

Chapter 25 Trust Region Method

1. Trust Region Model
2. Trust Region Radius and Subproblem



Trust Region Method

revisit Newton method

$$\therefore f(\mathbf{x}) \approx q_k(\mathbf{x}) = f(\mathbf{x}^{(k)}) + (\mathbf{g}^{(k)})^\top (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{F}_k (\mathbf{x} - \mathbf{x}^{(k)}),$$

where $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$, $\mathbf{F}_k = \nabla^2 f(\mathbf{x}^{(k)})$.

$$\therefore \text{Newton method: given } \mathbf{x}^{(k)}, \mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} q_k(\mathbf{x}), \text{ i.e.,}$$
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{F}_k)^{-1} \mathbf{g}^{(k)}$$

Disadvantages:

- Newton method is locally convergent, i.e., $\mathbf{x}^{(0)}$ should be near to \mathbf{x}^* .
- Newton method may fail if $\mathbf{F} \neq 0$.

motivation: restrict line search in neighborhood of $\mathbf{x}^{(k)}$

trust region: $\Omega_k = \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}^{(k)}\| \leq r_k\}$, where $r_k > 0$ is trust region radius.

trust region method: given $\mathbf{x}^{(k)}$,
 $\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \Omega_k} q_k(\mathbf{x})$ and tune radius r_k dynamically at each step.

★ $\|\cdot\|$ can be any ℓ^p -norm. ℓ^2 -norm is popular.

★ \mathbf{F}_k in $q_k(\mathbf{x})$ can be replaced by any approximation, e.g., BFGS, DFP, ...



Trust Region Method

(w.l.o.g) we use

$$f(\mathbf{x}) \approx q_k(\mathbf{x}) = f(\mathbf{x}^{(k)}) + (\mathbf{g}^{(k)})^\top (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{B}_k (\mathbf{x} - \mathbf{x}^{(k)}),$$

where $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$, $\mathbf{B}_k \approx \nabla^2 f(\mathbf{x}^{(k)})$.

trust region subproblem with general approximation

$$\begin{aligned} \min q_k(\mathbf{x}) &= f(\mathbf{x}^{(k)}) + (\mathbf{g}^{(k)})^\top (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{B}_k (\mathbf{x} - \mathbf{x}^{(k)}) \\ \text{s.t. } \|\mathbf{x} - \mathbf{x}^{(k)}\| &\leq r_k \end{aligned}$$

By defining $\mathbf{s} = \mathbf{x} - \mathbf{x}^{(k)}$, trust region subproblem reduces to:

$$\min (\mathbf{g}^{(k)})^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \mathbf{B}_k \mathbf{s}, \quad \text{s.t. } \|\mathbf{s}\| \leq r_k. \quad (1)$$

properties of trust region method

- ① rapid locally convergent (rate), and ideal globally convergent.
- ② Instead of first direction then stepsize in linear search, trust region method first limit stepsize and then determine direction.
- ③ trust region method can be efficient for nonconvex problem.

Trust Region Subproblem

Quadratic problem: $\min f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{b}^\top \mathbf{x}$, s.t. $\|\mathbf{x}\|_2 \leq 1$.

- If $\mathbf{b} = 0$, then the optimum is

$$\mathbf{x}^* = \begin{cases} 0, & \text{if } \lambda_{\min}(\mathbf{Q}) \geq 0 \\ \text{unit eigenvector of } \mathbf{Q} \text{ for } \lambda_{\min}(\mathbf{Q}), & \text{if } \lambda_{\min}(\mathbf{Q}) < 0 \end{cases}$$

- If $\mathbf{b} \neq 0$, the KKT condition is

$$\mathbf{Q}\mathbf{x}^* + \mathbf{b} + \lambda^* \mathbf{x}^* = 0, (\|\mathbf{x}^*\|^2 - 1)\lambda^* = 0, \lambda^* \geq 0, \mathbf{Q} + \lambda^* \mathbf{I} \succeq 0. \quad (2)$$

proof. Optimum \mathbf{x}^* always exists (why?). If $\|\mathbf{x}^*\| < 1$, then \mathbf{x}^* is an unconstrained local minimizer of f . Moreover, it must be a global minimizer (why?), KKT condition (2) holds with $\lambda^* = 0$.

If $\|\mathbf{x}^*\| = 1$, let $h(\mathbf{x}) = \frac{1}{2} \max\{0, \|\mathbf{x}\|^2 - 1\}$, $\alpha > 0$ and $f_k(\mathbf{x}) = f(\mathbf{x}) + k|h(\mathbf{x})|^2 + \frac{\alpha}{2}\|\mathbf{x} - \mathbf{bx}^*\|^2$.



Trust Region Subproblem

- ③ Optimum x^* always exists (why?). If $\|x^*\| < 1$, then x^* is an unconstrained local minimizer of f . Moreover, it must be a global minimizer (why?), KKT condition (2) holds with $\lambda^* = 0$.
- ④ If $\|x^*\| = 1$, let $h(x) = \frac{1}{2} \max\{0, \|x\|^2 - 1\}$, $\alpha > 0$ and $f_k(x) = f(x) + k|h(x)|^2 + \frac{\alpha}{2}\|x - bx^*\|^2$.
- ⑤ Let $x^k = \arg \min\{f_k(x) \mid \|x - x^*\| \leq 1\}$. We will show that x^k is an unconstrained local min of f_k for all large k .
- ⑥ Taking limit $k \rightarrow \infty$ of $f_k(x^k) = f(x^k) + k|h(x^k)|^2 + \frac{\alpha}{2}\|x^k - x^*\|^2 \leq f_k(x^*) = f(x^*)$, along any convergent subsequence of $\{x^k\}$, we get $h(\bar{x}) = \lim_{k \rightarrow \infty} h(x^k) \rightarrow 0$.
- ⑦ Furthermore, taking limit of $f(x^k) + \frac{\alpha}{2}\|x^k - x^*\|^2 \leq f(x^*)$ shows $f(\bar{x}) + \frac{\alpha}{2}\|\bar{x} - x^*\|^2 \leq f(x^*)$
- ⑧ Since $h(\bar{x}) = 0$, it follows that $f(x^*) \leq f(\bar{x})$. Thus, we have $\bar{x} = x^*$ and $f(x^*) = f(\bar{x})$.
- ⑨ Since \bar{x} is any limit point, we have $x^k \rightarrow x^*$, so $\|x^k - x^*\| < 1$ for large k , $\Rightarrow x^k$ is an unconstrained local min of f_k , $\nabla f_k(x^k) = 0$, $\nabla^2 f_k(x^k) \succeq 0$.



Trust Region Subproblem

- 10 Taking limit of $0 = \nabla f(x^k) + 2kh(x^k)\nabla h(x^k) + \alpha(x^k - x^*) = Qx^k + b + 2kh(x^k)x^k + \alpha(x^k - x^*)$. shows
 $2kh(x^k) = -\frac{(x^k)^\top(Qx^k + b + \alpha(x^k - x^*))}{\|x^k\|^2} \rightarrow -(x^*)^\top(Qx^* + b) \equiv \lambda^*$. Taking limit in (3) yields $Qx^* + b + \lambda^*x^* = 0$, $\lambda^* \geq 0$.
- 11 The remaining condition $Q + \lambda^*I \succeq 0$ follows from a contra-positive argument: if $Q + \lambda^*I \not\succeq 0$, then there exists a $u \neq 0$ (e.g., the eigenvector of $Q + \lambda^*I$ corresponding to a negative eigenvalue, perturbed if necessary) s.t.

$$u^\top(Q + \lambda^*I)u < 0, \quad \langle u, x^* \rangle < 0.$$

Then, perturb x^* by u to derive a contradiction to the optimality of x^* .

- 12 Finally, check the sufficiency: suppose (2) holds for some x^* , λ^* . Consider the following quadratic function

$$f_{\lambda^*}(x) = f(x) + \frac{1}{2}\lambda^*(\|x\|^2 - 1)$$

Claim: x^* is a global min of f_{λ^*} . So x^* is a global min of f over the sphere $\|x\|^2 = 1$.



Solving Trust Region Subproblem

- Observation: if $Q + \lambda I \succ 0$, then

$$Qx + b + \lambda x = 0 \quad \Rightarrow \quad x(\lambda) = (Q + \lambda I)^{-1}b.$$

Moreover, $\|x(\lambda)\|$ decreases w.r.t λ .

- Note: $x(-\lambda_{\min}(Q))$ is not unique, and $\|x(-\lambda_{\min}(Q))\|$ can be made arbitrarily large (why?).
- Strategy: binary search of λ over $[\lambda_b, \infty)$ with $\lambda_b = \min\{0, -\lambda_{\min}(Q)\}$,
 - ★ Check if $\|x(\lambda_b)\| \leq 1$. If yes, then $x(\lambda_b)$ is an optimum (why?). Stop.
 - ★ Else, we must have $\|x(\lambda_b)\| > 1$. Let $\lambda_a = \min\{1, \lambda_b\}$, and check if $\|x(\lambda_a)\| \leq 1$. If not, update $\lambda_a := 2\lambda_a$ until $\|x(\lambda_a)\| \leq 1$.
 - ★ Let $\lambda = (\lambda_a + \lambda_b)/2$. If $\|x(\lambda)\| \leq 1$, then update $\lambda_a = \lambda$. Else, we update $\lambda_b = \lambda$. Stop when $|\lambda_b - \lambda_a| \leq \epsilon$.



Trust Region Method

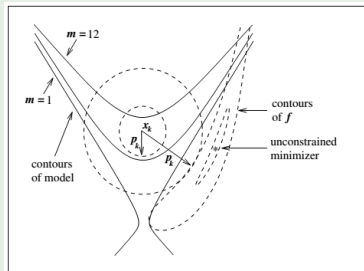
Example $(\min_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) = 10(x_2 - x_1^2)^2 + (1 - x_1)^2)$

At the current point $\mathbf{x}^{(k)} = (0, 1)$, the gradient and Hessian of f at $\mathbf{x}^{(k)}$ are

$$\nabla f(\mathbf{x}^{(k)}) = \begin{pmatrix} -2 \\ 20 \end{pmatrix}, \quad \nabla^2 f(\mathbf{x}^{(k)}) = \begin{pmatrix} -38 & 0 \\ 0 & 20 \end{pmatrix}.$$

Two possible trust regions (circles) and their corresponding steps s_k .

The solid lines are contours of the model function in (1).



how to tune radius r_k ?

when there is good agreement between the model $q_k(\mathbf{x}^{(k)})$ and the objective function $f(\mathbf{x}^{(k)})$, one should set r_k as large as possible.

Trust Region Radius

Definition (ratio for measuring agreement)

- actual reduction: $\text{Ared}_k = f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k+1)})$,
- predicted reduction: $\text{Pred}_k = q_k(\mathbf{x}^{(k)}) - q_k(\mathbf{x}^{(k+1)})$,
- ratio for measuring agreement: $\rho_k = \frac{\text{Ared}_k}{\text{Pred}_k}$.

procedure for tuning radius r_k

- if $\rho_k \approx 1$, it is good agreement between f and q_k . Then, dilate trust region;
- if $\rho_k \approx 0$ or $\rho_k \leq 0$, then reduce the trust region;
- otherwise, don't alter the trust region.

For example: let $\bar{r} > 0$ and $x(r_k)$ be the optimum of (1), denote

$$r_{k+1} := \begin{cases} 0.5r_k, & \text{if } \rho < 0.25 \\ r_k, & \text{if } \rho \in [0.25, 0.75] \\ \min\{2r_k, \bar{r}\}, & \text{if } \rho > 0.75 \end{cases}$$

Solve Trust Region Subproblem

By ignoring index, trust region subproblem (3) amounts to

$$\min q(s) = \mathbf{g}^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \mathbf{B} \mathbf{s}, \quad \text{s.t. } \|\mathbf{s}\| \leq r. \quad (3)$$

By KKT conditions, \mathbf{s}^* is an optimum of (3) $\iff \exists$ scalar $\lambda^* \geq 0$ such that

$$(\mathbf{B} + \lambda^* \mathbf{I}) \mathbf{s}^* = -\mathbf{g}, \quad \|\mathbf{s}^*\|_2 \leq r, \quad \lambda^*(r - \|\mathbf{s}^*\|_2) = 0, \quad \mathbf{B} + \lambda^* \mathbf{I} \succ 0.$$

Cauchy point

Let $\mathbf{s}^G = -\frac{r\mathbf{g}}{\|\mathbf{g}\|_2}$ (why?) be the optimum of linear version of (3), i.e.,

$$\mathbf{s}^G = \arg \min_{\|\mathbf{s}\| \leq r} q(\mathbf{s}) = \mathbf{g}^\top \mathbf{s}.$$

The Cauchy point is defined by $\mathbf{s}^c = \tau \mathbf{s}^G = -\tau \frac{r\mathbf{g}}{\|\mathbf{g}\|_2}$, where

$$\tau = \begin{cases} 1, & \mathbf{g}^\top \mathbf{B} \mathbf{g} \leq 0; \\ \min\{\|\mathbf{g}\|_2^3 / (r\mathbf{g}^\top \mathbf{B} \mathbf{g}), 1\} & \text{otherwise.} \end{cases}$$

Cauchy point (version 2)

- Cauchy point \mathbf{x}^c is defined by

$$\begin{cases} \mathbf{z} := \arg \min \{ \mathbf{g}^\top (\mathbf{x} - \mathbf{x}^{(k)}) \mid \|\mathbf{x} - \mathbf{x}^{(k)}\| \leq r \}, \\ \mathbf{x}^c := \arg \min_{\tau \in (0,1)} \left\{ \tau \mathbf{g}^\top (\mathbf{z} - \mathbf{x}^{(k)}) + \frac{\tau^2}{2} (\mathbf{z} - \mathbf{x}^{(k)})^\top \mathbf{B} (\mathbf{z} - \mathbf{x}^{(k)}) \right\}. \end{cases}$$

- Optima of \mathbf{z} and \mathbf{x}^c are

$$\begin{cases} \mathbf{z} = \mathbf{x}^{(k)} - \frac{r\mathbf{g}}{\|\mathbf{g}\|}, \\ \mathbf{x}^c = \mathbf{x}^{(k)} - \frac{\tau r \mathbf{g}}{\|\mathbf{g}\|}, \end{cases} \quad \text{where } \tau = \begin{cases} 1, & \text{if } \mathbf{g}^\top \mathbf{B} \mathbf{g} \leq 0 \\ \min \left\{ \frac{\|\mathbf{g}\|^3}{r \mathbf{g}^\top \mathbf{B} \mathbf{g}}, 1 \right\}, & \text{else} \end{cases}$$

- Sufficient descent property: $\exists c_1 \in (0, 1]$ s.t. $q(\mathbf{x}^c) \leq -c_1 \|\mathbf{g}\| \min \left(r, \frac{\|\mathbf{g}\|}{\|\mathbf{F}\|} \right)$.
- Hessian information \mathbf{F} can be used to accelerate convergence, e.g., find the Cauchy point \mathbf{x}^c by

$$\min_{\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq r} \left\{ \mathbf{g}^\top (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{F} (\mathbf{x} - \mathbf{x}^{(k)}) \mid \mathbf{x} - \mathbf{x}^{(k)} \in V \right\},$$

where $V := \text{span}\{\mathbf{g}, \mathbf{F}^{-1}\mathbf{g}\}$. It is simple to solve (reducible to root-finding of a 4th order polynomial which admits closed form solution).



Convergence of Trust Region Method

Theorem

Suppose the level set $\mathcal{S} := \{\mathbf{x} \mid f(\mathbf{x}) \leq f(\mathbf{x}^{(0)})\}$ is bounded, Hessian \mathbf{F} is bounded on \mathcal{S} . Assume $\mathbf{x}^{(k)}$ is at least as good as Cauchy point, i.e., $q(\mathbf{x}^{(k)}) \leq -c_1 \|\mathbf{g}\| \min\left(r, \frac{\|\mathbf{g}\|}{\|\mathbf{F}\|}\right)$. Then, $\mathbf{g} \rightarrow 0$ as $k \rightarrow \infty$. In practice, \mathbf{F} can be replaced by uniformly bounded approximation \mathbf{B} .

