

Nonlinear Least-squares problems

(10)

So far $F(x_*) = 0$, $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$

Will consider $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$

Say $m > n \Rightarrow$ overdetermined } regardless of existence
 $m < n \Rightarrow$ underdetermined } uniqueness of solution

Tonight, overdetermined case \sim nonlinear LS problem

Assume $\|\cdot\| = \|\cdot\|_2$ throughout

Linear LS problems:

- suppose τ depends on $\alpha_1, \dots, \alpha_n$
- postulate linear model $\tau = x_1 \alpha_1 + \dots + x_n \alpha_n$
- want to determine x_i 's from observations

$$i=1, \dots, m \quad \begin{cases} b_i = i^{\text{th}} \text{ observed value of } \tau \\ a_{ij} = i^{\text{th}} \text{ observed value of } \alpha_j \quad 1 \leq j \leq n \end{cases}$$

Method of least squares: Determine x_i 's to

$$\text{minimize } \sum_{i=1}^m (b_i - \sum_{j=1}^n a_{ij} x_j)^2 = \|B - Ax\|^2$$

$$\text{where } B = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix} \quad A = (a_{ij}) \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$B - Ax \sim$ residual (vector)

$\|B - Ax\|^2 \sim$ sum of squared residuals

solution is characterized by $A^T A x = A^T B \dots (*)$

normal equation
of least-squares

$\Leftrightarrow A^T (B - Ax) = 0$ i.e. residuals are orthogonal (normal)
to range(A)

(*) always has a solution - unique $\Leftrightarrow \text{rank}(A) = n$

Also $A^T A$ is symmetric and positive semidefinite

(SPD) $\Leftrightarrow \text{rank}(A) = n$

- If A is full rank, can form $A^T A$ $A^T b$ and solve using Cholesky. But there are alternatives based on QR decomposition or SVD of A

- Extend to nonlinear LS problem.

Obvious extension: nonlinear model

$$b = f(a, x) \quad \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} \quad \text{independent variables}$$

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \sim \text{parameters}$$

Then given observations on b and a choose x to

$$\text{minimize } \|b - g(a, x)\|^2$$

where

$$b = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}$$

$b_i = i^{\text{th}}$ observation on b

$$g = \begin{pmatrix} f(a_1, x) \\ \vdots \\ f(a_m, x) \end{pmatrix}$$

$a_i = i^{\text{th}}$ observation of a

Example: $b = x_1 e^{x_2 x}$

Consider even more general problem

$$\text{minimize } f(x) = \frac{1}{2} \|F(x)\|^2 \quad F: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Natural: apply Newton-based methods to find x^* such that $\nabla f(x^*) = 0$

This problem has special structure

Can show: $\nabla f = J^T F$, $J = F'$ (are evaluated at "current" x)

$$\nabla^2 f = J^T J + \sum_{i=1}^m F_i \nabla^2 F_i, \text{ where } F_i = i^{\text{th}} \text{ component of } F$$

See: $J^T J$ is symmetric positive semidefinite always, and positive definite $\Leftrightarrow \text{rank}(J) = n$

Also once J is evaluated to form ∇f , have $J^T J$ with some additional arithmetic (and no additional function evaluations)

$\sum_{i=1}^m F_i \nabla^2 F_i$ requires additional evaluations of m $n \times n$ Hessians (usually infeasible)

Also $\sum F_i \nabla^2 F_i$ may be such that $\nabla^2 f$ is indefinite \Rightarrow Newton step may not be a descent direction

Gauss-Newton method

given an initial x

until termination do

$$\text{solve } J^T J s = -J^T F \quad (**)$$

update $x \leftarrow x + s$

Remarks:

(1) $(**)$ ~ normal equation for linear LS problem
 $\min \|F - J s\|^2$

usually solve using QR or SVD, sometimes form $(**)$ and solve using Cholesky

(2) Gauss-Newton step $S = -(J^T J)^{-1} J^T F$ is a descent direction for f if $\text{rank}(J) = n$

$$\begin{aligned} \nabla f^T S &= (J^T F)^T \left(-(J^T J)^{-1} J^T F \right) \\ &= - (J^T F)^T (J^T J)^{-1} (J^T F) < 0 \end{aligned}$$

for $J^T F \neq 0$,

since

$\text{rank}(J) = n \Rightarrow J^T J$ is SPD $\Rightarrow (J^T J)^{-1}$ is SPD

(3) termination? As before

stop if iteration number $> \text{it}_{\max}$

stop if $\|S\| \leq \epsilon_x$

can't expect $\|F(x_*)\| = 0$, so stop if

$$\|\nabla f(x)\| = \|J(x)^T F(x)\| \leq \epsilon_f$$

Convergence

Recall: Given $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$, consider

General iteration

Given x_0 and B_0

for $k=0, 1, \dots$

$$x_{k+1} = x_k - B_k^{-1} F(x_k)$$

determine B_{k+1}

Theorem: Suppose F is continuously differentiable near x_* such that $F(x_*) = 0$ and $F'(x_*)$ nonsingular.

Suppose that the general iteration produces $\{x_k\}$ with $\{B_k\}$ such that for all k

$$1) \|B_k^{-1}\| \leq M$$

$$2) \|I - B_k^{-1} F'(x_k)\| \leq \epsilon_{\max} < 1$$

Then, given any ϵ , $0 < \epsilon_{\max} < 1$ there is a $\delta > 0$ such that $\|x_0 - x_*\| \leq \delta$, then $x_k \rightarrow x_*$ with $\|x_{k+1} - x_*\| \leq \epsilon \|x_k - x_*\|$ for all k

Apply to Gauss-Newton:

$$F \leftarrow \nabla f = J^T F \quad (\text{not same } F\text{'s})$$

$$B_k \leftarrow J(x_k)^T J(x_k)$$

Assume $f = \frac{1}{2} \|F\|^2$ is twice continuously differentiable near x_* such that $\nabla f(x_*) = J^T(x_*) F(x_*) = 0$ and $J(x_*)$ is full-rank

see: $J(x_*)^T J(x_*)$ is SPD $\Rightarrow J(x)^T J(x)$ and $(J(x)^T J(x))^{-1}$ exist and are continuous and are SPD and

$$\|(J(x)^T J(x))^{-1}\| \leq M \text{ for } x \text{ near } x_*$$

So (1) holds.

$$\text{Assume: } \|I - (J(x_*)^T J(x_*))^{-1} \nabla^2 f(x)\| \leq \epsilon_* < 1$$

then for any ϵ_{\max} in $(\epsilon_*, 1)$ have

$$\|I - (J(x)^T J(x))^{-1} \nabla^2 f(x)\| \leq \epsilon_{\max} < 1 \quad \dots (***)$$

for x sufficiently near x_* \Rightarrow (2) also holds for x_k near x_*

Then Theorem \Rightarrow For any ϵ , $\epsilon_{\max} < \epsilon < 1$, $\exists \delta > 0$ such that

$$\|x_0 - x_*\| \leq \delta \Rightarrow x_k - x_*$$

$$\text{with } \|x_{k+1} - x_*\| \leq \epsilon \|x_k - x_*\| \text{ for all } k$$

so can assume $x_k \rightarrow x_*$

For x near x_* get

$$x_+ = x - (J^T J)^{-1} J^T F$$

set $e = x - x_*$

$$e_+ = x_+ - x_*$$

Have $e_+ = e - (J^T J)^{-1} J^T F$

Also $F(x) = F(x_*) + \int_0^1 \frac{d}{dt} F(x_* + te) dt$

$$= \left\{ \int_0^1 J(x_* + te) dt \right\} e \pm J(x_*) e$$

$$= J(x_*) e + \underbrace{\left\{ \int_0^1 [J(x_* + te) - J(x_*)] dt \right\}}_E e$$

$$= J(x_*) e + E(x) e$$

Note: $\|E(x)\| \leq \int_0^1 \|J(x_* + te) - J(x_*)\| dt$
 \wedge
 $t \leq \|e\|$

$$= \frac{\gamma}{2} \|e\|$$

So

$$e_+ = e - (J^T J)^{-1} J^T [J e + (J_* - J) e + E e]$$

$$= e - e - \underbrace{(J^T J)^{-1} J^T}_{\text{bounded near } x_*} \left[\underbrace{(J_* - J) e}_{\| \cdot \| \leq \gamma \| e \|} + \underbrace{E e}_{\| \cdot \| \leq \frac{\gamma}{2} \| e \|} \right]$$

bounded near x_*
(say $\leq M$)

$$\|e_+\| \leq M \frac{3\gamma}{2} \|e\|^2 \Rightarrow x_k \rightarrow x_* \text{ quadratically}$$

What about globalization?

Criterion for accepting steps is as before:

$$\text{ared} \geq t \text{pred}$$

where

$$\text{ared} = f(x) - f(x+s)$$

$$\text{pred} = f(x) - g(s)$$

where

$$g(s) = f(x) + \nabla f^T s + \frac{1}{2} s^T \underbrace{J^T J}_{\text{not } \nabla^2 f} s$$

$$\text{pred} = f(x) - g(s)$$

$$= -\nabla^2 f s - \frac{1}{2} s^T J^T J s$$

note:

$$g(s) = \frac{1}{2} \|F\|^2 + F^T J s + \frac{1}{2} \|J s\|^2$$

$$= \frac{1}{2} \|F + J s\|^2$$

$$\Rightarrow \text{pred} = \frac{1}{2} \left(\|F\|^2 - \|F + J s\|^2 \right)$$

and have $\text{ared} \geq t \text{pred}$

$$\|F\|^2 - \|F(x+s)\|^2 \geq t \left(\|F\|^2 - \|F + J s\|^2 \right)$$

(almost like $\|F\| - \|F(x+s)\| \geq t (\|F\| - \|F + J s\|)$ when $w=u$)

How to modify steps if necessary?

Know GN step is a descent direction for f , so can implement backtracking as before.

Can also implement trust region methods as before

$$\text{Have } g(s) = \frac{1}{2} \|F + J s\|^2$$

so given TR radius $\delta > 0$

$$\text{TR step} = \underset{\|s\| \leq \delta}{\text{argmin}} g(s) = \underset{\|s\| \leq \delta}{\text{argmin}} \|F + J s\|$$

and as before have

$$\text{TR step} = s(\mu) = - \left[J^T J + \mu I \right]^{-1} J^T F \quad \dots (***)$$

for a unique $\mu > 0$, as follows

$$\|s^{GH}\| \leq \delta \Rightarrow \text{TR step} = s^{GH} \leftarrow - (J^T J)^{-1} J^T F$$

$$\|s^{GH}\| > \delta \Rightarrow \text{TR step} = s(\mu) \text{ for the unique } \mu > 0 \text{ such that } \|s(\mu)\| = \delta$$

Remark: Methods for nonlinear LS with step (***) are called Levenberg-Marquardt methods and date back to Levenberg (1944) and Marquardt (1963). In the old days LM algorithms controlled μ . Now considered much better to control δ .

Remark: Can find TR step (or dogleg approximation) as before. Note that model Hessian is $J^T J$ not $\nabla^2 f$ so can't see & exploit negative curvature.

What if (e.g. in "large residual" or "highly nonlinear case")

$$\sum F_i \nabla^2 F_i \quad \text{is not small relative to } J^T J$$

Evaluating $\sum F_i \nabla^2 F_i$ to get $\nabla^2 f$ is very likely infeasible

Idea: Use $J^T J + A$, where A is maintained by quasi-Newton (secant) updating.

Usual situation: Have B and step s (possibly $-B^{-1} \nabla f$) at "current" x

set $x_+ = x + s$ and B_+ "least change" update satisfying $B_+ s = y \equiv f(x_+) - f(x)$

How have

$$B = J^T J + A$$

and want

$$B_+ = J_+^T J_+ + A_+$$

where

A_+ = "least change" update satisfying an appropriate secant condition (?)

straight forward condition: Want

$$B_+ s = f(x_+) - f(x) \equiv y$$

so require

$$A_+ s = y - J_+^T J_+ s$$

Dennis-Schwarz suggest replacing y by $y^\#$

$$y^\# = J(x_+)^T f(x_+) - J(x)^T f(x)$$

Also suggest symmetry-preserving update analogous to Davidon-Fletcher-Powell update

$$A_+ = A + \frac{(y^\# - As)y^T + y(y^\# - As)^T}{y^T s} - \frac{(y^\# - As)^T s y y^T}{(y^T s)^2}$$

DGW $y^\#$ is "slight but clear winner" along several choices that all give local superlinear convergence (D-S page 232)

More recent updates analogous to Broyden-Fletcher-Goldfarb-Shanno update are now considered somewhat better (AC Broyden-Fletcher)