
Lectures 7 and 8: Trust-region methods for unconstrained optimization

Coralia Cartis, Mathematical Institute, University of Oxford

C6.2/B2: Continuous Optimization

Linesearch versus trust-region methods

(UP): minimize $f(x)$ subject to $x \in \mathbb{R}^n$.

Linesearch methods: ‘liberal’ in the choice of search direction, keeping bad behaviour in control by choice of α^k .

- choose descent direction s^k ,
- compute stepsize α^k to reduce $f(x^k + \alpha s^k)$,
- update $x^{k+1} := x^k + \alpha^k s^k$.

Trust region (TR) methods: ‘conservative’ in the choice of search direction, so that a full stepsize along it may really reduce the objective.

- pick direction s^k to reduce a “local model” of $f(x^k + s^k)$,
- accept $x^{k+1} := x^k + s^k$ if decrease in the model is also achieved by $f(x^k + s^k)$,
- else set $x^{k+1} := x^k$ and “refine” the model.

Trust-region models for unconstrained problems

Approximate $f(x^k + s)$ by:

- linear model $l_k(s) := f(x^k) + s^\top \nabla f(x^k)$ or

- quadratic model

$$q_k(s) := f(x^k) + s^\top \nabla f(x^k) + \frac{1}{2} s^\top \nabla^2 f(x^k) s.$$

Impediments:

models may not resemble $f(x^k + s)$ when s is large,

models may be unbounded from below,

- * $l_k(s)$ always unbounded below (unless $\nabla f(x^k) = 0$)

- * $q_k(s)$ is always unbounded below if $\nabla^2 f(x^k)$ is negative definite or indefinite, and sometimes if $\nabla^2 f(x^k)$ is positive semidefinite.

Trust region models and subproblem

Prevent bad approximations by trusting the model only in a **trust region**, defined by the **trust region constraint**

$$\|s\| \leq \Delta_k, \quad (\text{R})$$

for some “appropriate” **radius** $\Delta_k > 0$.

The constraint (R) also prevents l_k, q_k from unboundedness!

\implies **the trust region subproblem**

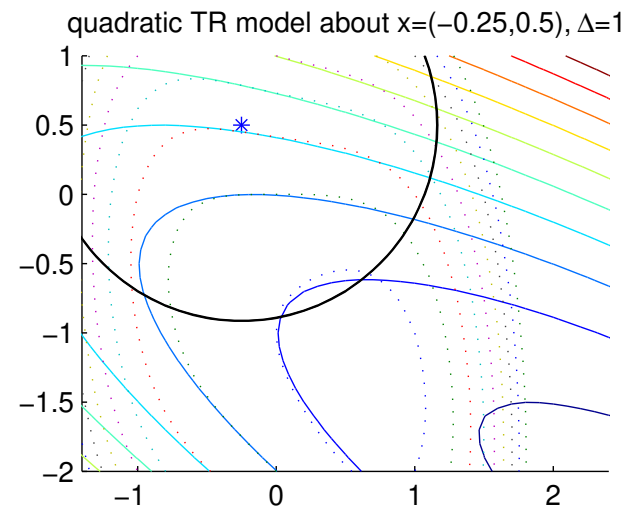
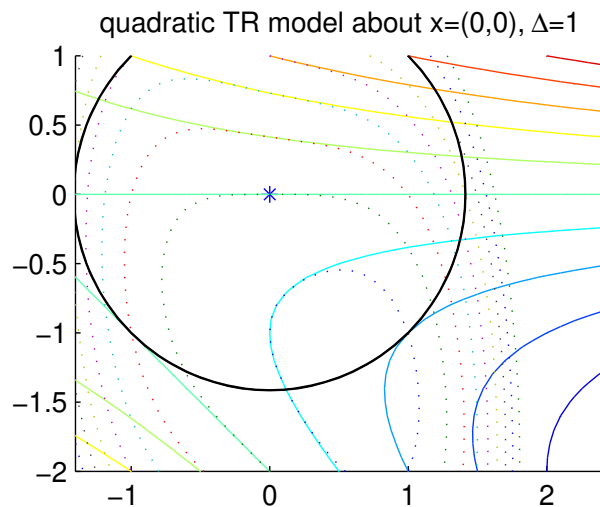
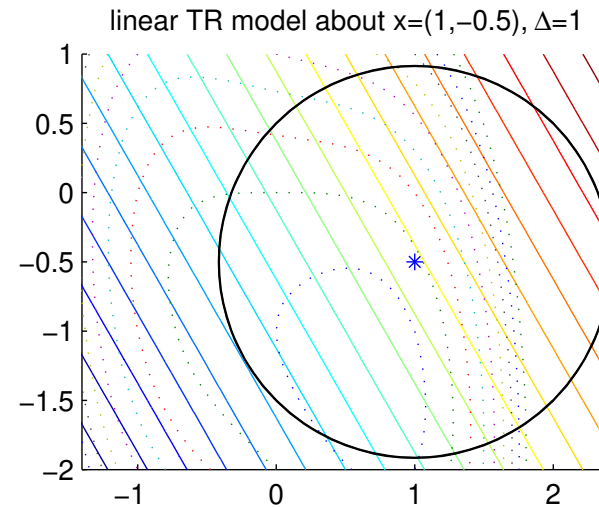
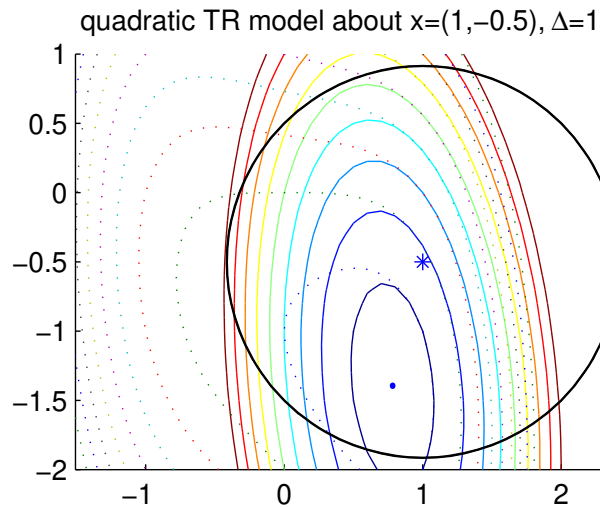
$$\min_{s \in \mathbb{R}^n} m_k(s) \quad \text{subject to} \quad \|s\| \leq \Delta_k, \quad (\text{TR})$$

where $m_k := l_k, k \geq 0$, or $m_k := q_k, k \geq 0$.

- **From now on, $m_k := q_k$.**

(TR) easier to solve than (P). May even solve (TR) only approximately.

Trust region models and subproblem - an example



Trust-region models of $f(x) = x_1^4 + x_1x_2 + (1 + x_2)^2$.

Generic trust-region method

Let s^k be a(n approximate) solution of (TR). Then

- predicted model decrease:

$$m_k(0) - m_k(s^k) = f(x^k) - m_k(s^k).$$

- actual function decrease: $f(x^k) - f(x^k + s^k)$.

The trust region radius Δ_k is chosen based on the value of

$$\rho_k := \frac{f(x^k) - f(x^k + s^k)}{f(x^k) - m_k(s^k)}.$$

If ρ_k is not too smaller than 1, $x^{k+1} := x^k + s^k$, $\Delta_{k+1} \geq \Delta_k$.

If ρ_k close to or ≥ 1 , Δ_k is increased.

If $\rho_k \ll 1$, $x^{k+1} = x^k$ and Δ_k is reduced.

A Generic Trust Region (GTR) method

Given $\Delta_0 > 0$, $x^0 \in \mathbb{R}^n$, $\epsilon > 0$. While $\|\nabla f(x^k)\| \geq \epsilon$, do:

1. Form the local quadratic model $m_k(s)$ of $f(x^k + s)$.
2. Solve (approximately) the (TR) subproblem for s^k with $m_k(s^k) < f(x^k)$ ("sufficiently").

Compute $\rho_k := [f(x^k) - f(x^k + s^k)]/[f(x^k) - m_k(s^k)]$.

3. If $\rho_k \geq 0.9$, then [very successful step]
set $x^{k+1} := x^k + s^k$ and $\Delta_{k+1} := 2\Delta_k$.

Else if $\rho_k \geq 0.1$, then [successful step]
set $x^{k+1} := x^k + s^k$ and $\Delta_{k+1} := \Delta_k$.

Else [unsuccessful step]
set $x^{k+1} = x^k$ and $\Delta_{k+1} := \frac{1}{2}\Delta_k$.

4. Let $k := k + 1$. □

Trust-region methods

- Other sensible values of the parameters of the GTR are possible.

“Solving” the (TR) subproblem

$$\min_{s \in \mathbb{R}^n} m_k(s) \quad \text{subject to} \quad \|s\| \leq \Delta_k, \quad (\text{TR})$$

... exactly or even approximately may imply work.

Want “minimal” condition of “sufficient decrease” in the model that ensures global convergence of the TR method (the Cauchy cond.). In practice, we (usually) do much better than this condition!

Example of applying a trust-region method: [Sartenaer, 2008].

- approximate solution of (TR) subproblem: better than Cauchy, but not exact.
- notation: $\Delta f / \Delta m_k \equiv \rho_k$.

The Cauchy point of the (TR) subproblem

- recall the steepest descent method has strong (theoretical) global convergence properties; same will hold for TR method with SD direction.

“minimal” condition of “sufficient decrease” in the model: require

$$m_k(s^k) \leq m_k(s_c^k) \text{ and } \|s^k\| \leq \Delta_k,$$

where $s_c^k := -\alpha_c^k \nabla f(x^k)$, with

$$\alpha_c^k := \arg \min_{\alpha > 0} m_k(-\alpha \nabla f(x^k)) \text{ subject to } \|\alpha \nabla f(x^k)\| \leq \Delta_k.$$

[i.e. a linesearch along steepest descent direction is applied to m_k at x^k and is restricted to the trust region.] Easy:

$$\alpha_c^k := \arg \min_{\alpha} m_k(-\alpha \nabla f(x^k)) \text{ subject to } 0 < \alpha \leq \frac{\Delta_k}{\|\nabla f(x^k)\|}.$$

- $y_c^k := x^k + s_c^k$ is the Cauchy point.
-

Global convergence of the GTR method

Theorem 11 (GTR global convergence)

Let $f \in \mathcal{C}^2(\mathbb{R}^n)$ and bounded below on \mathbb{R}^n . Let ∇f be Lipschitz continuous on \mathbb{R}^n . Let $\{x^k\}$ be generated by the generic trust region (GTR) method, and let the computation of s^k be such that $m_k(s^k) \leq m_k(s_c^k)$ for all k . Then either

there exists $k \geq 0$ such that $\nabla f(x^k) = 0$

or

$$\lim_{k \rightarrow \infty} \|\nabla f(x^k)\| = 0.$$

We (only) sketch the proof of $\liminf_{k \rightarrow \infty} \|\nabla f(x^k)\| = 0$ (which also implies finite termination of GTR) next.

Computation of the Cauchy point

Computation of the Cauchy point: find α_c^k global solution of

$$\min_{\alpha > 0} m_k(-\alpha \nabla f(x^k)) \text{ subject to } \|\alpha \nabla f(x^k)\| \leq \Delta_k,$$

where $m_k(s) = f(x_k) + s^T \nabla f(x^k) + \frac{1}{2} s^T \nabla^2 f(x^k) s$, & $\nabla f(x^k) \neq 0$.

■ $\|\alpha \nabla f(x^k)\| \leq \Delta_k \text{ \& } \alpha > 0 \Leftrightarrow 0 < \alpha \leq \frac{\Delta_k}{\|\nabla f(x^k)\|} := \bar{\alpha}.$

■ $\phi(\alpha) := m_k(-\alpha \nabla f(x^k)) = f(x^k) - \alpha \|\nabla f(x^k)\|^2 + \frac{\alpha^2}{2} h^k,$

where $h^k := \nabla f(x^k)^T \nabla^2 f(x^k) \nabla f(x^k).$

■ $\phi'(0) = -\|\nabla f(x^k)\|^2 < 0$ so ϕ decreasing from $\alpha = 0$ for suff. small α ; thus $\alpha_c^k > 0$.

■ $h^k > 0$: $\alpha_{\min} := \frac{\|\nabla f(x^k)\|^2}{h^k} = \arg \min_{\alpha > 0} \phi(\alpha).$

$\implies \alpha_c^k = \min(\alpha_{\min}, \bar{\alpha}).$

■ $h^k \leq 0$: $\phi(\alpha)$ unbounded below on \mathbb{R} and so $\alpha_c^k = \bar{\alpha}.$

Proof of global convergence of the GTR method

Lemma 12: (Cauchy model decrease) In GTR with Cauchy decrease $m_k(s^k) \leq m_k(s_c^k)$ for all k , we have the model decrease for each k ,

$$\begin{aligned} f(x^k) - m_k(s^k) &\geq f(x^k) - m_k(s_c^k) \\ &\geq \frac{1}{2} \|\nabla f(x^k)\| \min \left\{ \Delta_k, \frac{\|\nabla f(x^k)\|}{1 + \|\nabla^2 f(x^k)\|} \right\} \end{aligned}$$

Proof of Lemma 12. (Recall Computation of the Cauchy point)

If $h^k \leq 0$, then $m_k(-\alpha_c^k \nabla f(x^k)) \leq f(x^k) - \alpha_c^k \|\nabla f(x^k)\|^2$. In this

case, we also have $\alpha_c^k = \bar{\alpha} = \frac{\Delta_k}{\|\nabla f(x^k)\|}$ and so

$$f(x^k) - m_k(s_c^k) \geq \Delta_k \|\nabla f(x^k)\|.$$

Else, $h^k > 0$; then $\alpha_c^k = \min\{\alpha_{\min}, \bar{\alpha}\}$ where $\alpha_{\min} = \|\nabla f(x^k)\|^2 / h^k$.

Assume first that $\alpha_c^k = \bar{\alpha}$. Then $\alpha_c^k h^k \leq \|\nabla f(x^k)\|^2$ and

$$f(x^k) - m_k(s_c^k) = \alpha_c^k \|\nabla f(x^k)\|^2 - \frac{(\alpha_c^k)^2}{2} h^k \geq \frac{\alpha_c^k}{2} \|\nabla f(x^k)\|^2,$$

Proof of global convergence of the GTR method

Proof of Lemma 12 (continued).

and using the expression of $\bar{\alpha}$,

$$f(x^k) - m_k(s_c^k) \geq \frac{\Delta_k}{2\|\nabla f(x^k)\|} \|\nabla f(x^k)\|^2 = \frac{1}{2} \Delta_k \|\nabla f(x^k)\|.$$

Finally, let $\alpha_c^k = \alpha_{\min} = \|\nabla f(x^k)\|^2 / h^k$. Replacing this value in the model decrease we get

$$f(x^k) - m_k(s_c^k) = \alpha_c^k \|\nabla f(x^k)\|^2 - \frac{(\alpha_c^k)^2}{2} h^k = \frac{\|\nabla f(x^k)\|^4}{2h^k},$$

and further, by Cauchy-Schwarz and Rayleigh quotient inequalities,

$$\begin{aligned} \frac{\|\nabla f(x^k)\|^4}{2h^k} &= \frac{\|\nabla f(x^k)\|^4}{2(\nabla f(x^k))^T \nabla^2 f(x^k) \nabla f(x^k)} \\ &\geq \frac{\|\nabla f(x^k)\|^2}{2\|\nabla^2 f(x^k)\|} \geq \frac{\|\nabla f(x^k)\|^2}{2(1+\|\nabla^2 f(x^k)\|)} \quad (*) \end{aligned}$$

Thus $f(x^k) - m_k(s_c^k) \geq \frac{\|\nabla f(x^k)\|^2}{2(1+\|\nabla^2 f(x^k)\|)}$. \square

[(*) '+1' is only needed to cover the case $H^k = 0$.]

Proof of global convergence of the GTR method

Lemma 13: (Lower bound on TR radius) Let $f \in \mathcal{C}^2(\mathbb{R}^n)$ and ∇f be Lipschitz continuous on \mathbb{R}^n with Lipschitz constant L . In GTR with Cauchy decrease $m_k(s^k) \leq m_k(s_c^k)$ for all k , suppose that

there exists $\epsilon > 0$ such that $\|\nabla f(x^k)\| \geq \epsilon$ for all k .

Then, there exists a constant $c \in (0, 1)$ (independent of k) such that

$$\Delta_k \geq \frac{c}{L}\epsilon \quad \text{for all } k \geq 0.$$

Remarks:

(1) The proof of Lemma 13 relies on first showing that if

$\Delta_k \leq \frac{2c}{L}\epsilon$, then iteration k is successful and $\Delta_{k+1} \geq \Delta_k$.

(2) If GTR takes finitely many successful iterations, then we can show that the last successful iterate has zero gradient.

[$\Delta_k \rightarrow 0$ which contradicts L13 if gradient not zero.]

Proof of global convergence of the GTR method

Theorem 14: (At least one limit point is stationary) Let $f \in \mathcal{C}^2(\mathbb{R}^n)$ and and bounded below on \mathbb{R}^n . Let ∇f be Lipschitz continuous on \mathbb{R}^n with Lipschitz constant L . Let $\{x^k\}$ be generated by the generic trust region (GTR) method, and let the computation of s^k be such that $m_k(s^k) \leq m_k(s_c^k)$ for all k . Then either there exists $k \geq 0$ such that $\nabla f(x^k) = 0$ or $\liminf_{k \rightarrow \infty} \|\nabla f(x^k)\| = 0$.

Proof of Theorem 14. If there exists k such that $\nabla f(x^k) = 0$, then GTR terminates (this includes the case of having finitely many successful iterations). Assume there exists $\epsilon > 0$ such that $\|\nabla f(x^k)\| \geq \epsilon$ for all k . Then using that there are infinitely many successful iterations $k \in \mathcal{S}$, and def. of GTR/ ρ_k ,

$$\begin{aligned} f(x^k) - f(x^{k+1}) &\geq 0.1(f(x^k) - m_k(s^k)) \\ &\geq \frac{0.1}{2} \|\nabla f(x^k)\| \min\left\{\frac{\|\nabla f(x^k)\|}{1 + \|\nabla^2 f(x^k)\|}, \Delta_k\right\} \end{aligned}$$

for all $k \in \mathcal{S}$, where we also used Lemma 12.

Proof of global convergence of the GTR method

Proof of Theorem 14 (continued).

∇f Lipschitz cont. with Lips const $L \implies \|\nabla^2 f(x)\| \leq L \forall x$.

Thus since $\|\nabla f(x^k)\| \geq \epsilon$ for all k , we have for all $k \in \mathcal{S}$ that

$$f(x^k) - f(x^{k+1}) \geq 0.05\epsilon \min \left\{ \frac{\epsilon}{1+L}, \Delta_k \right\} \geq 0.05\epsilon \min \left\{ \frac{\epsilon}{1+L}, \frac{c}{L}\epsilon \right\},$$

where we also used Lemma 13. Thus

$$\text{for all } k \in \mathcal{S}: f(x^k) - f(x^{k+1}) \geq \frac{0.05c}{1+L}\epsilon^2. \quad (*)$$

Since $f(x^k) \geq f_{\text{low}}$ for all k , we deduce

$$\begin{aligned} f(x^0) - f_{\text{low}} &\geq f(x^0) - \lim_{k \rightarrow \infty} f(x^k) \geq \sum_{i=0}^{\infty} (f(x^i) - f(x^{i+1})) \\ &= \sum_{i \in \mathcal{S}} (f(x^i) - f(x^{i+1})) \geq |\mathcal{S}| \frac{0.05c}{1+L}\epsilon^2 \quad (**) \end{aligned}$$

where in '=' we used $f(x^k) = f(x^{k+1})$ on all unsuccessful k , and in the last ' \geq ', we used (*) and $|\mathcal{S}| = \text{no. of successful iterations}$. But LHS of (**) is finite while RHS of (**) is infinite since $|\mathcal{S}| = \infty$. Thus there must exist k such that $\|\nabla f(x^k)\| < \epsilon$. \square

Solving the (TR) subproblem

On each TR iteration we compute or approximate the solution of

$$\min_{s \in \mathbb{R}^n} m_k(s) = f(x^k) + s^\top \nabla f(x^k) + \frac{1}{2} s^\top \nabla^2 f(x^k) s$$

subject to $\|s\| \leq \Delta_k$.

- also, s^k must satisfy the Cauchy condition $m_k(s^k) \leq m_k(s_c^k)$, where $s_c^k := -\alpha_c^k \nabla f(x^k)$, with

$$\alpha_c^k := \arg \min_{\alpha > 0} m_k(-\alpha \nabla f(x^k)) \text{ subject to } \|\alpha \nabla f(x^k)\| \leq \Delta_k.$$

[Cauchy condition ensures global convergence]

- solve (TR) exactly (i.e., compute global minimizer of TR)
 \implies TR akin to Newton-like method.
 - solve (TR) approximately (i.e., an approximate global minimizer) \implies large-scale problems.
-

Solving the (TR) subproblem exactly

For $h \in \mathbb{R}$, $\Delta > 0$, $g \in \mathbb{R}^n$, H $n \times n$ symm. matrix, consider

$$\min_{s \in \mathbb{R}^n} m(s) := h + s^\top g + \frac{1}{2} s^\top H s, \text{ s. t. } \|s\| \leq \Delta. \quad (\text{TR})$$

Characterization result for the solution of (TR):

Theorem 15

Any **global** minimizer s^* of (TR) satisfies the equation

$$(H + \lambda^* I) s^* = -g,$$

where $H + \lambda^* I$ is positive semidefinite, $\lambda^* \geq 0$,

$$\lambda^* (\|s^*\| - \Delta) = 0 \quad \text{and} \quad \|s^*\| \leq \Delta.$$

If $H + \lambda^* I$ is positive definite, then s^* is unique.

- The above Theorem gives necessary **and** sufficient **global** optimality conditions for a **nonconvex** optimization problem!

Solving the (TR) subproblem exactly

Computing the global solution s^* of (TR):

Case 1. If H is positive definite and $Hz = -g$ satisfies $\|z\| \leq \Delta \implies s^* := z$ (unique), $\lambda^* := 0$ (by Theorem 15).

Case 2. If H is positive definite but $\|z\| > \Delta$, or H is not positive definite, Theorem 15 implies s^* satisfies

$$(H + \lambda I)s = -g, \quad \|s\| = \Delta, \quad (*)$$

for some $\lambda \geq \max\{0, -\lambda_{\min}(H)\} := \underline{\lambda}$.

Let $s(\lambda) = -(H + \lambda I)^{-1}g$, for any $\lambda > \underline{\lambda}$. Then $s^* = s(\lambda^*)$ where $\lambda^* \geq \underline{\lambda}$ solution of

$$\|s(\lambda)\| = \Delta, \quad \lambda \geq \underline{\lambda}.$$

\longrightarrow nonlinear equation in one variable λ . Use Newton's method to solve it. We discuss the system (*) in detail next.

Solving the (TR) subproblem exactly ...

$$(H + \lambda I)s = -g, \quad s^\top s = \Delta^2. \quad (*)$$

H symmetric \implies spectral decomposition: $H = U^\top \Lambda U$,
with U orthonormal matrix of the eigenvectors of H and Λ
diagonal mat. of eigenvalues of H , $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$; $\lambda_1 = \lambda_{\min}(H)$

Th. 15 $\implies H + \lambda I = U^\top (\Lambda + \lambda I)U$ positive semidefinite \implies
 $\lambda_1 + \lambda \geq 0 \implies \lambda \geq -\lambda_1 \implies \lambda \geq \max\{0, -\lambda_1\}$.

$\lambda \longrightarrow s(\lambda) := -(H + \lambda I)^{-1}g$, provided $H + \lambda I$ nonsingular.

$$\psi(\lambda) := \|s(\lambda)\|^2 = \|U^\top (\Lambda + \lambda I)^{-1}Ug\|^2 = g^\top U^\top (\Lambda + \lambda I)^{-2}Ug$$

• $g = U^\top \gamma$, for some $\gamma = (\gamma_1, \dots, \gamma_n) \in \mathbb{R}^n$. As $UU^\top = U^\top U = I$,

$$\psi(\lambda) = \gamma^\top (\Lambda + \lambda I)^{-2} \gamma = \sum_{i=1}^n \frac{\gamma_i^2}{(\lambda + \lambda_i)^2} \stackrel{(*)}{=} \Delta^2.$$

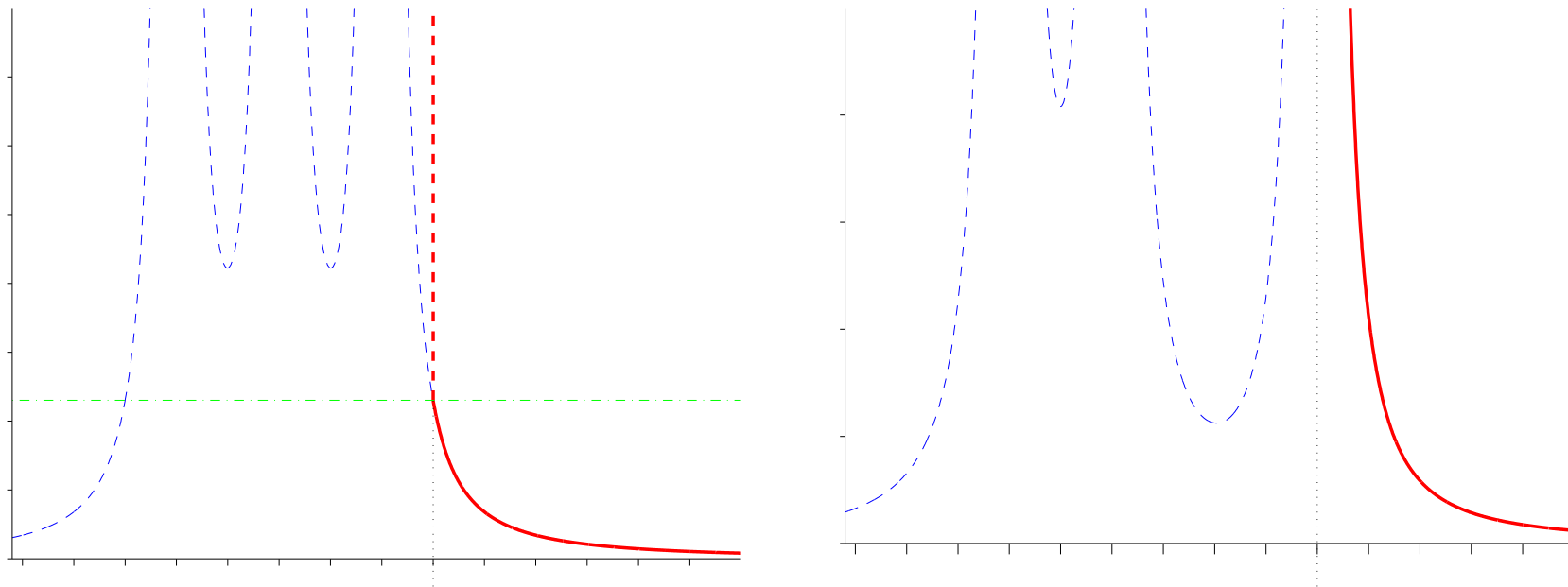
The secular equation

Consider

$$\psi(\lambda) := \|s(\lambda)\|^2 = \Delta^2$$

for $\lambda \in (\max\{0, -\lambda_1\}, \infty)$.

[see Pb Sheet 3]



‘Easy’ cases: Plots of λ vs. $\psi(\lambda)$; $H \succ 0$ (LHS) and H indef (RHS).

The secular equation

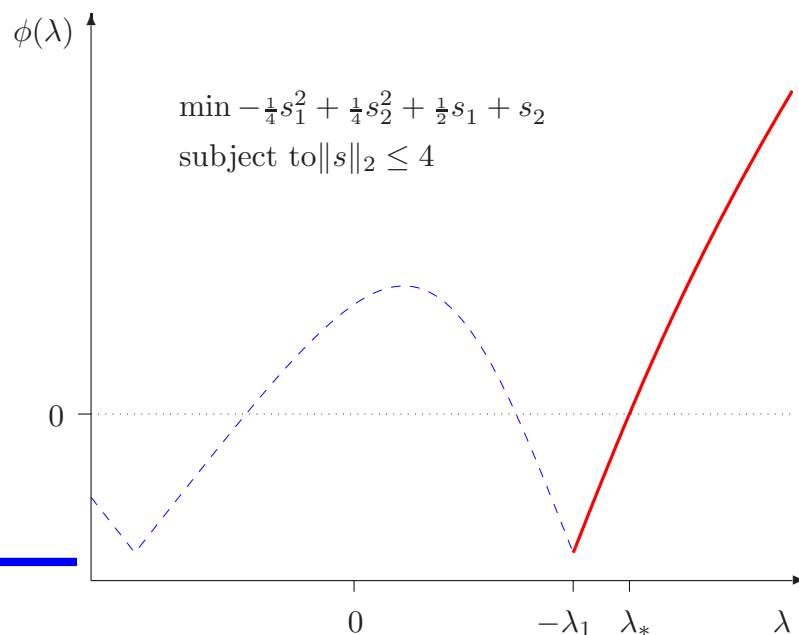
DON'T solve $\psi(\lambda) := \|s(\lambda)\|^2 = \Delta^2$.

Solve instead the **secular equation**

$$\phi(\lambda) := \frac{1}{\|s(\lambda)\|} - \frac{1}{\Delta} = 0 \text{ for } \lambda \in (\max\{0, -\lambda_1\}, \infty). \quad (\dagger)$$

- ϕ has no poles; it is analytic on $(-\lambda_1, \infty)$
 \implies ideal for Newton's mthd (exc. in the 'hard' case).

[globally convergent and locally quadratic if $\lambda^0 \in [-\lambda_1, \lambda_*]$; else safeguard with linesearch]



Solving the (TR) subproblem for large-scale problems

- Newton's mthd for (†): Cholesky factorization of $H + \lambda I$ for various $\lambda \rightarrow$ expensive or impossible for large problems.

No computation of the complete eigenvalue decomposition of H !

Solving the **large-scale** (TR) subproblem:

- Use iterative methods to **approximate** the global minimizer of (TR).

Use the Cauchy point (i.e. steepest descent):
impractical.

Use conjugate-gradient or Lanczos method (as the first step is a steepest descent, and thus our requirement of “sufficient decrease” in m_k will be satisfied).

Nonlinear least-squares/inverse problems

- $r : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $m \geq n$; r smooth.

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{1}{2} \sum_{j=1}^m [r_j(x)]^2 = \frac{1}{2} \|r(x)\|^2. \quad (\text{NLS})$$

The Levenberg-Marquardt method: replace linesearch in Gauss-Newton with trust-region

$$\implies \min_{s \in \mathbb{R}^n} \frac{1}{2} \|J(x^k)s + r(x^k)\|^2 \text{ subject to } \|s\| \leq \Delta_k.$$

- useful when $J(x^k)$ is rank-deficient (i.e., not full-rank); overcomes weakness of Gauss-Newton.

- s^k solves TR subproblem iff $\exists \lambda^k \geq 0$ such that

$$(J(x^k)^T J(x^k) + \lambda^k I) s^k = -J(x^k)^T r(x^k)$$

$$\text{and } \lambda^k (\|s^k\| - \Delta_k) = 0.$$

Linesearch vs. trust-region methods

Quasi-Newton methods/approximate derivatives also possible in the trust-region framework; no need for positive definite updates for the Hessian! Replace $\nabla^2 f(x^k)$ with approximation B^k in the quadratic local model $m_k(s)$.

Conclusions: state-of-the-art software for unconstrained problems implements linesearch or TR methods; both approaches have been made competitive (more heuristics needed by linesearch methods to deal with negative curvature). Choosing between the two is mostly a matter of “taste”.

Information on existing software can be found at the NEOS Center: <http://www.neos-guide.org>

→ look under [Optimization Guide](#) and [Optimization Tree](#), etc.
State-of-the-art NLO software: KNITRO, IPOPT, GALAHAD,...
