

Trust Region Method

Lectures for PHD course on
Unconstrained Numerical Optimization

Enrico Bertolazzi
DIMS – Università di Trento
May 2008

Outline

- 1 The Trust Region method
- 2 Convergence analysis
- 3 The exact solution of trust region step
- 4 The dogleg trust region step
- 5 The double dogleg trust region step
- 6 Two dimensional subspace minimization

- Newton and quasi-Newton methods approximate a solution iteratively by choosing at each step a search direction and minimize in this direction.
- An alternative approach is to find a direction and a step-length, then if the step is successful in some sense the step is accepted. Otherwise another direction and step-length is chosen.
- The choice of the step-length and direction is algorithm dependent but a successful approach is the one based on trust region.

- Newton and quasi-Newton at each step (approximately) solve the minimization problem

$$\arg \min_s m_k(s)$$

$$m_k(s) = f(x_k) + \nabla f(x_k)s + \frac{1}{2}s^T H_k s$$

in the case H_k is symmetric and positive definite (SPD).

- If H_k is SPD the minimum is

$$s = -H_k^{-1}g_k, \quad g_k = \nabla f(x_k)^T$$

and s is the quasi-Newton step.

- If $H_k = \nabla^2 f(x_k)$ and is SPD, then $s = -\nabla^2 f(x_k)^{-1} \nabla f(x_k)^T$ is the Newton step.

- If H_k is not positive definite, the search direction $-H_k^{-1}g_k$ may fail to be a descent direction and the previous minimization problem can have no solution.
- The problem is that the model $m_k(s)$ is an approximation of $f(x)$

$$m_k(s) \approx f(x_k + s)$$

and this approximation is valid only in a small neighbors of x_k .

- So that an alternative minimization problem is the following

$$\arg \min_s m_k(s) \quad \text{subject to } \|s\| \leq \Delta_k$$

Δ_k is the radius of the trust region of the model $m_k(s)$, i.e. the region where we trust the model is valid.



Algorithm (Generic trust region algorithm)

```

x assigned; Δ assigned;
while ‖∇f(x)‖ > ε do
  — setup the model
  m(s) = f(x) + ∇f(x)s + ½sTHS;
  — compute the step
  s ← arg min‖s‖ ≤ Δ m(s);
  xnew ← x + s;
  — check the reduction
  if is xnew acceptable? then
    x ← xnew;
    update Δ;
  else
    reduce Δ;
  end if
end while

```



When accept the step?

- The point x_{new} in the previous algorithm can be accepted or rejected. The acceptance criterium can be the Armijo criterium of sufficient decrease

$$f(x_{new}) \leq f(x) + \beta_0 \nabla f(x)(x_{new} - x)$$

where $\beta_0 \in (0, 1)$ is a small constant (typically 10^{-4}).

- Alternatively compute the expected and actual reduction with the ratio ρ :

$$p_{red} = m(0) - m(s), \quad a_{red} = f(x) - f(x + s),$$

$$\rho = a_{red}/p_{red}$$

If the ratio ρ is near 1 the match of the model with the real function is good. We accept the step if $\rho > \beta_1$ where $\beta_1 \in (0, 1)$ normally $\beta_1 \approx 0.1$.



If the step is rejected how to reduce the trust radius ?

- We construct the parabola $p(t)$ such that ($s = x_{new} - x$)

$$p(0) = f(x), \quad p'(0) = \nabla f(x)s, \quad p(\Delta) = f(x_{new}),$$

the solution is

$$p(t) = f(x) + (\nabla f(x)s)t + Ct^2$$

$$C = \frac{f(x_{new}) - f(x) - (\nabla f(x)s)\Delta}{\Delta^2}$$

- The new radius is on the minimum of the parabola:

$$\Delta_{new} = -\frac{(\nabla f(x)s)}{2C} = \frac{\Delta^2(\nabla f(x)s)}{2[f(x) + (\nabla f(x)s)\Delta - f(x_{new})]}$$

- A safety interval is normally assumed; if the new radius is outside $[\Delta/10, \Delta/2]$ then it is put again in this interval.



If the step is accepted how to modify the trust radius ?

- Compute the expected and actual reduction

$$p_{red} = m(\mathbf{0}) - m(\mathbf{s})$$

$$a_{red} = f(\mathbf{x}) - f(\mathbf{x} + \mathbf{s})$$

- Compute the ratio of expected and actual reduction

$$\rho = \frac{a_{red}}{p_{red}}$$

- Compute the new radius

$$\Delta_{new} = \begin{cases} \max\{2\|\mathbf{s}\|, \Delta\} & \text{if } \rho \geq \beta_2 \\ \Delta & \text{if } \rho \in (\beta_1, \beta_2) \\ \|\mathbf{s}\|/\Delta & \text{if } \rho \leq \beta_1 \end{cases}$$



Algorithm (Check reduction algorithm)

CheckReduction($\mathbf{x}, \mathbf{s}, \Delta$):

$\mathbf{x}_{new} \leftarrow \mathbf{x} + \mathbf{s}$

$\alpha \leftarrow \nabla f(\mathbf{x})\mathbf{s}$

$a_{red} \leftarrow f(\mathbf{x}) - f(\mathbf{x}_{new})$

$p_{red} \leftarrow -\alpha - \mathbf{s}^T \mathbf{H} \mathbf{s} / 2$

$\rho \leftarrow a_{red} / p_{red}$

$r_{new} \leftarrow \begin{cases} \max\{2\|\mathbf{s}\|, r\} & \text{if } \rho \geq \beta_2 \\ r & \text{if } \rho \in (\beta_1, \beta_2) \\ \|\mathbf{s}\|/2 & \text{if } \rho \leq \beta_1 \end{cases}$

if $\rho < \beta_1$ **then**
 — *reject the step*

$\mathbf{x}_{new} \leftarrow \mathbf{x}$

end if



Lemma

Consider the following constrained quadratic problem where $\mathbf{H} \in \mathbb{R}^{n \times n}$ *symmetric and positive definite*.

$$\text{Minimize } f(\mathbf{s}) = f_0 + \mathbf{g}^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H} \mathbf{s},$$

$$\text{Subject to } \|\mathbf{s}\| \leq \Delta$$

Then the following curve

$$\mathbf{s}(\mu) \doteq -(\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g},$$

for any $\mu \geq 0$ defines a descent direction for $f(\mathbf{s})$. Moreover

- there exists a unique μ_* such that $\|\mathbf{s}(\mu_*)\| = \Delta$ and $\mathbf{s}(\mu_*)$ is the solution of the constrained problem;
- or $\|\mathbf{s}(0)\| < \Delta$ and $\mathbf{s}(0)$ is the solution of the constrained problem.



Proof.

(1/2).

If $\|\mathbf{s}(0)\| \leq \Delta$ then $\mathbf{s}(0)$ is the global minimum of $f(\mathbf{s})$ which is inside the trust region. Otherwise consider the Lagrangian

$$\mathcal{L}(\mathbf{s}, \mu) = f_0 + \mathbf{g}^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H} \mathbf{s} + \frac{1}{2} \mu (\mathbf{s}^T \mathbf{s} - \Delta^2),$$

Then we have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{s}}(\mathbf{s}, \mu) = \mathbf{H} \mathbf{s} + \mu \mathbf{s} + \mathbf{g} = 0 \quad \Rightarrow \quad \mathbf{s} = -(\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}$$

and $\mathbf{s}^T \mathbf{s} = \Delta^2$. Remember that if \mathbf{H} is SPD then $\mathbf{H} + \mu \mathbf{I}$ is SPD for all $\mu \geq 0$. Moreover the inverse of an SPD matrix is SPD. From

$$\mathbf{g}^T \mathbf{s} = -\mathbf{g}^T (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g} < 0 \quad \text{for all } \mu \geq 0$$

follows that $\mathbf{s}(\mu)$ is a descent direction for all $\mu \geq 0$.



Proof.

(2/2).

To prove the uniqueness expand the gradient \mathbf{g} with the eigenvectors of \mathbf{H}

$$\mathbf{g} = \sum_{i=1}^n \alpha_i \mathbf{u}_i$$

\mathbf{H} is SPD so that \mathbf{u}_i can be chosen orthonormal. It follows

$$(\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g} = (\mathbf{H} + \mu \mathbf{I})^{-1} \sum_{i=1}^n \alpha_i \mathbf{u}_i = \sum_{i=1}^n \frac{\alpha_i}{\lambda_i + \mu} \mathbf{u}_i$$

$$\|(\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}\|^2 = \sum_{i=1}^n \frac{\alpha_i^2}{(\lambda_i + \mu)^2}$$

and $\|(\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}\|$ is a monotonically decreasing function of μ . \square

Lemma

Consider the following constrained quadratic problem where $\mathbf{H} \in \mathbb{R}^{n \times n}$ is **symmetric** with $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ its eigenvalues.

$$\arg \min_{\|\mathbf{s}\| \leq \Delta} f(\mathbf{s}), \quad f(\mathbf{s}) = f_0 + \mathbf{g}^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H} \mathbf{s},$$

Then the following curve

$$\mathbf{s}(\mu) \doteq -(\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g},$$

for any $\mu > -\lambda_1$ defines a descent direction for $f(\mathbf{s})$ and $\mathbf{H} + \mu \mathbf{I}$ is positive definite. Moreover

- or $\|\mathbf{s}(0)\| < \Delta$ with $\mathbf{g}^T \mathbf{s}(0) < 0$ and $\mathbf{s}(0)$ is a **local minima** of the problem;
- or there exists a μ_* $> -\lambda_n$ such that $\|\mathbf{s}(\mu_*)\| = \Delta$ and $\mathbf{s}(\mu_*)$ is a **local minima** of the problem;

Remark

As a consequence of the previous Lemma we have:

- as the radius of the trust region becomes smaller as the scalar μ becomes larger. This means that the search direction become more and more oriented toward the gradient direction.
- as the radius of the trust region becomes larger as the scalar μ becomes smaller. This means that the search direction become more and more oriented toward the Newton direction.

Thus a trust region technique not only change the size of the step-length but also its direction. This results in a more robust numerical technique. The price to pay is that the solution of the minimization is more costly than the inexact line search.

but what happen when \mathbf{H} is not positive definite ?

Proof.

(1/6).

Consider the Lagrangian

$$\begin{aligned} \mathcal{L}(\mathbf{s}, \mu, \epsilon) &= f_0 + \mathbf{g}^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H} \mathbf{s} \\ &\quad + \frac{1}{2} \mu (\mathbf{s}^T \mathbf{s} + \epsilon^2 - \Delta^2) + \omega (\mathbf{g}^T \mathbf{s} + \delta^2), \end{aligned}$$

where

$$\mathbf{s}^T \mathbf{s} + \epsilon^2 - \Delta^2$$

is the constraint $\|\mathbf{s}\| \leq \Delta^2$ on the length of the step and

$$\mathbf{g}^T \mathbf{s} + \delta^2$$

is the constraint $\mathbf{g}^T \mathbf{s} \leq 0$ on the step that must be descent

Proof.

(2/6).

Then we must solve the nonlinear system:

$$\partial_s \mathcal{L}(s, \mu, \omega, \epsilon, \delta) = \mathbf{H}s + \mu s + (1 + \omega)g = 0$$

$$2\partial_\mu \mathcal{L}(s, \mu, \omega, \epsilon, \delta) = s^T s + \epsilon^2 - \Delta^2 = 0$$

$$\partial_\omega \mathcal{L}(s, \mu, \omega, \epsilon, \delta) = g^T s + \delta^2 = 0$$

$$\partial_\epsilon \mathcal{L}(s, \mu, \omega, \epsilon, \delta) = \mu\epsilon = 0$$

$$\partial_\delta \mathcal{L}(s, \mu, \omega, \epsilon, \delta) = 2\delta\omega = 0$$

from the first equation we have:

$$s = \frac{-1}{1 + \omega}(\mathbf{H} + \mu\mathbf{I})^{-1}g$$

and if we want a descent direction $g^T s < 0$ which imply $\omega = 0$.

Proof.

(3/6).

So that we must solve the reduced non linear system

$$s = -(\mathbf{H} + \mu\mathbf{I})^{-1}g$$

$$s^T s + \epsilon^2 - \Delta^2 = 0$$

$$g^T s = -\delta^2$$

$$\mu\epsilon = 0$$

combining the first and third equation we have

$$g^T(\mathbf{H} + \mu\mathbf{I})^{-1}g = \delta^2 \geq 0$$



Proof.

(4/6).

If $\epsilon \neq 0$ then we must have $\mu = 0$ and

$$\|-\mathbf{H}^{-1}g\| = \|s\| \leq \Delta$$

with $g^T \mathbf{H}^{-1}g \geq 0$. If $\epsilon = 0$ then we must have

$$\|-(\mathbf{H} + \mu\mathbf{I})^{-1}g\| = \|s\| = \Delta$$

with $g^T(\mathbf{H} + \mu\mathbf{I})^{-1}g \geq 0$. Expand $g = \sum_{i=1}^n \alpha_i u_i$ with an orthonormal base of eigenvectors of \mathbf{H} it follows

$$\|(\mathbf{H} + \mu\mathbf{I})^{-1}g\| = \sum_{i=1}^n \frac{\alpha_i^2}{(\lambda_i + \mu)^2}$$

$$g^T(\mathbf{H} + \mu\mathbf{I})^{-1}g = \sum_{i=1}^n \frac{\alpha_i^2}{\lambda_i + \mu}$$

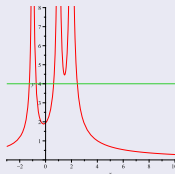


Proof.

(5/6).

$\|(\mathbf{H} + \mu\mathbf{I})^{-1}g\|$ is a monotonically decreasing function of μ for $\mu > -\lambda_k$ where k is the first index such that $\alpha_k \neq 0$. For example

$$\|(\mathbf{H} + \mu\mathbf{I})^{-1}g\| = (\mu + 1)^{-2} + 2(\mu - 1)^{-2} + 3(\mu - 2)^{-2}$$



Proof.

(6/6).

Thus, or

$$\| -H^{-1}g \| = \|s\| \leq \Delta \text{ with } g^T H^{-1}g > 0.$$

or let be k the first index such that $\alpha_k \neq 0$, we can find a $\mu > -\lambda_k$ such that

$$\| -(H + \mu I)^{-1}g \| = \sum_{i=k}^n \frac{\alpha_i^2}{(\lambda_i + \mu)^2} = \Delta$$

$$g(H + \mu I)^{-1}g = \sum_{i=k}^n \frac{\alpha_i^2}{\lambda_i + \mu} > 0$$



Outline

- 1 The Trust Region method
- 2 **Convergence analysis**
- 3 The exact solution of trust region step
- 4 The dogleg trust region step
- 5 The double dogleg trust region step
- 6 Two dimensional subspace minimization



Algorithm (Basic trust region algorithm)

x_0 assigned; Δ_0 assigned; $k \leftarrow 0$;

while $\| \nabla f(x_k) \| \neq 0$ **do**

$m_k(s) = f(x_k) + \nabla f(x_k)s + \frac{1}{2}s^T H_k s$; — *setup the model*

$s_k \leftarrow \arg \min_{\|s\| \leq \Delta_k} m_k(s)$; — *compute the step*

$x_{k+1} \leftarrow x_k + s_k$;

$\rho_k \leftarrow (f(x_k) - f(x_{k+1})) / (m_k(0) - m_k(s_k))$;

— *check the reduction*

if $\rho_k > \beta_2$ **then**

$\Delta_{k+1} \leftarrow 2\Delta_k$; — *very successful*

else if $\rho_k > \beta_1$ **then**

$\Delta_{k+1} \leftarrow \Delta_k$; — *successful*

else

$\Delta_{k+1} \leftarrow \Delta_k/2$; $x_{k+1} \leftarrow x_k$; — *failure*

end if

$k \leftarrow k + 1$;

end while



Cauchy point

Definition

Consider the quadratic

$$m(s) = f_0 + g^T s + \frac{1}{2}s^T H s$$

and the minimization problem

$$s^c(\Delta) = \arg \min_{s \in \{-tg \mid t \geq 0, \| -tg \| \leq \Delta\}} m(s)$$

The point $s^c(\Delta)$ is called
Cauchy point or step.



Estimate the length of the Cauchy step

Lemma

For the Cauchy step the following characterization is valid:

$$s^c(\Delta) = -\tau(\Delta) \frac{\mathbf{g}}{\|\mathbf{g}\|}$$

$$\tau(\Delta) = \begin{cases} \Delta & \text{if } \mathbf{g}^T \mathbf{H} \mathbf{g} \leq 0 \\ \min \left\{ \frac{\|\mathbf{g}\|^3}{\mathbf{g}^T \mathbf{H} \mathbf{g}}, \Delta \right\} & \text{if } \mathbf{g}^T \mathbf{H} \mathbf{g} > 0 \end{cases}$$

Moreover

$$\tau(\Delta) \geq \min \left\{ \frac{\|\mathbf{g}\|}{\varrho(\mathbf{H})}, \Delta \right\}$$

where $\varrho(\mathbf{H})$ is the spectral radius of \mathbf{H}



Proof.

Consider an orthonormal base of eigenvectors for \mathbf{H} and write \mathbf{g} if this coordinate:

$$\mathbf{g} = \sum_{i=1}^n \alpha_i \mathbf{u}_i$$

so that

$$\frac{\mathbf{g}^T \mathbf{H} \mathbf{g}}{\mathbf{g}^T \mathbf{g}} = \frac{\sum_{i=1}^n \lambda_i \alpha_i^2}{\sum_{i=1}^n \alpha_i^2} \leq \frac{\sum_{i=1}^n |\lambda_i| \alpha_i^2}{\sum_{i=1}^n \alpha_i^2} \leq \varrho(\mathbf{H})$$

and finally

$$\frac{\|\mathbf{g}\|^3}{\mathbf{g}^T \mathbf{H} \mathbf{g}} = \|\mathbf{g}\| \frac{\mathbf{g}^T \mathbf{g}}{\mathbf{g}^T \mathbf{H} \mathbf{g}} \geq \frac{\|\mathbf{g}\|}{\varrho(\mathbf{H})}$$



Proof.

Consider

$$h(t) = m(-t\mathbf{g}/\|\mathbf{g}\|) = f_0 - t\|\mathbf{g}\| + \frac{t^2 \mathbf{g}^T \mathbf{H} \mathbf{g}}{2\|\mathbf{g}\|^2}$$

$h(t)$ is a parabola in t and if $\mathbf{g}^T \mathbf{H} \mathbf{g} \leq 0$ then the parabola decrease monotonically for $t \geq 0$. In this case the point is on the boundary of the trust region ($t = \Delta$).

If $\mathbf{g}^T \mathbf{H} \mathbf{g} > 0$ the parabola is decreasing until the global minima at

$$t = \frac{\|\mathbf{g}\|^3}{\mathbf{g}^T \mathbf{H} \mathbf{g}}$$

Otherwise we separate the case if the minimum of the parabola is inside or outside the trust region. (cont.)



Estimate the reduction obtained by the Cauchy step

In the convergence analysis is important to obtain estimation of the reduction of the function to be minimized.

A first step in this direction is the estimation of the reduction of the model quadratic function.

Lemma

Consider the quadratic

$$m(\mathbf{s}) = f_0 + \mathbf{g}^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H} \mathbf{s}$$

then for the Cauchy step we have:

$$m(\mathbf{0}) - m(s^c(\Delta)) \geq \frac{1}{2} \|\mathbf{g}\| \min \left\{ \Delta, \frac{\|\mathbf{g}\|}{\varrho(\mathbf{H})} \right\}$$



Proof.

Compute

$$m(\mathbf{0}) - m(\mathbf{s}^c(\Delta)) = \tau(\Delta) \|g\| - \frac{\tau(\Delta)^2}{2 \|g\|^2} g^T H g$$

If $g^T H g \leq 0$ for lemma on slide N.25 we have $\tau(\Delta) = \Delta$

$$\begin{aligned} m(\mathbf{0}) - m(\mathbf{s}^c(\Delta)) &= \Delta \|g\| - \frac{\Delta^2}{2 \|g\|^2} g^T H g \\ &= \Delta \left(\|g\| - \frac{\Delta g^T H g}{2 \|g\|^2} \right) \\ &\geq \Delta \|g\| \end{aligned}$$

(cont.)

Proof.

If $g^T H g > 0$ we have

$$\tau(\Delta) = \min \left\{ \|g\|^3 / (g^T H g), \Delta \right\}$$

and

$$\begin{aligned} m(\mathbf{0}) - m(\mathbf{s}^c(\Delta)) &= \tau(\Delta) \left(\|g\| - \frac{1}{2} \min \left\{ \|g\|, \Delta \frac{g^T H g}{\|g\|^2} \right\} \right) \\ &\geq \tau(\Delta) \left(\|g\| - \frac{1}{2} \|g\| \right) \\ &\geq \tau(\Delta) \frac{1}{2} \|g\| \end{aligned}$$

so that in general $m(\mathbf{0}) - m(\mathbf{s}^c(\Delta)) \geq \tau(\Delta) \frac{1}{2} \|g\|$. \square

- A successful step in trust region algorithm imply that the ratio

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(\mathbf{0}) - m_k(\mathbf{s}_k)}$$

is greater than a constant $\beta_1 > 0$.

- Any reasonable step in a trust region algorithm should be no (asymptotically) worse than a Cauchy step. So we require

$$m_k(\mathbf{0}) - m_k(\mathbf{s}_k) \geq \eta [m_k(\mathbf{0}) - m_k(\mathbf{s}^c(\Delta_k))]$$

for a constant $\eta > 0$.

- Using lemma on slide N.28

$$\begin{aligned} f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k) &= \rho_k (m_k(\mathbf{0}) - m_k(\mathbf{s}_k)) \\ &\geq \rho_k \eta [m_k(\mathbf{0}) - m_k(\mathbf{s}^c(\Delta_k))] \\ &\geq \frac{\eta \beta_1}{2} \|\nabla f(\mathbf{x}_k)\| \min \left\{ \Delta_k, \frac{\|\nabla f(\mathbf{x}_k)\|}{\varrho(H_k)} \right\} \end{aligned}$$

- Thus any reasonable trust region numerical scheme satisfy

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq \frac{\eta \beta_1}{2} \|\nabla f(\mathbf{x}_k)\| \min \left\{ \Delta_k, \frac{\|\nabla f(\mathbf{x}_k)\|}{\varrho(H_k)} \right\}$$

for any successful step (for unsuccessful step $\mathbf{x}_{k+1} = \mathbf{x}_k$).

- Let \mathcal{S} the index set of successful step, then

$$\begin{aligned} f(\mathbf{x}_0) - \lim_{k \in \mathcal{S}} f(\mathbf{x}_k) &\geq \\ &\frac{\eta \beta_1}{2} \sum_{k \in \mathcal{S}} \|\nabla f(\mathbf{x}_k)\| \min \left\{ \Delta_k, \frac{\|\nabla f(\mathbf{x}_k)\|}{\varrho(H_k)} \right\} \end{aligned}$$

thus we can use arguments similar to Zoutendijk theorem to prove convergence.

- To complete the argument we must set conditions that guarantees that $\Delta_k \neq 0$ as $k \rightarrow \infty$ and that cardinality of \mathcal{S} is not finite.

Technical assumption

The following assumptions permits to characterize a class of convergent trust region algorithm.

Assumption

For any **successful** step in trust region algorithm, the ratio

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(\mathbf{0}) - m_k(\mathbf{s}_k)}$$

is greater than a constant $\beta_1 > 0$.

Assumption

For any step in trust region algorithm, the model reduction for a constant $\eta > 0$ satisfy the inequality:

$$m_k(\mathbf{0}) - m_k(\mathbf{s}_k) \geq \eta [m_k(\mathbf{0}) - m_k(\mathbf{s}^c(\Delta_k))]$$



The following lemma permits to estimate the reduction ratio ρ_k and conclude that there exists a positive trust ray Δ_k for which the step is accepted!

Lemma

Let be $f \in C^1(\mathbb{R}^n)$ with Lipschitz continuous gradient

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq \gamma \|\mathbf{x} - \mathbf{y}\|$$

and apply basic trust region algorithm of slide N.23 with assumption of slide N.33 then we have

$$\Delta_k \geq \frac{(1 - \beta_2)\eta \|\nabla f(\mathbf{x}_k)\|}{2(\varrho(\mathbf{H}_k) + \gamma)}$$

for any accepted step.



Proof.

By using Taylor's theorem

$$\begin{aligned} f(\mathbf{x}_k + \mathbf{s}_k) &= f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k) \mathbf{s}_k \\ &\quad + \int_0^1 [\nabla f(\mathbf{x}_k + t\mathbf{s}_k) - \nabla f(\mathbf{x}_k)] \mathbf{s}_k dt \end{aligned}$$

so that

$$\begin{aligned} m_k(\mathbf{s}_k) - f(\mathbf{x}_k + \mathbf{s}_k) &= (\mathbf{s}_k^T \mathbf{H}_k \mathbf{s}_k) / 2 \\ &\quad - \int_0^1 [\nabla f(\mathbf{x}_k + t\mathbf{s}_k) - \nabla f(\mathbf{x}_k)] \mathbf{s}_k dt \end{aligned}$$

and

$$|m_k(\mathbf{s}_k) - f(\mathbf{x}_k + \mathbf{s}_k)| \leq \frac{\mathbf{s}_k^T \mathbf{H}_k \mathbf{s}_k}{2} + \frac{\gamma}{2} \|\mathbf{s}_k\|^2 \leq \frac{\varrho(\mathbf{H}_k) + \gamma}{2} \|\mathbf{s}_k\|^2$$

(cont.)



Proof.

using these inequalities we can estimate the ratio

$$\begin{aligned} \left| \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(\mathbf{0}) - m_k(\mathbf{s}_k)} - 1 \right| &= \frac{|m_k(\mathbf{s}_k) - f(\mathbf{x}_k + \mathbf{s}_k)|}{|m_k(\mathbf{0}) - m_k(\mathbf{s}_k)|} \\ &\leq \frac{1}{2\eta} \frac{(\varrho(\mathbf{H}_k) + \gamma) \|\mathbf{s}_k\|^2}{|m_k(\mathbf{0}) - m_k(\mathbf{s}^c(\Delta))|} \\ &\leq \frac{(\varrho(\mathbf{H}_k) + \gamma) \Delta^2}{\eta \|\nabla f(\mathbf{x}_k)\| \min \left\{ \Delta, \frac{\|\nabla f(\mathbf{x}_k)\|}{\varrho(\mathbf{H}_k)} \right\}} \end{aligned}$$

(cont.)



Proof.

If $\Delta \leq \|\nabla f(\mathbf{x}_k)\| / \varrho(\mathbf{H}_k)$ we obtain

$$|\rho_k - 1| \leq \frac{(\varrho(\mathbf{H}_k) + \gamma)\Delta}{\eta \|\nabla f(\mathbf{x}_k)\|}$$

so that when $\Delta_k \leq \Delta$:

$$\Delta = \frac{(1 - \beta_2)\eta \|\nabla f(\mathbf{x}_k)\|}{(\varrho(\mathbf{H}_k) + \gamma)}$$

than $\rho_k \geq 1 - \beta_2$ and the step is accepted \square

Corollary

Apply basic trust region algorithm of slide N.23 with assumption of slide N.33 to $f \in \mathcal{C}^1(\mathbb{R}^n)$ with Lipschitz continuous gradient

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq \gamma \|\mathbf{x} - \mathbf{y}\|$$

then we have

$$f(\mathbf{x}_0) - \lim_{k \in \mathcal{S}} f(\mathbf{x}_k) \geq \frac{\eta^2 \beta_1 (1 - \beta_2)}{4} \sum_{k \in \mathcal{S}} \frac{\|\nabla f(\mathbf{x}_k)\|^2}{\varrho(\mathbf{H}_k) + \gamma}$$

moreover if $\varrho(\mathbf{H}_k) \leq C$ for all k we have

$$f(\mathbf{x}_0) - \lim_{k \in \mathcal{S}} f(\mathbf{x}_k) \geq \frac{\eta^2 \beta_1 (1 - \beta_2)}{4(C + \gamma)} \sum_{k \in \mathcal{S}} \|\nabla f(\mathbf{x}_k)\|^2$$

Convergence theorem

Theorem (Convergence to stationary points)

Apply basic trust region algorithm of slide N.23 with assumption of slide N.33 to $f \in \mathcal{C}^1(\mathbb{R}^n)$ with Lipschitz continuous gradient

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq \gamma \|\mathbf{x} - \mathbf{y}\|$$

if the set

$$\mathcal{K} = \{\mathbf{x} \mid f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$$

is compact and $\varrho(\mathbf{H}_k) \leq C$ for all k we have

$$\lim_{k \rightarrow \infty} \nabla f(\mathbf{x}_k) = \mathbf{0}$$

Proof.

A trivial application of previous corollary. \square

Convergence theorem

Theorem (Convergence to minima)

Apply basic trust region algorithm of slide N.23 with assumption of slide N.33 to $f \in \mathcal{C}^2(\mathbb{R}^n)$. If $\mathbf{H}_k = \nabla^2 f(\mathbf{x}_k)$ and the set

$$\mathcal{K} = \{\mathbf{x} \mid f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$$

is compact then:

- 1 Or the iteration terminate at \mathbf{x}_k which satisfy second order necessary condition.
- 2 Or the limit point $\mathbf{x}_* = \lim_{k \rightarrow \infty} \mathbf{x}_k$ satisfy second order necessary condition.



J. J. Moré, D.C.Sorensen

Computing a Trust Region Step

SIAM J. Sci. Stat. Comput. 4, No. 3, 1983

Solving the constrained minimization problem

As for the line-search problem we have many alternative for solving the constrained minimization problem:

- We can solve **accurately** the constrained minimization problem. For example by an iterative method.
- We can **approximate** the solution of the constrained minimization problem.

as for the line search the accurate solution of the constrained minimization problem is not paying while a good cheap approximations is normally better performing.



Outline

- 1 The Trust Region method
- 2 Convergence analysis
- 3 The exact solution of trust region step
- 4 The dogleg trust region step
- 5 The double dogleg trust region step
- 6 Two dimensional subspace minimization



The Newton approach

(1/7)

- Consider the Lagrangian

$$\mathcal{L}(s, \mu) = a + \mathbf{g}^T s + \frac{1}{2} s^T \mathbf{H} s + \frac{1}{2} \mu (s^T s - \Delta^2),$$

where $a = f(\mathbf{x})$ and $\mathbf{g} = \nabla f(\mathbf{x})^T$.

- Then we can try to solve the nonlinear system

$$\frac{\partial \mathcal{L}}{\partial (s, \mu)}(s, \mu) = \begin{pmatrix} \mathbf{H} s + \mu s + \mathbf{g} \\ (s^T s - \Delta^2)/2 \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ 0 \end{pmatrix}$$

- Using Newton method we have

$$\begin{pmatrix} s_{k+1} \\ \mu_{k+1} \end{pmatrix} = \begin{pmatrix} s_k \\ \mu_k \end{pmatrix} - \begin{pmatrix} \mathbf{H} + \mu \mathbf{I} & s \\ s^T & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{H} s_k + \mu_k s_k + \mathbf{g} \\ (s_k^T s_k - \Delta^2)/2 \end{pmatrix}$$



The Newton approach

(2/7)

Lemma

let be $s(\mu)$ the solution of $(\mathbf{H} + \mu \mathbf{I})s(\mu) = -\mathbf{g}$ than we have

$$s'(\mu) = -(\mathbf{H} + \mu \mathbf{I})^{-1} s(\mu) \quad \text{and} \quad s''(\mu) = 2(\mathbf{H} + \mu \mathbf{I})^{-2} s(\mu)$$

Proof.

It enough to differentiate the relation

$$\mathbf{H} s(\mu) + \mu s(\mu) = \mathbf{g}$$

two times:

$$\mathbf{H} s'(\mu) + \mu s'(\mu) + s(\mu) = \mathbf{0}$$

$$\mathbf{H} s''(\mu) + \mu s''(\mu) + 2s'(\mu) = \mathbf{0}$$



The Newton approach

(3/7)

- A better approach to compute μ is given by solving $\Phi(\mu) = 0$ where

$$\Phi(\mu) = \|\mathbf{s}(\mu)\| - \Delta, \quad \text{and} \quad \mathbf{s}(\mu) = -(\mathbf{H} + \mu\mathbf{I})^{-1}\mathbf{g}$$

- To build Newton method we need to evaluate

$$\Phi'(\mu) = \frac{\mathbf{s}(\mu)^T \mathbf{s}'(\mu)}{\|\mathbf{s}(\mu)\|}, \quad \mathbf{s}'(\mu) = -(\mathbf{H} + \mu\mathbf{I})^{-1}\mathbf{s}(\mu)$$

- Putting all in a Newton step we obtain

$$\mu_{k+1} = \mu_k + \frac{\Delta - \|\mathbf{s}(\mu_k)\|}{\mathbf{s}(\mu_k)^T \mathbf{s}'(\mu_k)} \|\mathbf{s}(\mu_k)\|$$



The Newton approach

(5/7)

Lemma

If \mathbf{H} is SPD for all $\mu > 0$ we have:

$$\Phi'(\mu) < 0 \quad \text{and} \quad \Phi''(\mu) > 0$$

Proof.

If $\mu > 0$ then $\mathbf{s}(\mu) \neq \mathbf{0}$. Evaluating $\Phi'(\mu)$ and using lemma of slide N.44 we have

$$\|\mathbf{s}(\mu)\| \Phi'(\mu) = \mathbf{s}(\mu)^T \mathbf{s}'(\mu) = -\mathbf{s}(\mu)^T (\mathbf{H} + \mu\mathbf{I})^{-1} \mathbf{s}(\mu) < 0$$

Evaluating $\Phi''(\mu)$ and using lemma of slide N.44 we have

$$\Phi''(\mu) = \frac{\mathbf{s}'(\mu)^T \mathbf{s}'(\mu) + \mathbf{s}(\mu)^T \mathbf{s}''(\mu)}{\|\mathbf{s}(\mu)\|} - \frac{(\mathbf{s}(\mu)^T \mathbf{s}'(\mu))^2}{\|\mathbf{s}(\mu)\|^3}$$

(cont.)



The Newton approach

(4/7)

- Newton step can be reorganized as follows

$$\mathbf{a} = (\mathbf{H} + \mu_k \mathbf{I})^{-1} \mathbf{g}$$

$$\mathbf{b} = (\mathbf{H} + \mu_k \mathbf{I})^{-1} \mathbf{a}$$

$$\beta = \|\mathbf{a}\|$$

$$\mu_{k+1} = \mu_k + \beta \frac{\beta - \Delta}{\mathbf{a}^T \mathbf{b}}$$

- Thus Newton step require **two** linear system solution per step. However the coefficient matrix is the same so that **only one** LU factorization, thus the cost per step is essentially due to the LU factorization.



The Newton approach

(6/7)

Proof.

Using Cauchy-Schwartz inequality

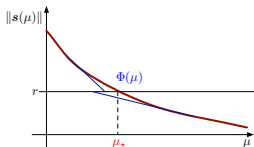
$$\begin{aligned} \Phi''(\mu) &\geq \frac{\mathbf{s}'(\mu)^T \mathbf{s}'(\mu) + \mathbf{s}(\mu)^T \mathbf{s}''(\mu)}{\|\mathbf{s}(\mu)\|} - \frac{\|\mathbf{s}(\mu)\|^2 \|\mathbf{s}'(\mu)\|^2}{\|\mathbf{s}(\mu)\|^3} \\ &= \frac{\mathbf{s}(\mu)^T \mathbf{s}''(\mu)}{\|\mathbf{s}(\mu)\|} \\ &= 2 \frac{\mathbf{s}(\mu)^T (\mathbf{H} + \mu\mathbf{I})^{-2} \mathbf{s}(\mu)}{\|\mathbf{s}(\mu)\|} > 0 \end{aligned}$$



The Newton approach

(7/7)

- From $\Phi''(\mu) > 0$ we have that Newton is monotonically convergent and steps underestimates μ .



- The model consists of two parameter α_k and β_k . To set this parameter we can impose

$$m_k(\mu_k) = \frac{\alpha_k}{\beta_k + \mu_k} - \Delta = \Phi(\mu_k)$$

$$m'_k(\mu_k) = -\frac{\alpha_k}{(\beta_k + \mu_k)^2} = \Phi'(\mu_k)$$

- solving for α_k and β_k we have

$$\alpha_k = -\frac{(\Phi(\mu_k) + \Delta)^2}{\Phi'(\mu_k)} \quad \beta_k = -\frac{\Phi(\mu_k) + \Delta}{\Phi'(\mu_k)} - \mu_k$$

where

$$\Phi(\mu_k) = \|s(\mu_k)\| - \Delta \quad \Phi'(\mu_k) = -\frac{s(\mu_k)^T (\mathbf{H} + \mu_k \mathbf{I})^{-1} s(\mu_k)}{\|s(\mu_k)\|^2}$$

- Having α_k and β_k it is possible to solve $m_k(\mu) = 0$ obtaining

$$\mu_{k+1} = \frac{\alpha_k}{\Delta} - \beta_k$$



- If we develop the vector g with the orthonormal bases given by the eigenvectors of \mathbf{H} we have

$$g = \sum_{i=1}^n \alpha_i \mathbf{u}_i$$

- Using this expression to evaluate $s(\mu)$ we have

$$s(\mu) = -(\mathbf{H} + \mu \mathbf{I})^{-1} g = \sum_{i=1}^n \frac{\alpha_i}{\mu + \lambda_i} \mathbf{u}_i$$

$$\|s(\mu)\| = \left(\sum_{i=1}^n \frac{\alpha_i^2}{(\mu + \lambda_i)^2} \right)^{1/2}$$

- This expression suggest to use as a model for $\Phi(\mu)$ the following expression

$$m_k(\mu) = \frac{\alpha_k}{\beta_k + \mu} - \Delta$$



- Substituting α_k and β_k the step become

$$\mu_{k+1} = \mu_k - \frac{\Phi(\mu_k)}{\Phi'(\mu_k)} - \frac{\Phi(\mu_k)^2}{\Phi'(\mu_k)\Delta} = \mu_k - \frac{\Phi(\mu_k)}{\Phi'(\mu_k)} \left(1 + \frac{\Phi(\mu_k)}{\Delta} \right)$$

- Comparing with the Newton step

$$\mu_{k+1} = \mu_k - \frac{\Phi(\mu_k)}{\Phi'(\mu_k)}$$

we see that this method perform larger step by a factor $1 + \Phi(\mu_k)\Delta^{-1}$.

- Notice that $1 + \Phi(\mu_k)\Delta^{-1}$ converge to 1 as $\mu_k \rightarrow \mu_*$. So that this iteration become the Newton iteration as μ_k becomes near the solution.



Algorithm (Exact trust region algorithm)

$$\text{exact_trust_region}(\Delta, g, H)$$

$$\mu \leftarrow 0;$$

$$s \leftarrow H^{-1}g;$$

while $\|s\| - \Delta > \epsilon$ and $\mu \geq 0$ **do**

— *compute the model*

$$s' \leftarrow -(H + \mu I)^{-1}s;$$

$$\Phi \leftarrow \|s\| - \Delta;$$

$$\Phi' \leftarrow (s^T s') / \|s\|$$

— *update μ and s*

$$\mu \leftarrow \mu - \frac{\Phi}{\Phi'} \frac{\|s\|}{\Delta};$$

$$s \leftarrow -(H + \mu I)^{-1}g;$$

end while

if $\mu < 0$ **then**

$$s \leftarrow -H^{-1}g;$$

end if



Outline

- 1 The Trust Region method
- 2 Convergence analysis
- 3 The exact solution of trust region step
- 4 **The dogleg trust region step**
- 5 The double dogleg trust region step
- 6 Two dimensional subspace minimization



The DogLeg approach

(1/3)

- The computation of the μ such that $\|s(\mu)\| = \Delta$ of the **exact** trust region computation can be very expensive.
- An alternative was proposed by Powell:



M.J.D. Powell

A hybrid method for nonlinear equations
in: Numerical Methods for Nonlinear Algebraic Equations
ed. Ph. Rabinowitz, Gordon and Breach, pages 87-114,
1970.

where instead of computing exactly the curve $s(\mu)$ a piecewise linear approximation $s_{dl}(\mu)$ is used in computation.

- This approximation also permits to solve $\|s_{dl}(\mu)\| = \Delta$ explicitly.



The DogLeg approach

(2/3)

- Form the definition of $s(\mu) = -(H + \mu I)^{-1}g$ and the relation $s'(\mu) = (H + \mu I)^{-2}g$ it follows

$$s(0) = -H^{-1}g, \quad \lim_{\mu \rightarrow \infty} \mu^2 s'(\mu) = -g$$

i.e. the curve start from the Newton step and reduce to zero in the direction opposite to the gradient step.

- The direction $-g$ is a descent direction, so that a first piece of the piecewise approximation should be a straight line from x to the minimum of $m_k(-\lambda g)$. The minimum λ_* is found at

$$\lambda_* = \frac{\|g\|^2}{g^T H g}$$

- Having reached the minimum if the $-g$ direction we can now go to the point $x + s(0) = x - H^{-1}g$ with another straight line.



The DogLeg approach

(3/3)

- We denote by

$$s_g = -g \frac{\|g\|^2}{g^T H g}, \quad s_n = -H^{-1} g$$

respectively the step due to the unconstrained minimization in the gradient direction and in the Newton direction.

- The piecewise linear curve connecting $x + s_n$, $x + s_g$ and x is the **DogLeg** curve¹ $x_{dl}(\mu) = x + s_{dl}(\mu)$ where

$$s_{dl}(\mu) = \begin{cases} \mu s_g + (1 - \mu) s_n & \text{for } \mu \in [0, 1] \\ (2 - \mu) s_g & \text{for } \mu \in [1, 2] \end{cases}$$

¹notice that $s(\mu)$ is parametrized in the interval $[0, \infty]$ while $s_{dl}(\mu)$ is parametrized in the interval $[0, 2]$



By using Kantorovich we can prove:

Lemma

We denote by

$$s_g = -g \frac{\|g\|^2}{g^T H g}, \quad s_n = -H^{-1} g, \quad \gamma_* = \frac{\|s_g\|^2}{s_n^T s_g}$$

then $\gamma_* \leq 1$, moreover if s_n is not parallel to s_g then $\gamma_* < 1$.

Proof.

Using

$$s_n^T s_g = \|g\|^2 \frac{g^T H^{-1} g}{g^T H g} \quad \text{and} \quad s_g^2 = \frac{\|g\|^6}{(g^T H g)^2}$$

we have $\gamma_* = \|g\|^4 / [(g^T H g)(g^T H^{-1} g)]$ and using Kantorovich inequality the lemma is proved. \square



Lemma (Kantorovich)

Let $A \in \mathbb{R}^{n \times n}$ an SPD matrix then the following inequality is valid

$$1 \leq \frac{(x^T A x)(x^T A^{-1} x)}{(x^T x)^2} \leq \frac{(M + m)^2}{4 M m}$$

for all $x \neq 0$. Where $m = \lambda_1$ is the smallest eigenvalue of A and $M = \lambda_n$ is the biggest eigenvalue of A .

this lemma can be improved a little bit for the first inequality

Lemma (Kantorovich (bis))

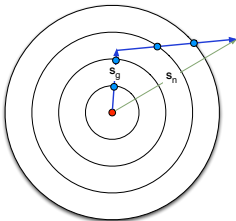
Let $A \in \mathbb{R}^{n \times n}$ an SPD matrix then the following inequality is valid

$$1 < \frac{(x^T A x)(x^T A^{-1} x)}{(x^T x)^2}$$

for all $x \neq 0$ and x not an eigenvector of A .



the Dogleg piecewise curve



Lemma

Consider the *dogleg* curve connecting $x + s_n$, $x + s_g$ and x . The curve can be expressed as $x_{dl}(\mu) = x + s_{dl}(\mu)$ where

$$s_{dl}(\mu) = \begin{cases} \mu s_g + (1 - \mu) s_n & \text{for } \mu \in [0, 1] \\ (2 - \mu) s_g & \text{for } \mu \in [1, 2] \end{cases}$$

for this curve if s_g is not parallel to s_n we have that the function

$$d(\mu) = \|x_{dl}(\mu) - x\| = \|s_{dl}(\mu)\|$$

is strictly monotone decreasing, moreover the direction $s_{dl}(\mu)$ is a descent direction for all $\mu \in [0, 2]$.



Proof.

(2/4).

Notice that $(2\mu - 4) < 0$ for $\mu \in [1, 2]$ so that we need only to check that

$$2\mu(s_g^2 + s_n^2 - 2s_g^T s_n) - 2s_n^2 + 2s_g^T s_n < 0 \quad \text{for } \mu \in [0, 1]$$

moreover

$$s_g^2 + s_n^2 - 2s_g^T s_n = \|s_g - s_n\|^2 \geq 0$$

Then it is enough to check the inequality for $\mu = 1$

$$2(s_g^2 + s_n^2 - 2s_g^T s_n) - 2s_n^2 + 2s_g^T s_n = 2s_g^2 - 2s_g^T s_n$$

i.e. we must check $s_g^2 - s_g^T s_n < 0$.



Proof.

(1/4).

In order to have a unique solution to the problem $\|s_{dl}(\mu)\| = \Delta$ we must have that $\|s_{dl}(\mu)\|$ is a monotone decreasing function:

$$\|s_{dl}(\mu)\|^2 = \begin{cases} \mu^2 s_g^2 + (1 - \mu)^2 s_n^2 + 2\mu(1 - \mu) s_g^T s_n & \mu \in [0, 1] \\ (2 - \mu)^2 s_g^2 & \mu \in [1, 2] \end{cases}$$

To check monotonicity we take first derivative

$$\begin{aligned} \frac{d}{d\mu} \|s_{dl}(\mu)\|^2 &= \begin{cases} 2\mu s_g^2 - 2(1 - \mu) s_n^2 + (2 - 4\mu) s_g^T s_n & \mu \in [0, 1] \\ (2\mu - 4) s_g^2 & \mu \in [1, 2] \end{cases} \\ &= \begin{cases} 2\mu(s_g^2 + s_n^2 - 2s_g^T s_n) - 2s_n^2 + 2s_g^T s_n & \mu \in [0, 1] \\ (2\mu - 4) s_g^2 & \mu \in [1, 2] \end{cases} \end{aligned}$$



Proof.

(3/4).

By using

$$\gamma_* = \frac{\|s_g\|^2}{s_n^T s_g} < 1$$

of the previous lemma

$$\begin{aligned} s_g^2 - s_g^T s_n &= \|s_g\|^2 \left(1 - \frac{s_n^T s_g}{\|s_g\|^2} \right) \\ &= \|s_g\|^2 \left(1 - \frac{1}{\gamma_*} \right) < 0 \end{aligned}$$



Proof.

(4/4).

To prove that $s_{dl}(\mu)$ is a descent direction it is enough to notice that

- for $\mu \in [0, 1]$ the direction $s_{dl}(\mu)$ is a convex combination of s_g and s_n .
- for $\mu \in [1, 2]$ the direction $s_{dl}(\mu)$ is parallel to s_g .

so that it is enough to verify that s_g and s_n are descent direction.

For s_g we have

$$s_g^T g = -\lambda_s g^T g < 0$$

For s_n we have

$$s_n^T g = -g^T H^{-1} g < 0$$

□

Solving

$$\alpha^2 \|s_g\|^2 + (1-\alpha)^2 \|s_n\|^2 + 2\alpha(1-\alpha)s_g^T s_n = \Delta^2$$

we have that if $\|s_g\| \leq \Delta \leq \|s_n\|$ the root in $[0, 1]$ is given by:

$$\Delta = \|s_g\|^2 + \|s_n\|^2 - 2s_g^T s_n = \|s_g - s_n\|^2$$

$$\alpha = \frac{\|s_n\|^2 - s_g^T s_n - \sqrt{(s_g^T s_n)^2 - \|s_g\|^2 \|s_n\|^2 + \Delta^2}}{\Delta}$$

to avoid cancellation the computation formula is the following

$$\begin{aligned} \alpha &= \frac{1}{\Delta} \frac{\|s_n\|^4 - 2s_g^T s_n \|s_n\|^2 + \|s_g\|^2 \|s_n\|^2 - \Delta^2}{\|s_n\|^2 - s_g^T s_n + \sqrt{(s_g^T s_n)^2 - \|s_g\|^2 \|s_n\|^2 + \Delta^2}} \\ &= \frac{\|s_n\|^2 - \Delta^2}{\|s_n\|^2 - s_g^T s_n + \sqrt{(s_g^T s_n)^2 - \|s_g\|^2 \|s_n\|^2 + \Delta^2} \|s_g - s_n\|^2} \end{aligned}$$

□

Using the previous Lemma we can prove

Lemma

If $\|s_{dl}(0)\| \geq \Delta$ then there is unique point $\mu \in [0, 2]$ such that $\|s_{dl}(\mu)\| = \Delta$.

Proof.

It is enough to notice that $s_{dl}(2) = 0$ and that $\|s_{dl}(\mu)\|$ is strictly monotonically descendent. □

The approximate solution of the constrained minimization can be obtained by this simple algorithm

- if $\Delta \leq \|s_g\|$ we set $s_{dl} = \Delta s_g / \|s_g\|$;
- if $\Delta \leq \|s_n\|$ we set $s_{dl} = \alpha s_g + (1-\alpha)s_n$; where α is the root in the interval $[0, 1]$ of:

$$\alpha^2 \|s_g\|^2 + (1-\alpha)^2 \|s_n\|^2 + 2\alpha(1-\alpha)s_g^T s_n = \Delta^2$$

- if $\Delta > \|s_n\|$ we set $s_{dl} = s_n$;

□

Algorithm (Computing DogLeg step)

DoglegStep(s_g, s_n, Δ);

if $\Delta \leq \|s_g\|$ **then**

$s \leftarrow \Delta \frac{s_g}{\|s_g\|}$;

else if $\Delta \geq \|s_n\|$ **then**

$s \leftarrow s_n$;

else

$a \leftarrow \|s_g\|^2$;

$b \leftarrow \|s_n\|^2$;

$c \leftarrow \|s_g - s_n\|^2$;

$d \leftarrow (a + b - c)/2$;

$\alpha \leftarrow \frac{b - d + \sqrt{d^2 - ab + \Delta^2 c}}{b - d + \sqrt{d^2 - ab + \Delta^2 c}}$;

$s \leftarrow \alpha s_g + (1-\alpha)s_n$;

end if

return s ;

□

Outline

- 1 The Trust Region method
- 2 Convergence analysis
- 3 The exact solution of trust region step
- 4 The dogleg trust region step
- 5 **The double dogleg trust region step**
- 6 Two dimensional subspace minimization

The Double DogLeg approach

- We denote by

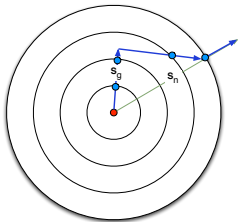
$$s_g = -g \frac{\|g\|^2}{g^T H g}, \quad s_n = -H^{-1} g, \quad \gamma_* = \frac{\|s_g\|^2}{s_g^T s_n}$$

respectively the step due to the unconstrained minimization in the gradient direction and in the Newton direction.

- The piecewise linear curve connecting $x + s_n$, $x + \gamma_* s_n$, $x + \gamma_* s_g$ and x is the **Double Dogleg curve** $x_{ddl}(\mu)$ where

$$s_{ddl}(\mu) = \begin{cases} (1 - \mu)\gamma_* s_n & \text{for } \mu \in [0, 1] \\ (\mu - 1)s_g + (2 - \mu)\gamma_* s_n & \text{for } \mu \in [1, 2] \\ (3 - \mu)s_g & \text{for } \mu \in [2, 3] \end{cases}$$

The Double Dogleg piecewise curve



Lemma

Consider the **double dogleg** curve connecting $x + s_n$, $x + \gamma_* s_n$, $x + s_g$ and x . The curve can be expressed as $x_{ddl}(\mu) = x + s_{ddl}(\mu)$ where

$$s_{ddl}(\mu) = \begin{cases} (1 - \mu)\gamma_* s_n & \text{for } \mu \in [0, 1] \\ (\mu - 1)s_g + (2 - \mu)\gamma_* s_n & \text{for } \mu \in [1, 2] \\ (3 - \mu)s_g & \text{for } \mu \in [2, 3] \end{cases}$$

for this curve if s_g is not parallel to s_n we have that the function

$$d(\mu) = \|s_{ddl}(\mu)\|$$

is strictly monotone decreasing, moreover the direction $s_{ddl}(\mu)$ is a descent direction for all $\mu \in [0, 3]$.

Proof.

(1/2).

In order to have a unique solution to the problem $\|s_{ddl}(\mu)\| = \Delta$ we must have that $\|s_{ddl}(\mu)\|$ is a monotone decreasing function. It is enough to prove for $\mu \in [1, 2]$:

$$\|s_{ddl}(1 + \alpha)\|^2 = \alpha^2 s_g^2 + (1 - \alpha)^2 \gamma_*^2 s_n^2 + 2\alpha(1 - \alpha)\gamma_* s_g^T s_n$$

To check monotonicity we take first derivative

$$\begin{aligned} \frac{d}{d\alpha} \|s_{ddl}(1 + \alpha)\|^2 &= 2\alpha s_g^2 - 2(1 - \alpha)\gamma_*^2 s_n^2 + (2 - 4\alpha)\gamma_* s_g^T s_n \\ &= 2\alpha(s_g^2 + \gamma_*^2 s_n^2 - 2\gamma_* s_g^T s_n) - 2\gamma_*^2 s_n^2 + 2\gamma_* s_g^T s_n \end{aligned}$$



Proof.

(2/2).

Notice that

$$s_g^2 + \gamma_*^2 s_n^2 - 2\gamma_* s_g^T s_n = \|s_g - \gamma_* s_n\|^2 > 0$$

because s_g and s_n are not parallel. Then it is enough to check the inequality for $\alpha = 1$

$$\begin{aligned} 2(s_g^2 + \gamma_*^2 s_n^2 - 2\gamma_* s_g^T s_n) - 2\gamma_*^2 s_n^2 + 2\gamma_* s_g^T s_n &= 2s_g^2 - 2\gamma_* s_g^T s_n \\ &= 0 \end{aligned}$$

The rest of the proof is similar as for the single dogleg step. \square



Using the previous Lemma we can prove

Lemma

If $\|s_{ddl}(0)\| \geq \Delta$ then there is unique point $\mu \in [0, 3]$ such that $\|s_{ddl}(\mu)\| = \Delta$.

The approximate solution of the constrained minimization can be obtained by this simple algorithm

- 1 if $\Delta \leq \|s_g\|$ we set $s_{ddl} = \Delta s_g / \|s_g\|$;
- 2 if $\Delta \leq \gamma_* \|s_n\|$ we set $s_{ddl} = \alpha s_g + (1 - \alpha)\gamma_* s_n$; where α is the root in the interval $[0, 1]$ of:

$$\alpha^2 \|s_g\|^2 + \gamma_*^2 (1 - \alpha)^2 \|s_n\|^2 + 2\gamma_* \alpha (1 - \alpha) s_g^T s_n = \Delta^2$$

- 3 if $\Delta \leq \|s_n\|$ we set $s_{ddl} = \Delta s_n / \|s_n\|$;
- 4 if $\Delta > \|s_n\|$ we set $s_{ddl} = s_n$;



Solving

$$\alpha^2 \|s_g\|^2 + \gamma_*^2 (1 - \alpha)^2 \|s_n\|^2 + 2\gamma_* \alpha (1 - \alpha) s_g^T s_n = \Delta^2$$

we have that if $\|s_g\| \leq \Delta \leq \gamma_* \|s_n\|$ the root in $[0, 1]$ is given by:

$$A = \gamma_*^2 \|s_n\|^2 - \|s_g\|^2$$

$$B = \Delta^2 - \|s_g\|^2$$

$$\alpha = \frac{A - B}{A + \sqrt{AB}}$$



Algorithm (Computing Double DogLeg step)

DoubleDoglegStep(s_g, s_n, Δ):

$\gamma_* \leftarrow \|s_g\|^2 / (s_g^T s_n)$;

if $\Delta \leq \|s_g\|$ **then**

$s \leftarrow \Delta s_g / \|s_g\|$;

else if $\Delta \leq \gamma_* \|s_n\|$ **then**

$A \leftarrow \gamma_*^2 \|s_n\|^2 - \|s_g\|^2$;

$B \leftarrow \Delta^2 - \|s_g\|^2$;

$\alpha \leftarrow (A - B) / (A + \sqrt{AB})$;

$s \leftarrow \alpha s_g + (1 - \alpha) s_n$;

else if $\Delta \leq \|s_n\|$ **then**

$s \leftarrow \Delta s_n / \|s_n\|$;

else

$s \leftarrow s_n$;

end if

return s ;



Outline

- 1 The Trust Region method
- 2 Convergence analysis
- 3 The exact solution of trust region step
- 4 The dogleg trust region step
- 5 The double dogleg trust region step
- 6 Two dimensional subspace minimization



Two dimensional subspace minimization

- When H is positive definite the dogleg step can be improved by widening the search subspace

$$s = \arg \min_{\|\alpha s_g + \beta s_n\| \leq \Delta} f(\alpha s_g + \beta s_n)$$

i.e. we must solve a two dimensional constrained problem.

- The 2D problem results:

$$\begin{aligned} f(\alpha s_g + \beta s_n) &= f_0 + \mathbf{g}^T (\alpha s_g + \beta s_n) \\ &\quad + \frac{1}{2} (\alpha s_g + \beta s_n)^T \mathbf{H} (\alpha s_g + \beta s_n) \\ &= f_0 + \alpha \mathbf{g}^T s_g + \beta \mathbf{g}^T s_n \\ &\quad + \frac{1}{2} \alpha^2 s_g^T \mathbf{H} s_g + \frac{1}{2} \beta^2 s_n^T \mathbf{H} s_n + \alpha \beta s_g^T \mathbf{H} s_n \end{aligned}$$



Two dimensional subspace minimization

The 2D problem written in matrix form:

$$f(\alpha, \beta) = f_0 + \mathbf{b}^T \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \frac{1}{2} \begin{pmatrix} \alpha & \beta \end{pmatrix} \mathbf{A} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$$\mathbf{b} = \begin{pmatrix} \mathbf{g}^T s_g \\ \mathbf{g}^T s_n \end{pmatrix}$$

$$\mathbf{A} = \begin{pmatrix} s_g^T \mathbf{H} s_g & s_g^T \mathbf{H} s_n \\ s_g^T \mathbf{H} s_n & s_n^T \mathbf{H} s_n \end{pmatrix}$$

and the constraint

$$\|\alpha s_g + \beta s_n\|^2 = \begin{pmatrix} \alpha & \beta \end{pmatrix} \mathbf{D} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$$\mathbf{D} = \begin{pmatrix} s_g^T s_g & s_g^T s_n \\ s_g^T s_n & s_n^T s_n \end{pmatrix}$$



Lemma

Consider the following constrained quadratic problem where $\mathbf{H} \in \mathbb{R}^{n \times n}$, $\mathbf{D} \in \mathbb{R}^{n \times n}$ are *symmetric and positive definite*.

$$\text{Minimize} \quad f(s) = f_0 + \mathbf{g}^T s + \frac{1}{2} s^T \mathbf{H} s,$$

$$\text{Subject to} \quad s^T \mathbf{D} s \leq r^2$$

Then the following curve




$$s(\mu) \doteq -(\mathbf{H} + \mu \mathbf{D})^{-1} \mathbf{g},$$

for any $\mu \geq 0$ defines a descent direction for $f(s)$. Moreover

- there exists a unique μ_* such that $\|s(\mu_*)\| = \Delta$ and $s(\mu_*)$ is the solution of the constrained problem;
- or $\|s(0)\| < \Delta$ and $s(0)$ is the solution of the constrained problem.



References

- 
 Jorge Nocedal, and Stephen J. Wright
 Numerical optimization
 Springer, 2006
- 
 J. Stoer and R. Bulirsch
 Introduction to numerical analysis
 Springer-Verlag, Texts in Applied Mathematics, 12, 2002.
- 
 J. E. Dennis, Jr. and Robert B. Schnabel
 Numerical Methods for Unconstrained Optimization and
 Nonlinear Equations
 SIAM, Classics in Applied Mathematics, 16, 1996.

