

13: Modal And Distributional Approximations

Taylor

University of Virginia

We mention:

- ① a few ways to find the posterior mode
- ② how to approximate a posterior using a mode
- ③ slightly more involved ways to approximate your posterior

Newton's Method aka the Newton-Raphson algorithm

Based on a first-order approximation of the first derivative of the log-likelihood.

Approximate $L'(\theta) = (\log p(\theta | y))'$ as

$$\mathbf{0} \stackrel{\text{set}}{=} L'(\theta + \delta\theta) \approx L'(\theta) + L''(\theta)(\delta\theta)$$

rearranges to

$$\delta\theta = -[L''(\theta)]^{-1}L'(\theta)$$

Newton's Method

Repeat the following iteration until convergence:

$$\theta^t = \theta^{t-1} - [L''(\theta^{t-1})]^{-1}L'(\theta^{t-1})$$

Newton's Method aka the Newton-Raphson algorithm

Newton's Method

Repeat the following iteration until convergence:

$$\theta^t = \theta^{t-1} - [L''(\theta^{t-1})]^{-1} L'(\theta^{t-1})$$

Notes:

- 1 easily handles unnormalized densities
- 2 starting value is important because it is not guaranteed to converge from everywhere
- 3 The derivatives can be determined analytically or numerically

Quasi-Newton and conjugate gradient methods

Notes:

- 1 Quasi-Newton methods (approximate second derivatives) are available when second derivatives are too costly or unavailable
- 2 "Broyden-Fletcher-Goldfarb-Shanno" is a common example of a Quasi-Newton method
- 3 in R: `optim(2.9,F,method="BFGS")`
- 4 conjugate-gradient methods only use gradient information, but they are for models of the form $\|A\theta - b\|_2$ (also handled by `optim()`)
- 5 compared with the two above, they generally require more iterations, but use less computation per iteration and less storage

Numerical computation of derivatives

In `optim`, if you don't provide a function to calculate the gradient, then it uses a finite-difference approximation:

$$L'_i(\theta) = \frac{dL}{d\theta_i} \approx \frac{L(\theta + \delta_i e_i) - L(\theta - \delta_i e_i)}{2\delta_i}$$

and

$$\begin{aligned} L''_{ij}(\theta) &= \frac{d^2 L}{d\theta_i d\theta_j} \\ &\approx \frac{L'_i(\theta + \delta_j e_j) - L'_i(\theta - \delta_j e_j)}{2\delta_j} \end{aligned}$$

where e_j is the vector of all zeros except for a 1 in the j th spot, and δ_j is a small number (`optim`'s default is $1e-3$)

Gaussian approximations

Once the mode or modes have been found (perhaps after including a boundary-avoiding prior distribution as discussed in section 13.2, or after transforming the parameters appropriately), we can construct an approximation based on the multivariate normal distribution.

Let $\hat{\theta}$ be the mode, then

$$p(\theta | y) \approx N(\hat{\theta}, V_{\theta})$$

where

$$V_{\theta} = \left[- \frac{d^2 \log p(\theta | y)}{d\theta^2} \Big|_{\theta=\hat{\theta}} \right]^{-1}$$

is calculated exactly or approximated using the formula from a few slides ago.

Example

From chapter 3:

① $p(y_i | \mu, \sigma^2) = \text{Normal}(\mu, \sigma^2)$

② $p(\mu, \sigma^2) \propto 1/\sigma^2$

$$p(\mu, \sigma^2 | y) \propto (\sigma^2)^{-(n+2)/2} \exp \left[-\frac{1}{2\sigma^2} \{ (n-1)s^2 + n(\bar{y} - \mu)^2 \} \right]$$

Example

$$p(\mu, \sigma^2 \mid y) \propto (\sigma^2)^{-(n+2)/2} \exp \left[-\frac{1}{2\sigma^2} \{ (n-1)s^2 + n(\bar{y} - \mu)^2 \} \right]$$

Letting $\theta = \log \sigma$, $p(\mu, \theta \mid y)$ is proportional to

$$\exp[-n\theta] \exp \left[-\frac{1}{2 \exp[2\theta]} \{ (n-1)s^2 + n(\bar{y} - \mu)^2 \} \right]$$

So $\log p(\mu, \theta \mid y)$ is

$$\text{constant} - n\theta - .5 \exp(-2\theta) [(n-1)s^2 + n(\bar{y} - \mu)^2]$$

$$\text{and } L'(\theta) = \begin{bmatrix} \exp(-2\theta)(\bar{y} - \mu)n \\ -n + \exp(-2\theta) [(n-1)s^2 + n(\bar{y} - \mu)^2] \end{bmatrix}$$

Example

Warning: `optim` *minimizes*, so we use $-\log p(\mu, \theta \mid y)$

$$n\theta + .5 \exp(-2\theta) [(n-1)s^2 + n(\bar{y} - \mu)^2]$$

and

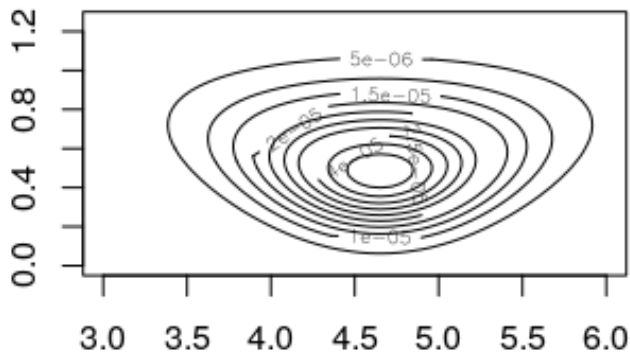
$$L'(\theta) = \begin{bmatrix} -\exp(-2\theta)(\bar{y} - \mu)n \\ n - \exp(-2\theta) [(n-1)s^2 + n(\bar{y} - \mu)^2] \end{bmatrix}$$

if `gr` left blank, finite difference approx. used

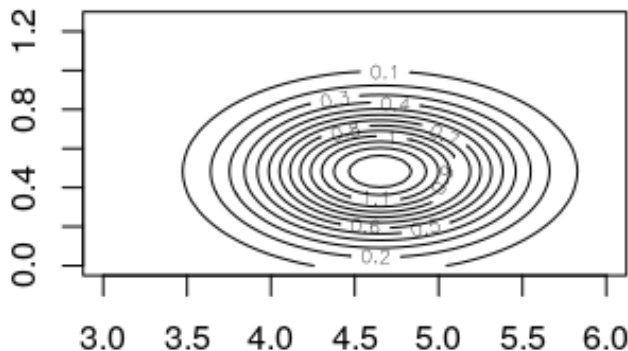
```
optim_results <- optim(par = c(5, 0),  
                        fn = neg_log_unnorm_post,  
                        gr = gradient,  
                        method = "BFGS",  
                        hessian = T)
```

See `mode_finding_examples.r`

Unnormalized true $p(\mu, \theta | y)$



Normal approx. $p(\mu, \theta \mid y)$



Gaussian approximations: Laplace's Method

If you want approximations to posterior *expectations* (say $E[h(\theta) | y]$), then you might consider Laplace's method, which is based on second-order Taylor approximations of the functions:

$$\textcircled{1} \quad u_1(\theta) = \log[h(\theta)q(\theta | y)]$$

$$\textcircled{2} \quad u_2(\theta) = \log q(\theta | y)$$

where $p(\theta | y) = q(\theta | y) / \int q(\theta | y) d\theta$.

Both are centered at maximizing values: θ_0^1, θ_0^2 , and this assumes h s are twice continuously differentiable.

Idea:

$$\frac{\int h(\theta)q(\theta | y)d\theta}{\int q(\theta | y)d\theta} = \frac{\int \exp[\log h(\theta) + \log q(\theta | y)] d\theta}{\int \exp[\log q(\theta | y)] d\theta}$$

Gaussian approximations: Laplace's Method

Exponentiating and integrating

$$\begin{aligned}u(\theta) &\approx u(\theta_0) + (\theta - \theta_0)^T u'(\theta_0) + \frac{1}{2}(\theta - \theta_0)^T u''(\theta_0)(\theta - \theta_0) \\&= u(\theta_0) + \frac{1}{2}(\theta - \theta_0)^T u''(\theta_0)(\theta - \theta_0)\end{aligned}$$

gives us

$$\begin{aligned}&\int \exp[u(\theta)] d\theta \\&\approx \int \exp[u(\theta_0) + (\theta - \theta_0)^T u'(\theta_0) + \frac{1}{2}(\theta - \theta_0)^T u''(\theta_0)(\theta - \theta_0)] d\theta \\&= \exp[u(\theta_0)] \int \exp\left[\frac{1}{2}(\theta - \theta_0)^T u''(\theta_0)(\theta - \theta_0)\right] d\theta \\&= \exp[u(\theta_0)] \int \exp\left[-\frac{1}{2}(\theta - \theta_0)^T \{-u''(\theta_0)\}(\theta - \theta_0)\right] d\theta\end{aligned}$$

The book has a few more generalizations that we don't address:

- ① approximating multimodal distributions with normal mixtures
- ② approximating multimodal distributions with t mixtures

The EM Algorithm

The **expectation-maximization algorithm** finds the argument that maximizes the marginal posterior. It's useful in situations where there is missing data in a model (e.g. factor models, hidden markov models, state space models, etc.).

It follows the following steps

- 1 replace missing values by their expectations given the guessed parameters,
- 2 estimate parameters assuming the missing data are equal to their estimated values,
- 3 re-estimate the missing values assuming the new parameter estimates are correct,
- 4 re-estimate parameters,

and so forth, iterating until convergence.

The EM Algorithm

Call $\theta = (\gamma, \phi)$. You're interested in the mode of $p(\phi | y)$.

$$\log p(\phi | y) = \log \frac{p(\gamma, \phi | y)}{p(\gamma | \phi, y)} = \log \underbrace{p(\gamma, \phi | y)}_{\text{joint posterior}} - \log \underbrace{p(\gamma | \phi, y)}_{\text{conditional posterior}}$$

The EM Algorithm

Call $\theta = (\gamma, \phi)$. You're interested in the mode of $p(\phi | y)$.

$$\log p(\phi | y) = \log \frac{p(\gamma, \phi | y)}{p(\gamma | \phi, y)} = \log \underbrace{p(\gamma, \phi | y)}_{\text{joint posterior}} - \log \underbrace{p(\gamma | \phi, y)}_{\text{conditional posterior}}$$

taking expectations on both sides with respect to $p(\gamma | \phi^{\text{old}}, y)$ yields:

$$\log p(\phi | y) = E \left[\log p(\gamma, \phi | y) | \phi^{\text{old}}, y \right] - E \left[\log p(\gamma | \phi, y) | \phi^{\text{old}}, y \right]$$

The EM Algorithm

We iteratively use the middle term in

$$\log p(\phi | y) = E [\log p(\gamma, \phi | y) | \phi^{\text{old}}, y] - E [\log p(\gamma | \phi, y) | \phi^{\text{old}}, y].$$

The Q quantity in the "E" step

$$Q(\phi | \phi^{\text{old}}) = E [\log p(\gamma, \phi | y) | \phi^{\text{old}}, y]$$

The EM Algorithm

We iteratively use the middle term in

$$\log p(\phi | y) = E [\log p(\gamma, \phi | y) | \phi^{\text{old}}, y] - E [\log p(\gamma | \phi, y) | \phi^{\text{old}}, y].$$

The Q quantity in the "E" step

$$Q(\phi | \phi^{\text{old}}) = E [\log p(\gamma, \phi | y) | \phi^{\text{old}}, y]$$

The EM algorithm

Repeat the following until convergence:

- 1 E-step: calculate $Q(\phi | \phi^{\text{old}})$
- 2 M-step: replace ϕ^{old} with $\arg \max Q(\phi | \phi^{\text{old}})$

The EM Algorithm

If we follow this strategy, $\log p(\phi \mid y)$ increases at every iteration:

$$\begin{aligned}\log p(\phi \mid y) &= E \left[\log p(\gamma, \phi \mid y) \mid \phi^{\text{old}}, y \right] - E \left[\log p(\gamma \mid \phi, y) \mid \phi^{\text{old}}, y \right] \\ &= Q(\phi \mid \phi^{\text{old}}) - E \left[\log p(\gamma \mid \phi, y) \mid \phi^{\text{old}}, y \right] \quad (\text{defn. } Q) \\ &\geq Q(\phi \mid \phi^{\text{old}}) - E \left[\log p(\gamma \mid \phi^{\text{old}}, y) \mid \phi^{\text{old}}, y \right] \quad (\text{HW})\end{aligned}$$

The EM Algorithm

If we follow this strategy, $\log p(\phi | y)$ increases at every iteration:

$$\begin{aligned}\log p(\phi | y) &= E \left[\log p(\gamma, \phi | y) | \phi^{\text{old}}, y \right] - E \left[\log p(\gamma | \phi, y) | \phi^{\text{old}}, y \right] \\ &= Q(\phi | \phi^{\text{old}}) - E \left[\log p(\gamma | \phi, y) | \phi^{\text{old}}, y \right] \quad (\text{defn. } Q) \\ &\geq Q(\phi | \phi^{\text{old}}) - E \left[\log p(\gamma | \phi^{\text{old}}, y) | \phi^{\text{old}}, y \right] \quad (\text{HW})\end{aligned}$$

So

$$\begin{aligned}\log p(\phi^{\text{new}} | y) - \log p(\phi^{\text{old}} | y) &= \log p(\phi^{\text{new}} | y) - \left\{ Q(\phi^{\text{old}} | \phi^{\text{old}}) - E \left[\log p(\gamma | \phi^{\text{old}}, y) | \phi^{\text{old}}, y \right] \right\} \\ &\geq Q(\phi | \phi^{\text{old}}) - E \left[\log p(\gamma | \phi^{\text{old}}, y) | \phi^{\text{old}}, y \right] \\ &\quad - \left\{ Q(\phi^{\text{old}} | \phi^{\text{old}}) - E \left[\log p(\gamma | \phi^{\text{old}}, y) | \phi^{\text{old}}, y \right] \right\} \\ &= Q(\phi | \phi^{\text{old}}) - Q(\phi^{\text{old}} | \phi^{\text{old}})\end{aligned}$$

The EM Algorithm

Notes:

- 1 The EM algo isn't inherently Bayesian. It can also be used to accomplish maximum likelihood estimation.
- 2 The expectation of $\log p(\gamma, \phi \mid y)$ is usually easy to compute because it is a sum, and might only depend on sufficient statistics
- 3 The EM algorithm implicitly deals with parameter constraints in the M-step
- 4 The EM algorithm is parameterization independent
- 5 The *Generalized* EM (GEM) just increases Q instead of maximizing it.
- 6 The book describes many generalizations in addition to this one
- 7 You can find multiple modes if you start from multiple starting points (using mixture approximations afterwards)
- 8 Debug by printing $\log p(\phi^i \mid y)$ at every iteration and make sure it increases monotonically

Variational inference approximates an intractable posterior $p(\theta \mid y)$ with some chosen distribution $g(\theta \mid \phi)$ (e.g. multivariate normal).

Variational inference approximates an intractable posterior $p(\theta \mid y)$ with some chosen distribution $g(\theta \mid \phi)$ (e.g. multivariate normal).

We will assume all J parameters are independent a posteriori. In other words

$$g(\theta \mid \phi) = \prod_{j=1}^J g_j(\theta_j \mid \phi_j) = g_j(\theta_j \mid \phi_j) g_{-j}(\theta_{-j} \mid \phi_{-j}).$$

We will find ϕ using an EM-like algorithm that minimizes Kullback-Leibler divergence.

Kullback-Leibler divergence is reversed this time:

$$\begin{aligned} KL(g||p) &= - \int \log \left(\frac{p(\theta | y)}{g(\theta | \phi)} \right) g(\theta | \phi) d\theta \\ &= - \int \log \left(\frac{p(\theta, y)}{g(\theta | \phi)} \right) g(\theta | \phi) d\theta + \int \log p(y) g(\theta | \phi) d\theta \\ &= - \underbrace{\int \log \left(\frac{p(\theta, y)}{g(\theta | \phi)} \right) g(\theta | \phi) d\theta}_{\text{variational lower bound}} + \log p(y) \end{aligned}$$

The term that we maximize (minimize the negative) is called the **variational lower bound** aka the **evidence lower bound** (ELBO).

Variational Inference

Every iteration, we cycle through all the hyper-parameters ϕ_1, \dots, ϕ_J , and change them until convergence is reached.

$$\begin{aligned} & \int \log \left(\frac{p(\theta, y)}{g(\theta | \phi)} \right) g(\theta | \phi) d\theta \\ &= \iint [\log p(\theta, y) - \log g_j(\theta_j | \phi_j) - \log g_{-j}(\theta_{-j} | \phi_{-j})] \\ & \quad g_j(\theta_j | \phi_j) g_{-j}(\theta_{-j} | \phi_{-j}) d\theta_j d\theta_{-j} \\ &= \int \left[\int \log p(\theta, y) g_{-j}(\theta_{-j} | \phi_{-j}) d\theta_{-j} \right] g_j(\theta_j | \phi_j) d\theta_j \\ & \quad - \int \log g_j(\theta_j | \phi_j) g_j(\theta_j | \phi_j) d\theta_j - \int \log g_{-j}(\theta_{-j} | \phi_{-j}) g_{-j}(\theta_{-j} | \phi_{-j}) d\theta_{-j} \\ &= \int \log \left(\frac{\tilde{p}(\theta_j)}{g_j(\theta_j | \phi_j)} \right) g_j(\theta_j | \phi_j) d\theta_j + \text{constant} \end{aligned} \quad (*)$$

Variational Inference

We think of $\tilde{p}(\theta_j)$ as an unnormalized density

$$\log \tilde{p}(\theta_j) = \int \log p(\theta, y) g_{-j}(\theta_{-j} \mid \phi_{-j}) d\theta_{-j}$$

if

$$\begin{aligned} \int \tilde{p}(\theta_j) d\theta_j &= \int \exp \left[\int \log p(\theta, y) g_{-j}(\theta_{-j} \mid \phi_{-j}) d\theta_{-j} \right] d\theta_j \\ &\leq \int \exp \left[\log \int p(\theta, y) g_{-j}(\theta_{-j} \mid \phi_{-j}) d\theta_{-j} \right] d\theta_j \quad (\text{Jensen's}) \\ &= \iint p(\theta, y) g_{-j}(\theta_{-j} \mid \phi_{-j}) d\theta_{-j} d\theta_j \\ &< \infty \end{aligned}$$

VI algorithm

For $j = 1, \dots, J$:

Set ϕ_j so that $\log g_j(\theta_j | \phi_j)$ is equal to

$$\log \tilde{p}(\theta_j) = \int \log p(\theta, y) g_{-j}(\theta_{-j} | \phi_{-j}) d\theta_{-j}$$