

Computational Aspects

of MLE:

If the likelihood equations are non-linear functions of the parameters, then usually cannot get a "closed-form" solution to these first-order conditions. Use a numerical approximation to the solution.

Lots of different techniques - illustrate with "Methods of Descent". (Max. $\log L$ equivalent to min. $-\log L$)

Exactly same problem encountered with non-linear least squares.

General approach : $\hat{\theta} = \theta_0 + s \cdot d(\theta_0)$

θ_0 = initial value

s = step-length (scalar, > 0)

$d(\cdot)$ = direction vector

1

Usually $d(\cdot)$ depends on gradient

vector (or maybe Hessian matrix), at θ_0 .

Sometimes $s = s$ (Hessian), too.

A specific member of family of descent methods is Newton-Raphson algorithm. First, let's look at a simple geometric example (scalar case). Let $g(\theta) = -\frac{\partial \log L(\theta)}{\partial \theta}$. We want to find $\hat{\theta}$ s.t. $g(\hat{\theta}) = 0$. Of course, also need to consider 2nd-order condition.

N.B.: Max. $\log L \Leftrightarrow$ min. $-\log L = f$. Start at an arbitrary point, θ_0 , move to θ_1 , etc. using algorithm, until convergence.

2



Clearly, this is just one possible approach where does this come from, mathematically? (knowing this will help us see when it may work well, or badly.)

Let $f(\underline{\theta}) = -\log L(\underline{\theta})$. Objective

is to min. $f(\underline{\theta})$. Let's approximate $f(\cdot)$ using a 2nd-order Taylor's series expansion

$$f(\underline{\theta}) \approx f(\underline{\tilde{\theta}}) + (\underline{\theta} - \underline{\tilde{\theta}})' \left(\frac{\partial f(\underline{\theta})}{\partial \underline{\theta}} \right) \Big|_{\underline{\theta}=\underline{\tilde{\theta}}}$$

How do we get from $\underline{\theta}_0$ to $\underline{\theta}_1$?

$$\text{Slope of tangent} = \left(\frac{g(\underline{\theta}_0)}{\underline{\theta}_0 - \underline{\theta}_1} \right) \\ = H(\underline{\theta}_0) \quad [= \frac{\partial g(\underline{\theta})}{\partial \underline{\theta}}]$$

(Approximation is valid in neighbourhood

of $\underline{\tilde{\theta}}$, the value that minimizes $f(\cdot)$)

More generally:

$$\underline{\theta}_{n+1} = \underline{\theta}_n - H^{-1}(\underline{\theta}_n) g(\underline{\theta}_n)$$

[In this form, also o.k. for vector, $\underline{\theta}$.]

$$f(\underline{\theta}) \approx f(\underline{\tilde{\theta}}) + (\underline{\theta} - \underline{\tilde{\theta}})' g(\underline{\tilde{\theta}}) \\ + \frac{1}{2} (\underline{\theta} - \underline{\tilde{\theta}})' H(\underline{\tilde{\theta}}) (\underline{\theta} - \underline{\tilde{\theta}})$$

Re-writing this:

Differentiate:

$$\begin{aligned} \left[\frac{\partial f(\underline{\theta})}{\partial \underline{\theta}} \right] &\approx 0 + g(\tilde{\underline{\theta}}) \\ &+ (\frac{1}{2})(2) H(\tilde{\underline{\theta}})(\underline{\theta} - \tilde{\underline{\theta}}) \\ &= H(\tilde{\underline{\theta}})(\underline{\theta} - \tilde{\underline{\theta}}) \end{aligned}$$

(because $g(\tilde{\underline{\theta}}) = 0$.)

Re-arranging :

$$\begin{aligned} \tilde{\underline{\theta}} &= \underline{\theta} - H^{-1}(\tilde{\underline{\theta}}) \left[\frac{\partial f(\underline{\theta})}{\partial \underline{\theta}} \right] \\ &= \underline{\theta} - H^{-1}(\tilde{\underline{\theta}}) g(\underline{\theta}) \end{aligned}$$

So, begin with $\underline{\theta} = \underline{\theta}_0$ & iterate:

$$\begin{aligned} \underline{\theta}_1 &= \underline{\theta}_0 - H^{-1}(\underline{\theta}_1) g(\underline{\theta}_0) \\ \underline{\theta}_2 &= \underline{\theta}_1 - H^{-1}(\underline{\theta}_2) g(\underline{\theta}_1) \\ &\vdots \\ &\vdots \\ \underline{\theta}_{n+1} &= \underline{\theta}_n - H^{-1}(\underline{\theta}_{n+1}) g(\underline{\theta}_n) \end{aligned}$$

or, approximately —

$$\underline{\theta}_{n+1} = \underline{\theta}_n - H^{-1}(\underline{\theta}_n) g(\underline{\theta}_n)$$

(i) Stop if $|\frac{\underline{\theta}_{n+1} - \underline{\theta}_n}{\underline{\theta}_n}| < \epsilon_i$

for $i = 1, 2, \dots, k$.

(ii) Algorithm corresponds to case of $s=1$.

(iii) $d(\underline{\theta}_n) = -H^{-1}(\underline{\theta}_n) g(\underline{\theta}_n)$

(iv) Algorithm will fail if $H(\cdot)$ singular.

(v) If $H(\cdot)$ becomes n.d., we'll

locate a max. of $f(\cdot)$ hence a min. If $\log \cdot$. So, need to ensure that $H(\tilde{\underline{\theta}})$ is p.d.

(vi) Algorithm may locate a local min.

(depending on initial value), or may oscillate & never converge.

Algorithm should work well if

$\log(\underline{\theta})$ is approximately quadratic at least in neighbourhood of max.

7

Example:

$$f(\theta) = \theta^3 - 2\theta^2 - 1$$

is immediate. (In this case, $g(\theta)$ is linear.)

$$\begin{aligned} f'(\theta) &= 3\theta^2 - 4\theta = \theta(3\theta - 4) \\ &= 0 \end{aligned}$$

$$\Rightarrow \theta = 0, \frac{4}{3}.$$

$$\begin{cases} f''(0) = -4 & (\text{max.}) \\ f''\left(\frac{4}{3}\right) = 4 & (\text{min.}) \end{cases}$$

Now, let's see how the algorithm goes:

Start at $\theta_0 = 1$:

$$g(\theta_0) = f'(\theta_0) = -4; H(\theta_0) = f''(\theta_0) = 2$$

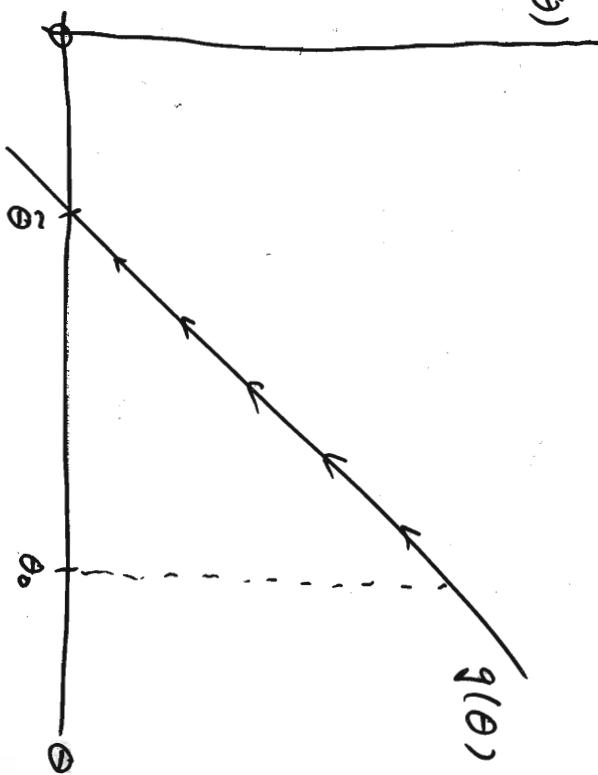
$$\begin{aligned} \text{So, } \theta_1 &= \theta_0 - H(\theta_0)^{-1} g(\theta_0) = 1 - (-\frac{1}{2}) = 1.5 \\ g(\theta_1) &= \frac{3}{4}; H(\theta_1) = 5 \end{aligned}$$

Need to experiment with different

initial values to ensure we locate a global max. of $\log L$.

Rapidly get to $\tilde{\theta} = 1.333$.

What happens if we choose $\theta_0 = -1$?



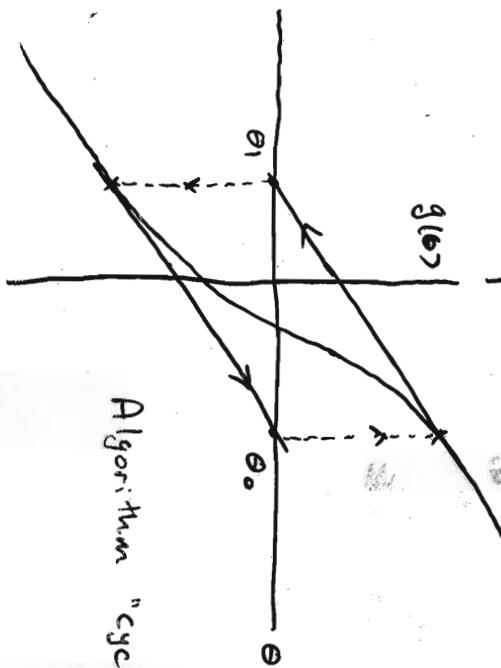
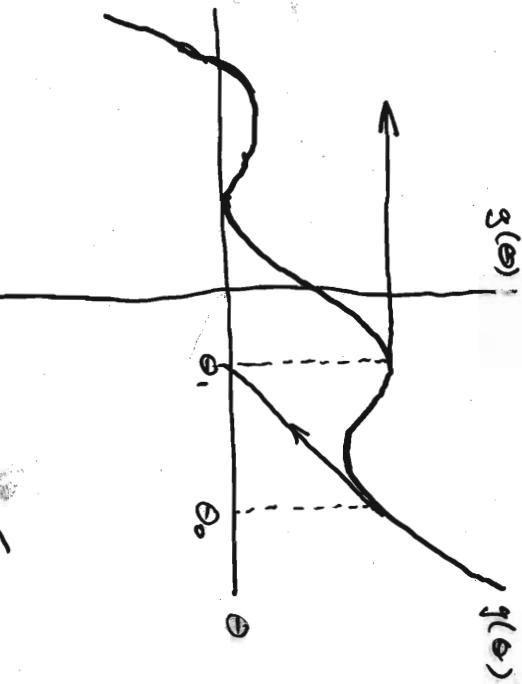
8

Things that can go wrong :

9

"Concentrating" the Likelihood Fctn.

$H(\cdot)$ singular
at $\underline{\theta}_1$.



Algorithm "cycles"

Sometimes the parameter vector falls into 2 natural parts : $\underline{\theta} = (\underline{\theta}_1, \underline{\theta}_2)$.
So, $L(\underline{\theta} | y) = L(\underline{\theta}_1, \underline{\theta}_2 | y)$. In addition, often the MLE for $\underline{\theta}_2$ can be written as a func. of the MLE for $\underline{\theta}_1$:
 $\hat{\underline{\theta}}_2 = f(\hat{\underline{\theta}}_1)$. In this case,

$$L(\underline{\theta}_1, \underline{\theta}_2 | y) = L[\underline{\theta}_1, f(\underline{\theta}_1) | y]$$

$$= L_c[\underline{\theta}_1 | y]$$

where $L_c(\cdot)$ is the concentrated likelihood fctn. The point is that we can then maximize $L_c(\cdot)$ with respect to just $\underline{\theta}_1$. If we can get $\hat{\underline{\theta}}_2 = f(\hat{\underline{\theta}}_1)$ analytically this means that the numerical maximization has to be done only w.r.t. $\underline{\theta}_1$, not $\underline{\theta}$, which simplifies matters significantly.

10.

N.B.: "Concentrating" the L.F. is not always an option - worth keeping in mind, though.

Example: (concentrating is unnecessary, but it still works)

$$y = X\beta + \varepsilon; \quad \varepsilon \sim N[0, \sigma^2 I]$$

$$p(y | \beta, \sigma) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left[-\frac{1}{2\sigma^2}(y - X\beta)'(y - X\beta)\right]$$

$$\log L(\beta, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2} \log 2\pi - \frac{1}{2\sigma^2} (y - X\beta)'(y - X\beta)$$

$$(\partial \log L / \partial \sigma^2) = \frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} (y - X\beta)'(y - X\beta) = 0$$

$$\log L_c(\beta | y) = -\frac{n}{2} \log \left[\frac{1}{n} (y - X\beta)'(y - X\beta) \right]$$

$$-\frac{n}{2} \log 2\pi - \frac{n}{2}$$

$$(\partial \log L_c / \partial \beta)_i = -\frac{n}{2} \frac{\partial}{\partial \beta_i} \frac{(y - X\beta)'(y - X\beta)}{(y - X\beta)'(y - X\beta)} \cdot \left[\frac{\partial}{\partial \beta_i} (y - X\beta)'(y - X\beta) \right]$$

$$= 0$$

$$\Rightarrow \frac{\partial}{\partial \beta} [(y - X\beta)'(y - X\beta)] = 0$$

$$\Rightarrow \tilde{\beta} = (X'X)^{-1} X'y$$

$$\Rightarrow \tilde{\sigma}^2 = \frac{1}{n} (y - X\tilde{\beta})'(y - X\tilde{\beta})$$

So, "concentrating" the L.F. certainly works. Imagine that we knew how to differentiate w.r.t. a scalar, but not a vector. Then, we could have obtained $L_c(\beta | y)$ above & then would have had to max. $\{L_c(\beta | y)\}$ w.r.t. β .

That is, maximize

$$L_c(\beta | y) = -\frac{n}{2} \log \left[\frac{1}{n} (y - X\beta)'(y - X\beta) \right] - \frac{n}{2} \log (2\pi) - \frac{n}{2}$$

$$\text{or, max. } \left\{ -\frac{n}{2} \log (y - X\beta)'(y - X\beta) \right\} \text{ w.r.t. } \beta.$$

We'll encounter situations where we can maximize partly analytically, but need to use numerical algorithm for the rest of problem. "Concentrating" will reduce dimension of space for numerical optimization.

1.3.

Inequality Constraints:

Sometimes we'll want to impose an inequality constraint on one or more of our MLE's. For example: $\tilde{\theta}_1 > 0$; $0 < \tilde{\theta}_2 < 1$. Most packages don't allow this. However, there are some tricks we can use to impose such constraints.

Example: We want $\tilde{\theta}_1 \geq 0$. So, replace θ_1 with $(\phi_1)^2$, which will be ≥ 0 .

No constraint is placed on ϕ_1 ; $\tilde{\theta}_1 = \tilde{\phi}_1^2$.

Example: We want $0 \leq \theta_2 \leq 1$. Replace θ_2 by $1/(1 + \phi_2)$, or replace θ_2 by $\exp(\phi)/(\lambda + \exp(\phi))$ & then $\tilde{\theta}_2 = 1/(1 + \tilde{\phi}^2)$, or $1/(1 + \exp(\tilde{\phi}))$ will satisfy the constraint.