

Trust region methods

Have trust region

$$N_{\delta}(x) = \{y : \|y - x\| \leq \delta\}$$

choose

$$\begin{aligned} s &= \text{argmin}_{\|w\| \leq \delta} \|F + Jw\|, & F &= F(x), \\ && J &= F'(x) \end{aligned}$$

- x is a stationary point for F if $\|F + Js\| \geq \|F\|$
for all s

- Saw: $\|s^*\| \leq \delta \Rightarrow s = s^*$
 $\|s^*\| > \delta \Rightarrow \|s\| = \delta$
etc.

- If J nonsingular, then $s = s(\mu) = -[J^T J + \mu I]^{-1} J^T F$
for unique $\mu > 0$

$$\|s^*\| \leq \delta \Rightarrow \mu = 0$$

$$\|s^*\| > \delta \Rightarrow \mu > 0 \text{ and such that } \|s(\mu)\| = \delta$$

Also, $s(\mu)$ and $\varphi(\mu) = \|s(\mu)\|^2$ are differentiable and
 $\varphi'(\mu) < 0 \Rightarrow \|s(\mu)\|$ monotone decreasing

$$\text{See: } \mu = 0 \Rightarrow s(\mu) = s^*$$

for μ range, $s(\mu) \approx -\frac{1}{\mu} J^T F$ - short step in
steepest descent direction



$$f(x) = \frac{1}{2} \|F(x)\|^2 = \frac{1}{2} F^T F$$

$$\nabla f = J^T F$$

$$\begin{aligned} f(x+s) &= \frac{1}{2} F(x+s)^T F(x+s) \\ &= \frac{1}{2} \|F\|^2 + \underbrace{F^T J s}_{(J^T F)^T s} + \Theta(s) \end{aligned}$$

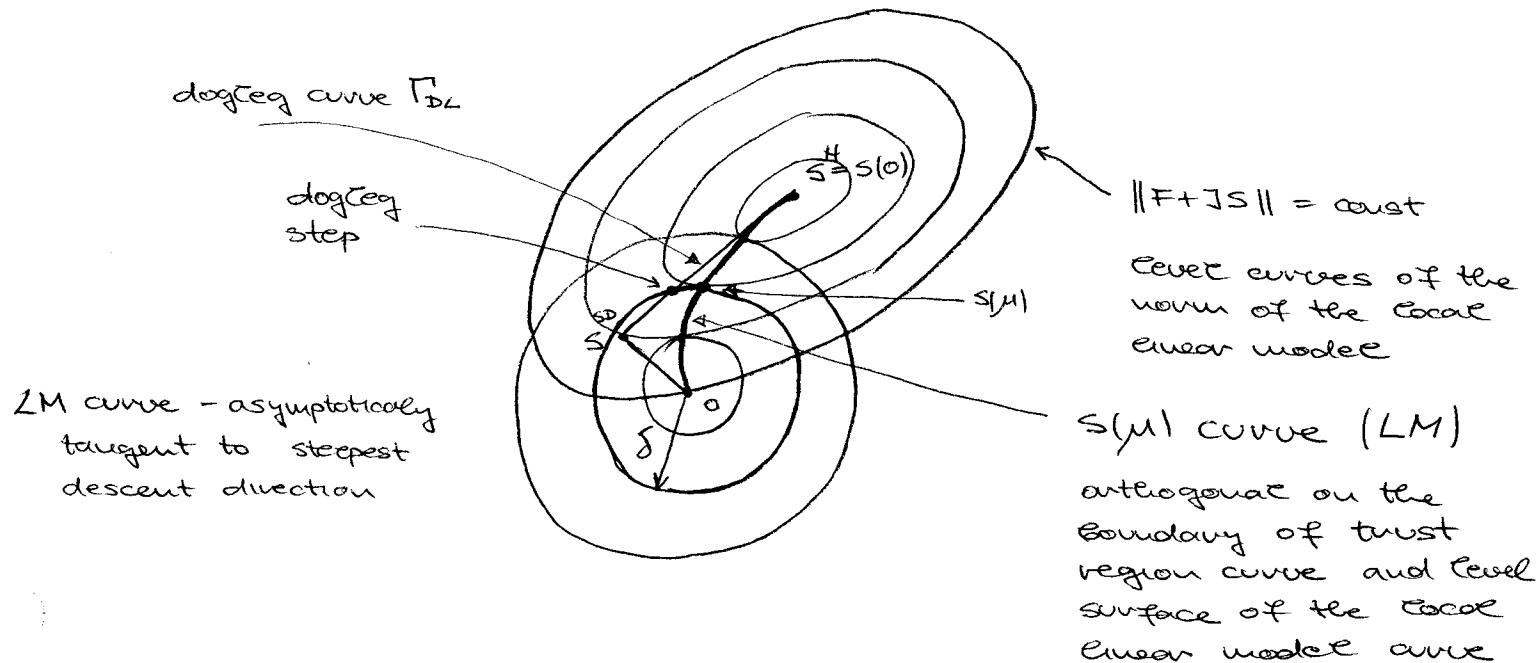


Assuming $\| \cdot \| = \| \cdot \|_2$ throughout

Fundamental practical difficulty:

can't compute $s(\mu)$ such that $\| s(\mu) \| = \delta$ both
accuracy and efficiency

Approach 1 (the Levenberg - Marquardt approach)



Compute $s(\mu)$ exactly on LM curve such that

$\| s(\mu) \|$ approximately equals δ

Set $\Phi(\mu) = \| s(\mu) \| - \delta$

Want μ_* such that $\Phi(\mu_*) = 0$

Problem: Computing $s(\mu) = -[J^T J + \mu I]^{-1} J^T F$

cost up to $O(n^3)$ arithmetic to factor

$J^T J + \mu I$ each time we change μ

Consider Newton's method. Have

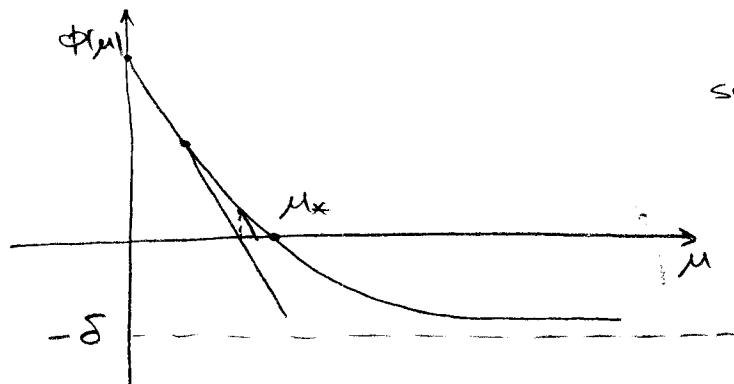
$$\dot{\Phi}(\mu) = \frac{-s(\mu)^T [J^T J + \mu I]^{-1} s(\mu)}{\| s(\mu) \|}$$

In computing $s(\mu)$ already have factored $J^T J + \mu I$
so can evaluate $[J^T J + \mu I]^{-1} s(\mu)$ and $\Phi'(\mu)$ in
 $\Theta(n^2)$ arithmetic.

- Already know $\|s(\mu)\|$ is monotone decreasing

$\Rightarrow \Phi(\mu)$ is monotone decreasing

Also, can show $|\Phi'(\mu)|$ is monotone decreasing



see: Newton step always undershoots μ^*

Newton step takes local linear model $\Phi + \Phi'$ s to zero

Consider new model

$$w(\delta) = \frac{\alpha}{\beta + \mu + \delta} - \delta$$

choose α, β such that

$$w(0) = \Phi(\mu)$$

$$w'(0) = \Phi'(\mu)$$

Get:

$$\alpha = - \frac{[\Phi(\mu) + \delta]^2}{\Phi'(\mu)} \quad \beta = \frac{-\Phi(\mu) + \delta}{\Phi'(\mu)} - \mu$$

choose δ such that $w(\delta) = 0$

Get

$$\delta = \frac{\alpha}{\beta} - (\beta + \mu) = \underbrace{\frac{\|s(\mu)\|}{\delta}}_{\text{Newton step } \delta^*} \left\{ - \frac{\Phi(\mu)}{\Phi'(\mu)} \right\}$$

see this has desired effect of

$$\delta = \begin{cases} \text{longer than } \delta^H & \text{if } \mu < \mu^* \\ \text{shorter than } \delta^H & \text{if } \mu > \mu^* \end{cases}$$

Also as $\mu \rightarrow \mu^*$ $\frac{\|s(\mu)\|}{\delta} \rightarrow 1$ and $\delta \rightarrow \delta^H$

can show quadratic convergence to μ^*

- Dennis & Schnabel recommend stopping once

$$\frac{3}{4} \delta \leq \|s(\mu)\| \leq \frac{3}{2} \delta$$

Approach 2 (The dogleg approach)

- Find step s such that $\|s\| = \delta$ exactly but s is only approximately on the $s(\mu)$ - curve
- Construct dogleg curve connecting $s=0$, $s=s^{SD}$, $s=s^H$

Here s^{SD} - "steepest descent step"

$$s^{SD} \text{ minimizer of } \|F + JS\| \text{ for } s = \lambda \left(-\nabla_{\frac{1}{2}} \|F\|^2 \right) \\ = -\lambda \underbrace{J^T F}_{\substack{\text{steepest} \\ \text{descent} \\ \text{dir. for } \|F\|}}$$

$$\text{see: } \Psi(\lambda) = \frac{1}{2} \|F + J(\lambda s)\|^2, \quad s = -J^T F \\ = \frac{1}{2} \|F\|^2 + \lambda F^T J s + \frac{1}{2} \lambda^2 \|Js\|^2$$

$$0 = \Psi'(\lambda) = F^T J s + \lambda \|Js\|^2 \\ \Rightarrow \lambda = -\frac{F^T J s}{\|Js\|^2} \sim \text{easy to compute}$$

see: Γ_{DL} ends at s^H

- A short step (corresponding to small δ) along Γ_{DL} is in steepest-descent direction (similar to $s(\mu)$ for large μ)

Can show, as s traverses Γ_{DL} from 0 to s^{SD} to s^H

(1) $\|s\|$ is monotone increasing

(2) $\|F+Js\|$ is monotone decreasing

So \exists unique point s^{DL} (the dogleg step) where Γ_{DL} intersects the TR boundary and

$$s^{DL} = \underset{s \in \Gamma_{DL}}{\operatorname{argmin}} \|F+Js\|$$

$$\|s\| \leq \delta$$

Computing s^{DL}

(1) If $\|s^H\| \leq \delta \Rightarrow s^{DL} = s^H$

(2) If $\|s^H\| > \delta$, compute s^{SD}

(3) If $\|s^{SD}\| > \delta$, then $s^{DL} = \frac{\delta}{\|s^{SD}\|} s^{DL}$

(4) If $\|s^{SD}\| < \delta$, then compute $s^{DL} = s^{SD} + \tau(s^H - s^{SD})$
for unique $\tau \in (0,1)$ such
that $\|s^{DL}\| = \delta$

Finding τ : Want

$$\|s^{SD} + \tau(s^H - s^{SD})\|^2 - \delta^2 = 0$$

$$(\|s^{SD}\|^2 - \delta^2) + 2s^{SD^T}(s^H - s^{SD})\tau + \|s^H - s^{SD}\|^2\tau^2 = 0$$

or

$$a\tau^2 + 2b\tau + c = 0$$

see

$a < 0, b > 0$ since $\|s\|$ is monotone decreasing
as s goes from 0 to s^{SD} to s^H

and $a > 0$

$$\text{so } \tau_{\pm} = \frac{-b \pm \sqrt{b^2 - ac}}{a}$$

want $\tau = \tau_+$ since we want $0 < \tau < 1$

$$\tau = \frac{-b + \sqrt{b^2 - ac}}{a}$$

to avoid cancellation, use

$$\tau = \frac{-b + \sqrt{b^2 - ac}}{a} \cdot \frac{b + \sqrt{b^2 - ac}}{b + \sqrt{b^2 - ac}} = \frac{-b^2 + b^2 - ac}{(b + \sqrt{b^2 - ac})^2}$$
$$\tau = \frac{-c}{b + \sqrt{b^2 - ac}}$$

Considerations for optimization

- Want to minimize $f: \mathbb{R}^n \rightarrow \mathbb{R}^1$
Use Newton + globalization to solve $\nabla f(x) = 0$
- Want to reduce f at each step so base acceptability etc. on local quadratic model

$$L(s) = f(x) + \nabla f(x)^T s + \frac{1}{2} s^T \nabla^2 f(x) s$$

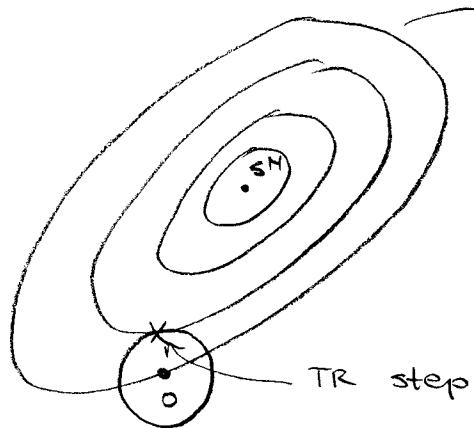
vecall: $s^H = -\nabla^2 f \nabla f$ is guaranteed to be descent step (for local quadratic model) only if $\nabla^2 f$ is SPD. If necessary perturb $\nabla^2 f$ to $B = \nabla^2 f + \mu I$ where $\mu > 0$ such that B is SPD and well-conditioned

For a TR method, choose

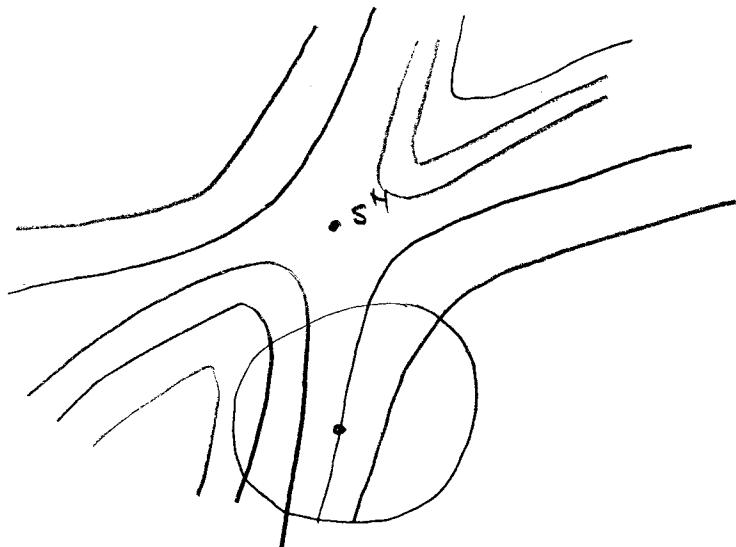
$$s = \underset{\|w\| \leq \delta}{\operatorname{argmin}} L(w)$$

New issue: If $\nabla^2 f$ indefinite (negative curvature - some of the eigenvalues are negative) then L has no global minimizer and local minimization within TR may be complicated

$\nabla^2 f$ - SPD



$\nabla^2 f$ - has one negative
and one positive
eigenvalue



$$\text{Can show: TR step} = -(\nabla^2 f + \mu I)^{-1} \nabla f$$

where $\mu > 0$ is such that $\nabla^2 f + \mu I$ is SPD

Traditional approach (see Dennis & Schnabel)

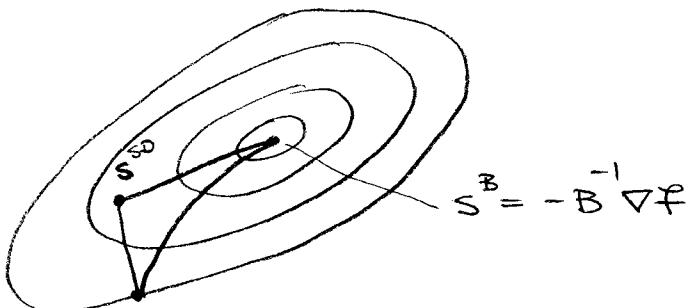
- Penalties if necessary to get $B = \nabla^2 f + \mu I$, $\mu > 0$ that is "safely" SPD.

Then take $g(s) = f + \nabla f^T s + \frac{1}{2} s^T B s$

Then have as before

$$\text{LM step } s(\mu) = -(\nabla^2 f + \mu I)^{-1} \nabla f$$

steepest descent
step does not
depend on Hessian
just on first
order derivative



Dogleg curve: $0 \rightarrow S = \underset{SD}{\text{argmin}} \mathcal{L}(S)$
 $S = -\lambda \nabla f$
for $\lambda > 0$
 $\rightarrow S^B$

(same important properties as before)

So as before, can compute

- LM step $s(\mu)$ such that $\|s(\mu)\| \approx \delta$
- dogleg step: S on dogleg curve with $\|S\| = \delta$

But this $(\nabla^2 f \leftarrow B = \nabla^2 f + \mu I)$ doesn't take advantage of negative curvature

How to do this?

- still have $s = s(\mu) = -(\nabla^2 f + \mu I)^{-1} \nabla f$ for $\mu > 0$ such that $\nabla^2 f + \mu I$ is SPD and $\|s(\mu)\| = \delta$
 - since $\nabla^2 f$ is symmetric, have $\nabla^2 f = U \Lambda U^T$ where $\Lambda = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$; $U^T U = I$ and $\lambda_1 > \dots > \lambda_n$
- λ_n - the most negative eigenvalue

$$\text{Have } \nabla^2 f + \mu I = U \Lambda U^T + \mu \underbrace{U U^T}_{I}$$

$$= U (\Lambda + \mu I) U^T$$

$$\Rightarrow (\nabla^2 f + \mu I)^{-1} = U \begin{pmatrix} \frac{1}{\lambda_1 + \mu} & & \\ & \ddots & \\ & & \frac{1}{\lambda_n + \mu} \end{pmatrix} U^T$$

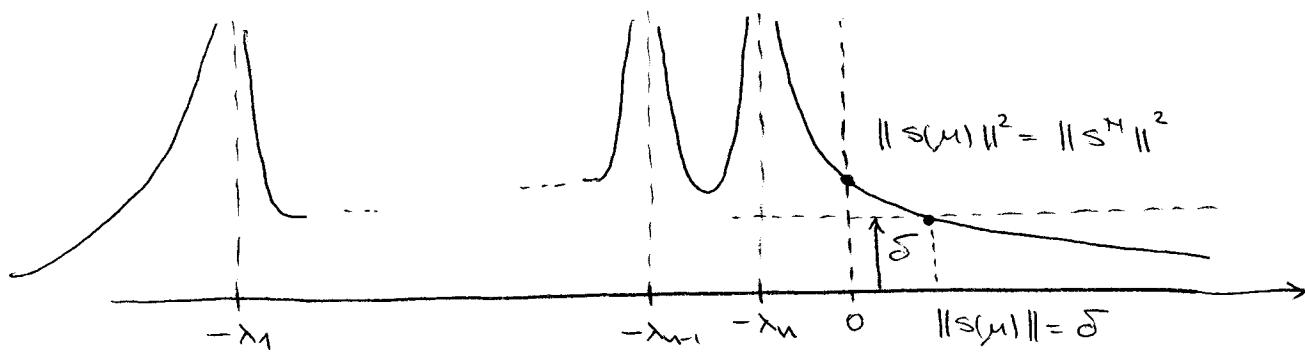
$$\Rightarrow s(\mu) = -(\nabla^2 f + \mu I)^{-1} \nabla f = U \begin{pmatrix} \frac{x_1}{\lambda_1 + \mu} \\ \vdots \\ \frac{x_n}{\lambda_n + \mu} \end{pmatrix}$$

$$\text{where } \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = -U^T \nabla f$$

so x_i is component of gradient in direction of i^{th} column of U

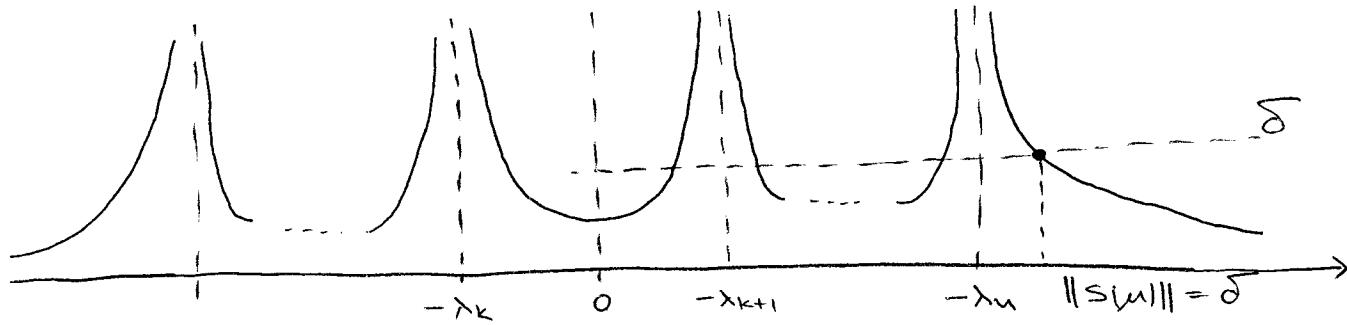
$$\|S(\mu)\|^2 = \|U \begin{pmatrix} \frac{\lambda_1}{\lambda_1 + \mu} \\ \vdots \\ \frac{\lambda_n}{\lambda_n + \mu} \end{pmatrix}\|^2 = \sum_{i=1}^n \frac{\lambda_i^2}{(\lambda_i + \mu)^2}$$

Suppose $\nabla^2 f$ is SPD. Then $\lambda_1 > \dots > \lambda_n > 0$ and $\|S(\mu)\|^2$ is a nice (analytic) function of μ for $\mu > 0$ but has poles at $-\lambda_1, \dots, -\lambda_n$.



Suppose $\nabla^2 f$ is indefinite

$$\lambda_1 > \dots > \lambda_k > 0 > \lambda_{k+1} > \dots > \lambda_n$$



Can compute $S(\mu)$ accurately but may be moderately expensive