{i2}, \dots, y_{in_i} | \alpha, \tau^2) \propto \int \prod_{j=1}^{n_i} p(y_{ij} | \mu_i, \sigma^2) p(\mu_i | \alpha, \tau^2) d\mu_i$$

$$p(y_i | \mu, \sigma, \alpha, \tau^2) \propto \prod_{j=1}^{n_i} p(y_{ij} | \mu_i, \sigma^2) p(\mu_i | \alpha, \tau^2)$$

$$\log p = \sum_{i=1}^G \sum_{j=1}^{n_i} \log p(y_{ij} | \mu_i, \sigma^2) + \log p(\mu_i | \alpha, \tau^2)$$

$$= \sum_{i=1}^G \left[ \sum_{j=1}^{n_i} \left( -\frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y_{ij} - \mu_i)^2 \right) - \frac{1}{2} \log \tau^2 - \frac{1}{2\tau^2} (\mu_i - \alpha)^2 \right]$$

$$\log p = -\frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_i \sum_j (y_{ij} - \mu_i)^2 - \frac{G}{2} \log \tau^2 - \frac{1}{2\tau^2} \sum_i (\mu_i - \alpha)^2$$

$$\frac{\partial}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{\sigma^4} \sum_i \sum_j (y_{ij} - \mu_i)^2$$

$$p(\mu_i | \alpha) \propto \left[ \prod_{j=1}^{n_i} p(y_{ij} | \mu_i, \sigma^2) \right] p(\mu_i | \alpha, \tau^2)$$

$$N\left(\alpha + \frac{\frac{1}{\sigma^2} \sum_j (y_{ij} - \alpha)}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}, \frac{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}{\frac{1}{\sigma^2} + \frac{1}{\tau^2}}\right)$$

$$\begin{aligned} y_{ij}^2 &= \mu_i^2 + 2\mu_i \varepsilon_{ij} + \varepsilon_{ij}^2 \\ \sum_j y_{ij}^2 &= \sum_j \mu_i^2 + 2\mu_i \sum_j \varepsilon_{ij} + \sum_j \varepsilon_{ij}^2 \\ &= n_i \mu_i^2 + 2\mu_i \sum_j \varepsilon_{ij} + \sum_j \varepsilon_{ij}^2 \end{aligned}$$

$$\begin{aligned} &-\frac{1}{2} \log \sigma^2 \\ &+ \frac{1}{2\sigma^2} \sum_j (y_{ij} - \mu_i)^2 \\ &= -\frac{1}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \sum_j (y_{ij}^2 - 2\mu_i y_{ij} + \mu_i^2) \end{aligned}$$

Ex. Mixture Models

$$Y_1 \rightarrow Y_n,$$

$$f(y_i) = \lambda \phi(y_i | \mu_1, \sigma_1^2) + (1-\lambda) \phi(y_i | \mu_2, \sigma_2^2)$$

$$\theta = (\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \lambda), \quad \lambda \in (0, 1)$$

$$\ell(\theta | Y) = \sum_{i=1}^n \log \left\{ \lambda \phi(y_i | \mu_1, \sigma_1^2) + (1-\lambda) \phi(y_i | \mu_2, \sigma_2^2) \right\}$$

Suppose  $Z \sim \text{Bernoulli}(\lambda)$ ,  $z_i \in \{0, 1\}$

$$z_i = 1, \text{ then } y_i \sim N(\mu_1, \sigma_1^2)$$

$$z_i = 0, \text{ then } y_i \sim N(\mu_2, \sigma_2^2).$$

CD L:

~~$$\ell(\theta | Y) = \prod_{i=1}^n \phi(y_i | \mu_1, \sigma_1^2)^{z_i} \phi(y_i | \mu_2, \sigma_2^2)^{1-z_i}$$~~

$$f(y_i | z_i) = \phi(y_i | \mu_1, \sigma_1^2)^{z_i} \phi(y_i | \mu_2, \sigma_2^2)^{1-z_i}$$

$$\beta(z_i) = \lambda^{z_i} (1-\lambda)^{1-z_i}$$

~~$$\ell(\theta) = \sum_{i=1}^n \log f(y_i | z_i) \beta(z_i)$$~~

$$\ell(\theta) = \sum_{i=1}^n \log f(y_i | z_i) \beta(z_i)$$

$$= \sum_{i=1}^n z_i \log \phi(y_i | \mu_1, \sigma_1^2) + (1-z_i) \log \phi(y_i | \mu_2, \sigma_2^2)$$

$$+ z_i \log \lambda + (1-z_i) \log (1-\lambda)$$

$$= \sum_{i=1}^n z_i \ell_1 + (1-z_i) \ell_2 + z_i \log \lambda + (1-z_i) \log (1-\lambda)$$

$$\frac{a}{x} =$$

$$\frac{b}{x} = 1 - \frac{a}{x}$$

$$\frac{b}{x} = \frac{x-a}{x}$$

$$x = a + b$$

$$p(z_i | y_i) \propto p(y_i | z_i) p(z_i)$$

$$= \lambda^{z_i} (1-\lambda)^{1-z_i} \lambda^{z_i} (1-\lambda)^{1-z_i}$$

$$= (\lambda e_1)^{z_i} (e_2 (1-\lambda))^{1-z_i}$$

$$= \text{Bernoulli} \left( \frac{\lambda e_1}{\lambda e_1 + (1-\lambda) e_2} \right)$$

$$\pi_i$$

$$E[z_i | y_i] = \pi_i = \frac{\lambda_0 e(y_i | \mu_{10}, \tau_{10}^2)}{\lambda_0 e(y_i | \mu_{10}, \tau_{10}^2) + (1-\lambda_0) e(y_i | \mu_{20}, \tau_{20}^2)}$$

$$Q(\theta | \theta_n) = \frac{1}{n} \sum_{i=1}^n E \left[ \sum_{i=1}^n z_i \log \lambda + (1-z_i) \log (1-\lambda) + c \right]$$

$$= \sum_{i=1}^n \pi_i \log \lambda + (1-\pi_i) \log (1-\lambda) + \pi_i \log \lambda + (1-\pi_i) \log (1-\lambda)$$

$$= \sum_{i=1}^n \pi_i \left[ -\frac{1}{2} \log 2\pi \sigma_1^2 - \frac{1}{2\sigma_1^2} (y_i - \mu_1)^2 \right] + (1-\pi_i) \left[ -\frac{1}{2} \log 2\pi \sigma_2^2 - \frac{1}{2\sigma_2^2} (y_i - \mu_2)^2 \right]$$

$$\Rightarrow \hat{\mu}_1 = \frac{\sum \pi_i y_i}{\sum \pi_i}$$

$$\hat{\sigma}_1^2 = \frac{\sum (y_i - \hat{\mu}_1)^2 \pi_i}{\sum \pi_i}$$

$$\hat{\mu}_2 = \frac{\sum (1-\pi_i) y_i}{\sum (1-\pi_i)}$$

$$\hat{\sigma}_2^2 = \frac{\sum (1-\pi_i) (y_i - \hat{\mu}_2)^2}{\sum (1-\pi_i)}$$

$$\hat{\lambda} = \frac{1}{n} \sum \pi_i$$

Complete  
data plf

For a regular exponential family,

$$g(x|\theta) = h(x) \exp(\theta^T t(x)) / a(\theta)$$

$$\log g(x|\theta) = \log h(x) + \theta^T t(x) - \log a(\theta)$$

$$Q(\theta|\tilde{\theta}) = \mathbb{E}[\log g(x|\theta)] = \theta^T \mathbb{E}[t(x)|\tilde{\theta}, \gamma] - \log a(\theta)$$

$$Q'(\theta|\tilde{\theta}) = \mathbb{E}[t(x)|\tilde{\theta}, \gamma] - \mathbb{E}_\theta[t(x)] = 0$$

$$\Rightarrow \mathbb{E}_\theta[t(x)] = \mathbb{E}[t(x)|\tilde{\theta}, \gamma]$$

Ex:  $y_1, y_2, \dots, y_m, z_{m+1}, \dots, z_n \sim N(\mu, \sigma^2)$

$$g(x|\mu, \sigma^2) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2\right)$$

$$\begin{aligned} \log g(x|\mu, \sigma^2) &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum (x_i^2 - 2x_i\mu + \mu^2) \\ &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \left[ \sum x_i^2 - 2\mu \sum x_i + n\mu^2 \right] \end{aligned}$$

$$\mathbb{E} \begin{pmatrix} \sum x_i^2 \\ \sum x_i \end{pmatrix} \bigg| \gamma = \begin{pmatrix} \sum_{i=1}^m y_i^2 + n(\tilde{\mu}^2 + \tilde{\sigma}^2) \\ \sum_{i=1}^m y_i + \tilde{\mu}(n-m+1) \end{pmatrix}$$

$$\mathbb{E} \begin{pmatrix} \sum x_i^2 \\ \sum x_i \end{pmatrix} = \begin{pmatrix} (\mu^2 + \sigma^2)n \\ \mu n \end{pmatrix}$$