



EECS 559 Optimization Methods for SIPML

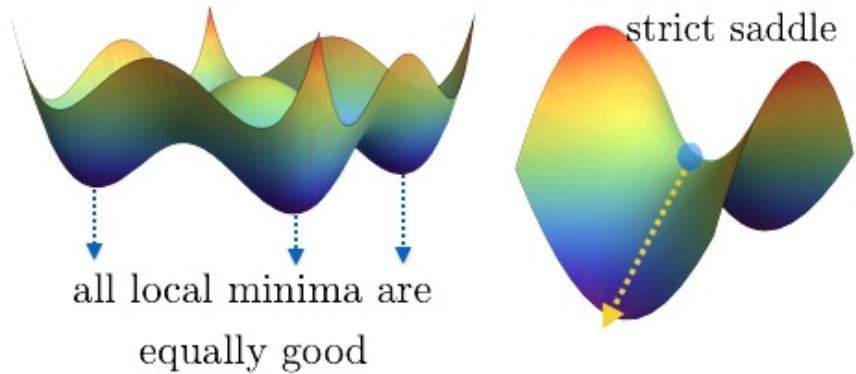
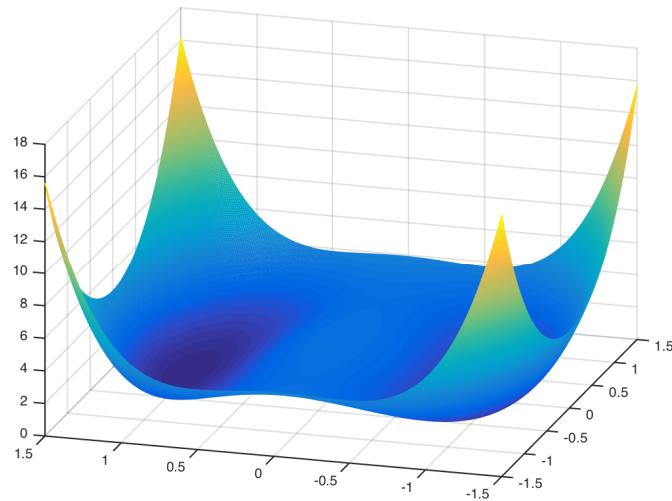
Lecture 12 – Trust-Region and Cubic Regularization

Instructor: Prof. Qing Qu (qingqu@umich.edu)

Lecture Agenda

- Trust-Region Method
 - Algorithmic Introduction
 - Solving the TRM Subproblem
 - Convergence for Strict Saddle Function
- Cubic Regularization of Newton's Method

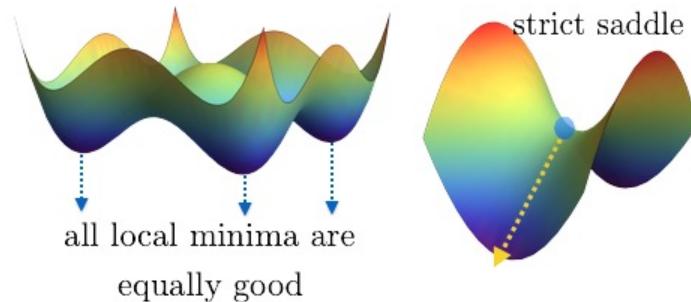
Example: Generalized Phase Retrieval



$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{i=1}^m (y_i - (\mathbf{a}_i^\top \mathbf{x})^2)^2$$

Unconstraint Nonconvex Problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$



- Assume the function $f(\mathbf{x}) : \mathbb{R}^n \mapsto \mathbb{R}$ is twice continuously differentiable.
- We aim to solve nonconvex/nonlinear problems, e.g. strict saddle functions.
- We need to *exploit the negative curvature* of functions to escape saddle points.

Idea: Optimizing Local Approximation

At every iteration \boldsymbol{x}_k , approximate $f(\boldsymbol{x}_k + \boldsymbol{d})$ by

- **Linear model:**

$$f_L(\boldsymbol{d}; \boldsymbol{x}_k) = f(\boldsymbol{x}_k) + \langle \nabla f(\boldsymbol{x}_k), \boldsymbol{d} \rangle$$

- **Quadratic model:**

$$f_Q(\boldsymbol{d}; \boldsymbol{x}_k) = f(\boldsymbol{x}_k) + \langle \nabla f(\boldsymbol{x}_k), \boldsymbol{d} \rangle + \frac{1}{2} \boldsymbol{d}^\top \boldsymbol{B}_k \boldsymbol{d}$$

for some symmetric matrix $\boldsymbol{B}_k \approx \nabla^2 f(\boldsymbol{x}_k)$
(e.g., quasi-Newton SR-1/SR-2 update).

Idea: Optimizing Local Approximation

- We want to optimize

$$f_Q(\mathbf{d}; \mathbf{x}_k) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d}$$

for some symmetric matrix

(e.g., quasi-Newton SR-1/SR-2 update).

- Difficulties for optimization:

- the approximation error can be large when $\|\mathbf{d}\|_2$ is large.
- minimizing the model might not be possible
(e.g., $f_Q(\mathbf{d}; \mathbf{x}_k)$ is unbounded when \mathbf{B}_k is *indefinite*)

Trust-region Methods (TRM)

Idea: add a trust-region constraint

$$\|\mathbf{d}\|_{\square} \leq \Delta_k,$$

for some “suitable” trust-region radius $\Delta_k > 0$ and solve

$$\mathbf{d}_k = \arg \min_{\mathbf{d} \in \mathbb{R}^n} f_Q(\mathbf{d}; \mathbf{x}_k), \quad \text{s.t.} \quad \|\mathbf{d}\|_{\square} \leq \Delta_k.$$

- We usually choose $\|\cdot\|_{\square} = \|\cdot\|_2$. Other common norms are $\|\cdot\|_1$ and $\|\cdot\|_{\infty}$.
- The global convergence does not depend on the choice of $\|\cdot\|_{\square}$, but the empirical performance do.

TRM Subproblem

Definition. (Trust-region subproblem)

The trust-region subproblem at the k -th iterate is

$$\min_{\mathbf{d} \in \mathbb{R}^n} f_Q(\mathbf{d}; \mathbf{x}_k) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d}, \text{ s.t. } \|\mathbf{d}\|_\square \leq \Delta_k,$$

where \mathbf{B}_k is a symmetric matrix and $\Delta_k > 0$ is the trust-region radius.

- Once the step \mathbf{d}_k is computed, update

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$$

Linesearch vs Trust-Region Methods

- Linesearch methods:
 - Descent direction \mathbf{d}_k first (e.g., steepest descent, Newton)
 - Then compute a stepsize τ_k based on \mathbf{d}_k to reduce $f(\mathbf{x}_k + \tau_k \mathbf{d}_k)$
 - Update $\mathbf{x}_{k+1} = \mathbf{x}_k + \tau_k \mathbf{d}_k$
- Trust-region methods:
 - Length first (trust-region radius Δ_k), direction second (“solve” subproblem for \mathbf{d}_k)
 - pick \mathbf{d}_k to reduce some (quadratic) model of $f(\mathbf{x}_k + \mathbf{d})$
 - Update $\mathbf{x}_{k+1} \leftarrow \begin{cases} \mathbf{x}_k + \mathbf{d}_k & \text{if } f(\mathbf{x}_k + \mathbf{d}_k) < f(\mathbf{x}_k) \\ \mathbf{x}_k & \text{otherwise} \end{cases}$

Choice of the Trust-Region Radius

How do we choose the trust-region radius to guarantee sufficient decrease $f_Q(\mathbf{d}_k; \mathbf{x}_k) < f_Q(\mathbf{0}; \mathbf{x}_k)$?

- **Measure of predicted model decrease:**

$$f_Q(\mathbf{0}; \mathbf{x}_k) - f_Q(\mathbf{d}; \mathbf{x}_k) = f(\mathbf{x}_k) - f_Q(\mathbf{d}; \mathbf{x}_k)$$

- **Measure of actual function decrease:**

$$f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{d}_k)$$

- **Measure of the quality of reduction:**

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{d}_k)}{f(\mathbf{x}_k) - f_Q(\mathbf{d}; \mathbf{x}_k)},$$

characterizes how well $f_Q(\mathbf{d}; \mathbf{x}_k)$ models $f(\mathbf{x}_k + \mathbf{d})$.

Choice of the Trust-Region Radius

Solve the trust-region subproblem for \mathbf{d}_k so that we have $f_Q(\mathbf{d}_k; \mathbf{x}_k) < f_Q(\mathbf{0}; \mathbf{x}_k)$. Set

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{d}_k)}{f_Q(\mathbf{0}; \mathbf{x}_k) - f_Q(\mathbf{d}; \mathbf{x}_k)},$$

- If $\rho_k \geq \eta_{vs}$, set $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \mathbf{d}_k$, $\Delta_{k+1} \leftarrow \gamma_i \Delta_k$
increase the radius
- If $\eta_s \leq \rho_k < \eta_{vs}$, set $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \mathbf{d}_k$, $\Delta_{k+1} \leftarrow \Delta_k$
- If $0 < \rho_k < \eta_s$, set $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k$, $\Delta_{k+1} \leftarrow \gamma_d \Delta_k$
decrease the radius

Algorithm 1 Trust-region Algorithm

Input: $0 < \gamma_d < 1 < \gamma_i$, and $0 < \eta_s \leq \eta_{vs} < 1$

Initialize \mathbf{x}_0 and $\Delta_0 = 1$.

for $k = 0, 1, 2, \dots$ **do**

 Build the 2nd-order model $f_Q(\mathbf{d})$ of $f(\mathbf{x}_k + \mathbf{d})$;

 Solve TRM subproblem for \mathbf{d}_k such that $f_Q(\mathbf{d}_k; \mathbf{x}_k) < f(\mathbf{x}_k)$

 Set $\rho_k \leftarrow \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{d}_k)}{f(\mathbf{x}_k) - f_Q(\mathbf{d}_k; \mathbf{x}_k)}$

if $\rho_k > \eta_{vs}$, **then**

 set $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \mathbf{d}_k$ and $\Delta_{k+1} \leftarrow \gamma_i \Delta_k$ (**very successful**)

else if $\rho_k \geq \eta_s$, **then**

 set $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \mathbf{d}_k$ and $\Delta_{k+1} \leftarrow \Delta_k$ (**successful**)

else, set $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k$ and $\Delta_{k+1} \leftarrow \gamma_d \Delta_k$ (**unsuccessful**)

end for

Extra Details: Trust-Region Methods

- **Termination criteria:** in practice, we often stop when

$$\|\nabla f(\mathbf{x}_k)\|_2 \leq 10^{-6} \max \{1, \|\nabla f(\mathbf{x}_0)\|_2\}$$

- **Typical parameter values:**

$$\eta_s = 0.1, \quad \eta_{vs} = 0.9, \quad \gamma_d = \frac{1}{2}, \quad \gamma_i = 2.$$

Further Readings

- *Numerical Optimization*. Jorge Nocedal and Stephen J. Wright.
(Chapter 4)
- *Trust Region Methods*. Andrew R. Conn, Nicholas I.M. Gould, and Philippe L. Toint. MOS-SIAM Series on Optimization, 2000.

Lecture Agenda

- Trust-Region Method
 - Algorithmic Introduction
 - Solving the TRM Subproblem
 - Convergence for Strict Saddle Function
- Cubic Regularization of Newton's Method

TRM Subproblem

Definition. (Trust-region subproblem)

The trust-region subproblem at the k -th iterate is

$$\min_{\mathbf{d} \in \mathbb{R}^n} f_Q(\mathbf{d}; \mathbf{x}_k) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d}, \text{ s.t. } \|\mathbf{d}\|_\square \leq \Delta_k,$$

where \mathbf{B}_k is a symmetric matrix and $\Delta_k > 0$ is the trust-region radius.

- Once the step \mathbf{d}_k is computed, update

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$$

The Effect of Different Norms

The choice of $\|\cdot\|_{\square}$ determines the behavior of \mathbf{d}_k as $\Delta_k \rightarrow 0$
If $\|\mathbf{d}\|_{\square} \ll 1$, then

$$\begin{aligned} f_Q(\mathbf{d}; \mathbf{x}_k) &= f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d} \\ &\approx f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle \end{aligned}$$

so that for $\Delta_k \ll 1$, we have

$$\min_{\|\mathbf{d}\|_{\square} \leq \Delta_k} f_Q(\mathbf{d}; \mathbf{x}_k) \approx \min_{\|\mathbf{d}\|_{\square} \leq \Delta_k} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle$$

The Effect of Different Norms

The solution approaches the *steepest-descent* direction of length Δ_k as $\Delta_k \rightarrow 0$

- **ℓ_2 -norm:**

$$\lim_{\Delta_k \rightarrow 0} \mathbf{d}_k \rightarrow -\frac{\Delta_k}{\|\nabla f(\mathbf{x}_k)\|_2} \nabla f(\mathbf{x}_k)$$

- **ℓ_∞ -norm:**

$$\lim_{\Delta_k \rightarrow 0} \mathbf{d}_k \rightarrow -\Delta_k \hat{\mathbf{e}}, \quad \text{with} \quad \hat{e}_j = \text{sign}([\nabla f(\mathbf{x}_k)]_j)$$

TRM Subproblem with ℓ_∞ -norm

For ℓ_∞ -norm, the subproblem is

$$\begin{aligned} & \min_{\mathbf{d} \in \mathbb{R}^n} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d}, \\ & \text{s.t. } -\Delta_k \mathbf{1} \leq \mathbf{d} \leq \Delta_k \mathbf{1}, \end{aligned}$$

- This is a quadratic program (QP) and possibly nonconvex (i.e., when \mathbf{B}_k is indefinite);
- Finding the global minimizer could be NP-hard.

TRM Subproblem with ℓ_∞ -norm

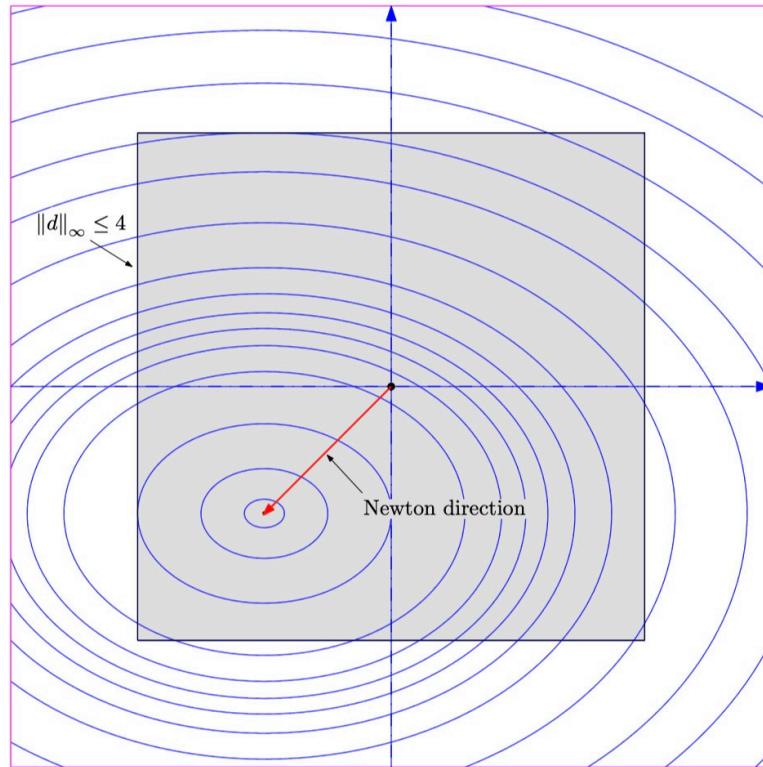
Example: Consider the following TRM subproblem

$$\min_{\mathbf{d}} f_Q(\mathbf{d}) = f + \mathbf{g}^\top \mathbf{d} + \frac{1}{2} \mathbf{d}^\top \mathbf{B} \mathbf{d}, \quad \text{s.t. } \|\mathbf{d}\|_\infty \leq 4,$$

with $f = 0$, $\mathbf{g} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$, and $\mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$.

The unique global minimizer $\mathbf{d}_* = \begin{bmatrix} -2 \\ -2 \end{bmatrix}$ lies inside the trust region with $f_Q(\mathbf{d}_*) = -6$.

TRM Subproblem with ℓ_∞ -norm



TRM Subproblem with ℓ_∞ -norm

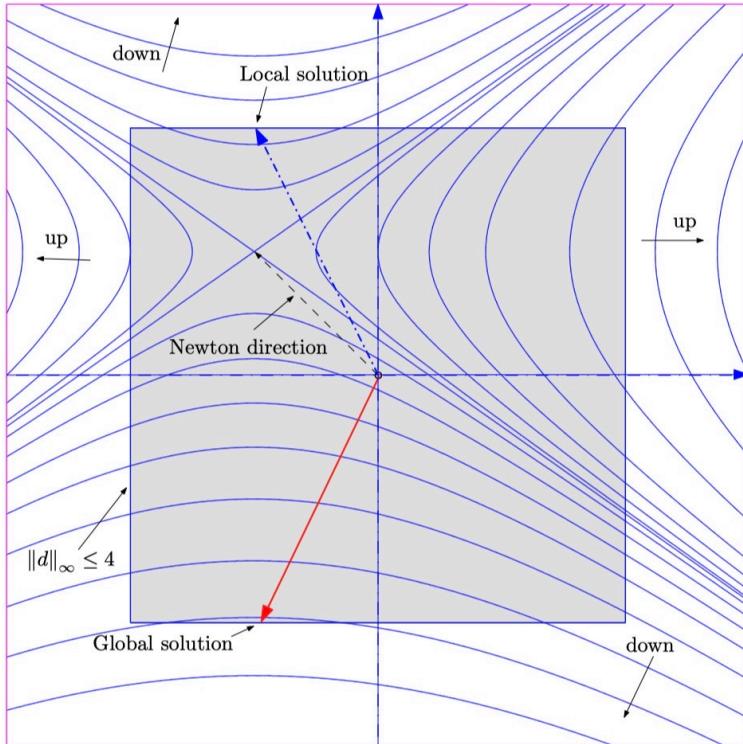
Example: Consider the following TRM subproblem

$$\min_{\mathbf{d}} f_Q(\mathbf{d}) = f + \mathbf{g}^\top \mathbf{d} + \frac{1}{2} \mathbf{d}^\top \mathbf{B} \mathbf{d}, \quad \text{s.t. } \|\mathbf{d}\|_\infty \leq 4,$$

$$\text{with } f = 0, \mathbf{g} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \text{ and } \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix}.$$

Then $f_Q(\mathbf{d})$ is unbounded below and $\mathbf{d}_N = \begin{bmatrix} -2 \\ 2 \end{bmatrix}$ is the step to a saddle point of $f_Q(\mathbf{d})$.

TRM Subproblem with ℓ_∞ -norm



- the unique global minimizer lies on the boundary of the trust-region;
- there are two local solutions.

TRM Subproblem with ℓ_2 -norm

For ℓ_2 -norm, the subproblem is

$$\begin{aligned} \min_{\mathbf{d} \in \mathbb{R}^n} \quad & f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d}, \\ \text{s.t. } \|\mathbf{d}\|_2 \leq & \Delta_k, \end{aligned}$$

This is a nonlinearly constrained optimization problem

- the global minimizer can be computed *efficiently*
- we will focus on the ℓ_2 -norm trust-region subproblem

The Cauchy point of TRM Subproblem

Definition (Cauchy Point) At point \mathbf{x} , we call $\mathbf{y}_C := \mathbf{x} + \mathbf{d}_C$, the Cauchy point if $\mathbf{d}_C = -\tau_C \nabla f(\mathbf{x})$ with

$$\tau_C := \arg \min_{\tau > 0} f_Q(-\tau \nabla f(\mathbf{x}); \mathbf{x}), \quad \text{s.t.} \quad \|\tau \nabla f(\mathbf{x})\|_2 \leq \Delta$$

- Cauchy point requires minimal condition of sufficient decrease in the model;
- Using Cauchy direction can only guarantee convergence to a critical point $\nabla f(\mathbf{x}) \rightarrow 0$;
- We need a direction that can escape saddle point.

TRM Subproblem with ℓ_2 -norm

For ℓ_2 -norm, the subproblem is

$$\begin{aligned} \min_{\mathbf{d} \in \mathbb{R}^n} \quad & f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d}, \\ \text{s.t. } \|\mathbf{d}\|_2 \leq & \Delta_k, \end{aligned}$$

This is a nonlinearly constrained optimization problem

- Cauchy step \mathbf{d}_C implies a steepest descent like method;
- exact solution implies Newton-like method

Exact Solution of TRM Subproblem

Case 1: $\|d_\star\|_2 \leq \Delta$

$$\nabla f_Q(\mathbf{d}; \mathbf{x}) = \mathbf{0}$$

This gives the following optimality condition

$$\nabla f(\mathbf{x}) + \mathbf{Bd} = \mathbf{0} \implies \mathbf{Bd}_\star = -\nabla f(\mathbf{x})$$

so that d_\star is the Newton direction.

Exact Solution of TRM Subproblem

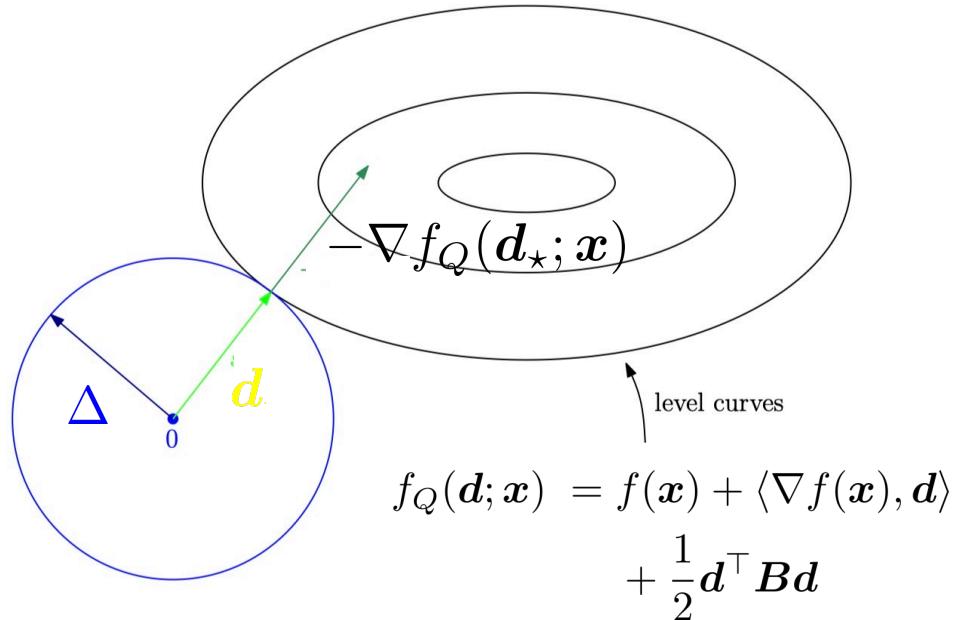
Case 2: $\|d_\star\|_2 = \Delta$

- d_\star and $-\nabla f_Q(d_\star; x)$ point in the same direction
- There exists $\lambda_\star > 0$ such that

$$\lambda_\star d_\star = -\nabla f_Q(d; x)$$

which is equivalent to

$$(B + \lambda_\star I)d_\star = -\nabla f(x)$$



Exact Solution of TRM Subproblem

Theorem. A vector \mathbf{d}_\star is a global minimizer of

$$\min_{\mathbf{d} \in \mathbb{R}^n} f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B} \mathbf{d}, \quad \text{s.t. } \|\mathbf{d}\|_2 \leq \Delta,$$

if and only if $\|\mathbf{d}_\star\|_2 \leq \Delta$ and there exists a scalar λ_\star such that

- $\lambda_\star \geq 0$;
- $(\mathbf{B} + \lambda_\star \mathbf{I}) \mathbf{d}_\star = -\nabla f(\mathbf{x})$
- $\mathbf{B} + \lambda_\star \mathbf{I}$ is positive semi-definite (PSD)
- $\lambda_\star (\|\mathbf{d}_\star\|_2 - \Delta) = 0$.

Moreover, if $\mathbf{B} + \lambda_\star \mathbf{I}$ is positive definite, then \mathbf{d}_\star is unique.

Approach I: Newton's Method

Thus, we need to consider three cases:

1. \mathbf{B} is PSD and

$$\mathbf{B}\mathbf{d}_\star = -\nabla f(\mathbf{x}), \quad \|\mathbf{d}_\star\|_2 \leq \Delta. \quad (\text{easy})$$

2. \mathbf{B} is PSD and $\lambda_\star > 0$ that

$$(\mathbf{B} + \lambda_\star \mathbf{I})\mathbf{d}_\star = -\nabla f(\mathbf{x}), \quad \|\mathbf{d}_\star\|_2 = \Delta. \quad (\text{challenging})$$

3. \mathbf{B} is *not* PSD (e.g., indefinite) and $\lambda_\star > 0$ that

$$(\mathbf{B} + \lambda_\star \mathbf{I})\mathbf{d}_\star = -\nabla f(\mathbf{x}), \quad \|\mathbf{d}_\star\|_2 = \Delta. \quad (\text{challenging})$$

We only need to consider Case 2 & Case 3 with $\|\mathbf{d}_\star\|_2 = \Delta$, solving nonlinear system in \mathbf{d} and λ .

Approach I: Newton's Method

We only need to consider Case 2 & Case 3 with $\|\mathbf{d}_\star\|_2 = \Delta$, solving nonlinear system in \mathbf{d} and λ such that

$$(\mathbf{B} + \lambda \mathbf{I})\mathbf{d} = -\nabla f(\mathbf{x}), \quad \|\mathbf{d}\|_2 = \Delta,$$

$$\mathbf{B} + \lambda \mathbf{I} \succeq \mathbf{0}, \quad \lambda \geq 0.$$

Goal. Find a scalar $\lambda_\star \geq \max \{0, -\lambda_n\}$ and a vector \mathbf{d}_\star such that

$$(\mathbf{B} + \lambda \mathbf{I})\mathbf{d}_\star = -\nabla f(\mathbf{x}), \quad \|\mathbf{d}_\star\|_2 = \Delta.$$

Approach I: Newton's Method

We only need to consider Case 2 & Case 3 with $\|\mathbf{d}_\star\|_2 = \Delta$, solving nonlinear system in \mathbf{d} and λ such that

$$(\mathbf{B} + \lambda \mathbf{I})\mathbf{d} = -\nabla f(\mathbf{x}), \quad \|\mathbf{d}\|_2 = \Delta,$$

$$\mathbf{B} + \lambda \mathbf{I} \succeq \mathbf{0}, \quad \lambda \geq 0.$$

- Given \mathbf{B} is symmetric, consider the eigen decomposition

$$\mathbf{B} = \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^\top, \quad \boldsymbol{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n),$$

so that $\mathbf{d}(\lambda) = -(\mathbf{B} + \lambda \mathbf{I})^{-1} \nabla f(\mathbf{x}) = -\sum_{i=1}^n \frac{\mathbf{v}_i^\top \nabla f(\mathbf{x})}{\lambda_i + \lambda} \mathbf{v}_i$

Approach I: Newton's Method

$$\mathbf{d}(\lambda) = -(\mathbf{B} + \lambda \mathbf{I})^{-1} \nabla f(\mathbf{x}) = -\sum_{i=1}^n \frac{\mathbf{v}_i^\top \nabla f(\mathbf{x})}{\lambda_i + \lambda} \mathbf{v}_i$$

Thus, we would like

$$\|\mathbf{d}(\lambda)\|_2^2 = \sum_{i=1}^n \frac{(\mathbf{v}_i^\top \nabla f(\mathbf{x}))^2}{(\lambda_i + \lambda)^2} = \Delta^2$$

- This implies that $\psi(\lambda) := \|\mathbf{d}(\lambda)\|_2$ has poles that $\lambda = -\lambda_i$ if $\mathbf{v}_i^\top \nabla f(\mathbf{x}) \neq 0$.

Approach I: Newton's Method

Question: How do we actually solve $\|\mathbf{d}(\lambda)\|_2 = \Delta$?

A: Solve the following secular equation

$$\phi(\lambda) := \frac{1}{\|\mathbf{d}(\lambda)\|_2} - \frac{1}{\Delta} = 0$$

using Newton's method (safeguarded).

- The method works except for the “hard” case.

Approach I: Newton's Method

Example: Consider the following TRM subproblem

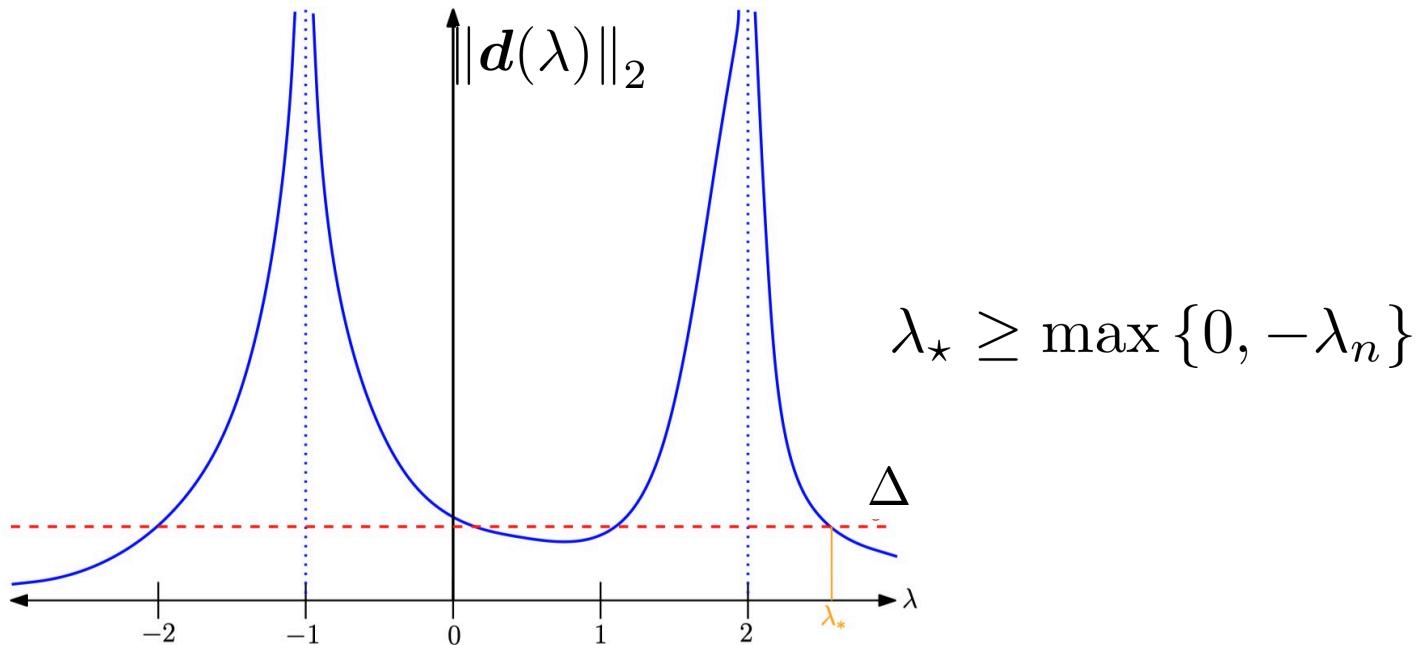
$$\min_{\mathbf{d}} f_Q(\mathbf{d}) = f + \mathbf{g}^\top \mathbf{d} + \frac{1}{2} \mathbf{d}^\top \mathbf{B} \mathbf{d}, \quad \text{s.t. } \|\mathbf{d}\|_2 \leq 4,$$

with $f = 0$, $\mathbf{g} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$, and $\mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix}$.

$$\mathbf{d}(\lambda) = \begin{bmatrix} -2/(\lambda + 1) \\ -4/(\lambda - 2) \end{bmatrix} \implies \|\mathbf{d}(\lambda)\|_2^2 = \frac{4}{(\lambda + 1)^2} + \frac{16}{(\lambda - 2)^2}$$

- It follows that $\psi(\lambda) = \|\mathbf{d}(\lambda)\|_2$ has two poles at $\lambda_1 = -1$ and $\lambda_2 = 2$.

Approach I: Newton's Method

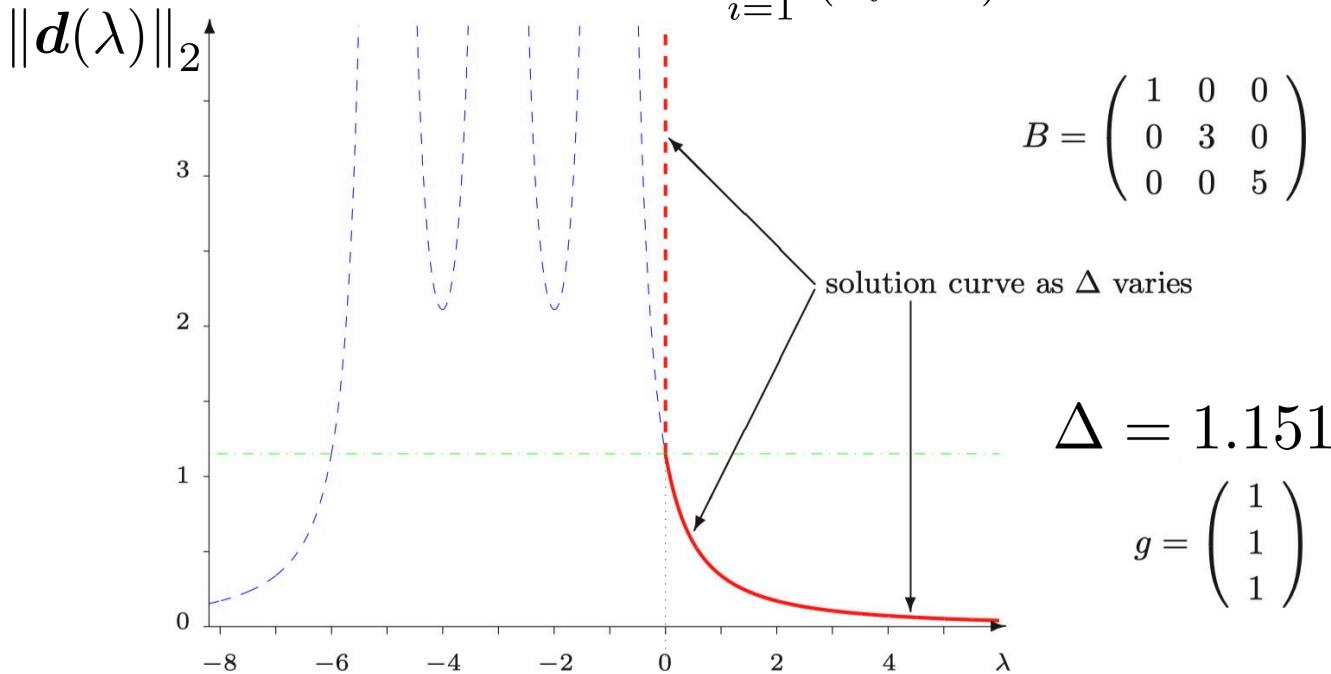


It is clear that for any $\Delta > 0$ there is a $\lambda_* > 2$ such that $\psi(\lambda_*) = \Delta$

Approach I: Newton's Method

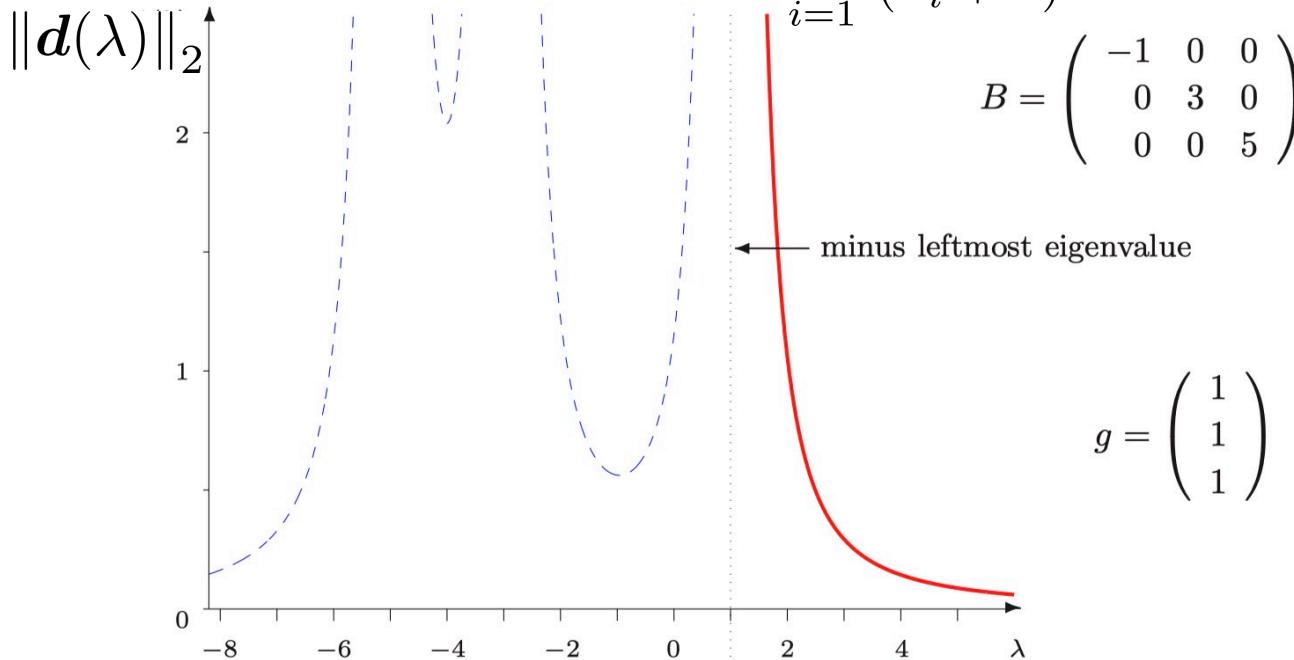
- convex case:

$$\|d(\lambda)\|_2^2 = \sum_{i=1}^n \frac{(v_i^\top g)^2}{(\lambda_i + \lambda)^2} = \Delta^2$$



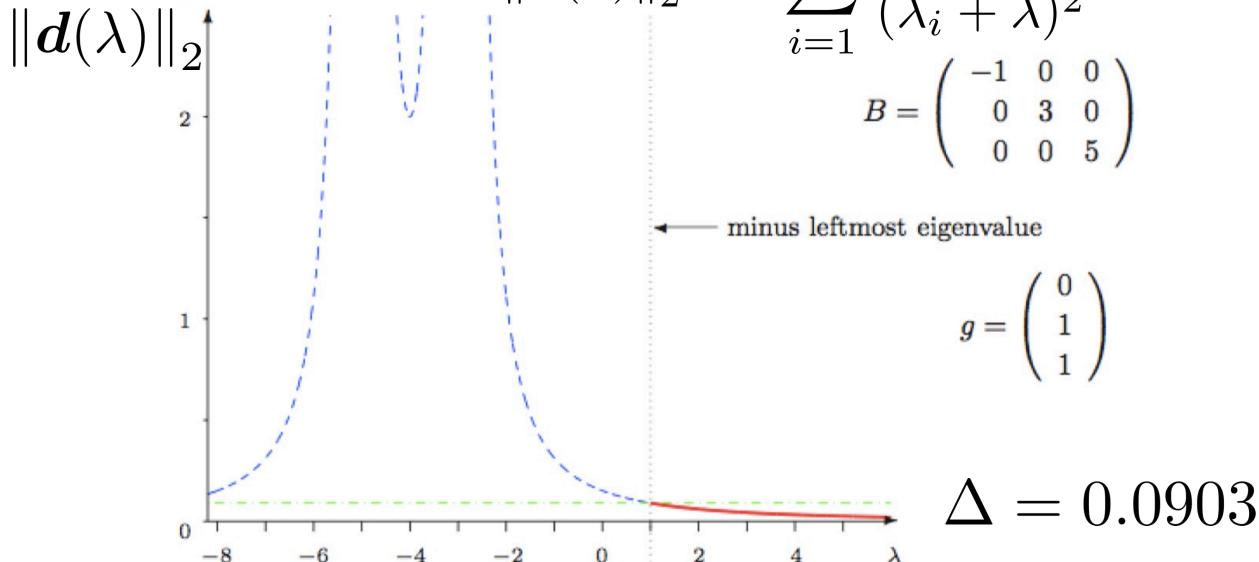
Approach I: Newton's Method

- nonconvex case: $\|d(\lambda)\|_2^2 = \sum_{i=1}^n \frac{(v_i^\top g)^2}{(\lambda_i + \lambda)^2} = \Delta^2$



Approach I: Newton's Method

- the “hard” case: $\|d(\lambda)\|_2^2 = \sum_{i=1}^n \frac{(v_i^\top g)^2}{(\lambda_i + \lambda)^2} = \Delta^2$



- $g^\top v_3 = 0$ so that no pole at $\lambda = -\lambda_3 = 1$.
- No obvious solution for $\Delta > 0.0903$. (of course, there is one!)

Approach I: Newton's Method

- the “hard” case: $\|d(\lambda)\|_2^2 = \sum_{i=1}^n \frac{(\mathbf{v}_i^\top \mathbf{g})^2}{(\lambda_i + \lambda)^2} = \Delta^2$

For indefinite \mathbf{B} , the hard case occurs when \mathbf{g} (i.e., $\nabla f(\mathbf{x})$) is orthogonal to the eigenvector \mathbf{v}_n associated with λ_n

- “Okay” if the radius Δ is “small enough”;
- No “obvious” solution when Δ is big, but in fact a solution is of the form

$$\mathbf{d}_{\lim} + \sigma \mathbf{v}_n$$

where $\mathbf{d}_{\lim} = \lim_{\lambda \rightarrow -\lambda_n} \mathbf{d}(\lambda)$, $\|\mathbf{d}_{\lim} + \sigma \mathbf{v}_n\|_2 = \Delta$

Approach II: SDP Relaxation

$$\begin{aligned} \min_{\mathbf{d} \in \mathbb{R}^n} \quad & f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \mathbf{B}_k \mathbf{d}, \\ \text{s.t.} \quad & \|\mathbf{d}\|_2 \leq \Delta_k, \end{aligned}$$

which can be *equivalently* rewritten as

$$\begin{aligned} \min_{\mathbf{d} \in \mathbb{R}^n} \quad & \frac{1}{2} \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{B}_k & \nabla f(\mathbf{x}_k) \\ \nabla f(\mathbf{x}_k)^\top & 2f(\mathbf{x}_k) \end{bmatrix} \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix} \\ \text{s.t.} \quad & \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix} \leq \Delta_k^2 + 1, \end{aligned}$$

Approach II: SDP Relaxation

$$\begin{aligned} \min_{\mathbf{d} \in \mathbb{R}^n} \quad & \frac{1}{2} \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{B}_k & \nabla f(\mathbf{x}_k) \\ \nabla f(\mathbf{x}_k)^\top & 2f(\mathbf{x}_k) \end{bmatrix} \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix} \\ \text{s.t.} \quad & \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix} \leq \Delta_k^2 + 1, \end{aligned}$$

Let $\mathbf{D} = \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix} \begin{bmatrix} \mathbf{d} \\ 1 \end{bmatrix}^\top$, then consider the convex relaxation via *lifting*

$$\begin{aligned} \min_{\mathbf{D} \in \mathbb{R}^{(n+1) \times (n+1)}} \quad & \left\langle \begin{bmatrix} \mathbf{B}_k & \nabla f(\mathbf{x}_k) \\ \nabla f(\mathbf{x}_k)^\top & 2f(\mathbf{x}_k) \end{bmatrix}, \mathbf{D} \right\rangle \\ \text{s.t.} \quad & \langle \mathbf{D}, \mathbf{I} \rangle \leq \Delta_k^2 + 1, \quad D_{n+1,n+1} = 1, \quad \mathbf{D} \succeq 0. \end{aligned}$$

Approach II: SDP Relaxation

$$\begin{aligned} \min_{\mathbf{D} \in \mathbb{R}^{(n+1) \times (n+1)}} \quad & \left\langle \begin{bmatrix} \mathbf{B}_k & \nabla f(\mathbf{x}_k) \\ \nabla f(\mathbf{x}_k)^\top & 2f(\mathbf{x}_k) \end{bmatrix}, \mathbf{D} \right\rangle \\ \text{s.t. } \langle \mathbf{D}, \mathbf{I} \rangle \leq \Delta_k^2 + 1, \quad & D_{n+1,n+1} = 1, \quad \mathbf{D} \succeq \mathbf{0}. \end{aligned}$$

- Once the optimal \mathbf{D}_* is obtained, find the direction \mathbf{d}_* via computing the leading eigenvector of \mathbf{D}_* ;
- However, solving the semidefinite programming (SDP) problem can be very *expensive* !

Further Readings

- *Numerical Optimization*. Jorge Nocedal and Stephen J. Wright.
(Chapter 4)
- *Trust Region Methods*. Andrew R. Conn, Nicholas I.M. Gould, and Philippe L. Toint. MOS-SIAM Series on Optimization, 2000.

Lecture Agenda

- **Trust-Region Method**
 - Algorithmic Introduction
 - Solving the TRM Subproblem
 - **Convergence for Strict Saddle Function**
- Cubic Regularization of Newton's Method

Assumptions

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x})$$

- The nonlinear function $f(\boldsymbol{x}) : \mathbb{R}^n \mapsto \mathbb{R}$ is twice continuously differentiable;
- The gradient is L_1 -Lipschitz continuous

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}')\|_2 \leq L_1 \|\boldsymbol{x} - \boldsymbol{x}'\|_2, \quad \forall \boldsymbol{x}, \boldsymbol{x}' \in \mathbb{R}^n.$$

- The Hessian is L_2 -Lipschitz continuous

$$\|\nabla^2 f(\boldsymbol{x}) - \nabla^2 f(\boldsymbol{x}')\|_2 \leq L_2 \|\boldsymbol{x} - \boldsymbol{x}'\|_2, \quad \forall \boldsymbol{x}, \boldsymbol{x}' \in \mathbb{R}^n.$$

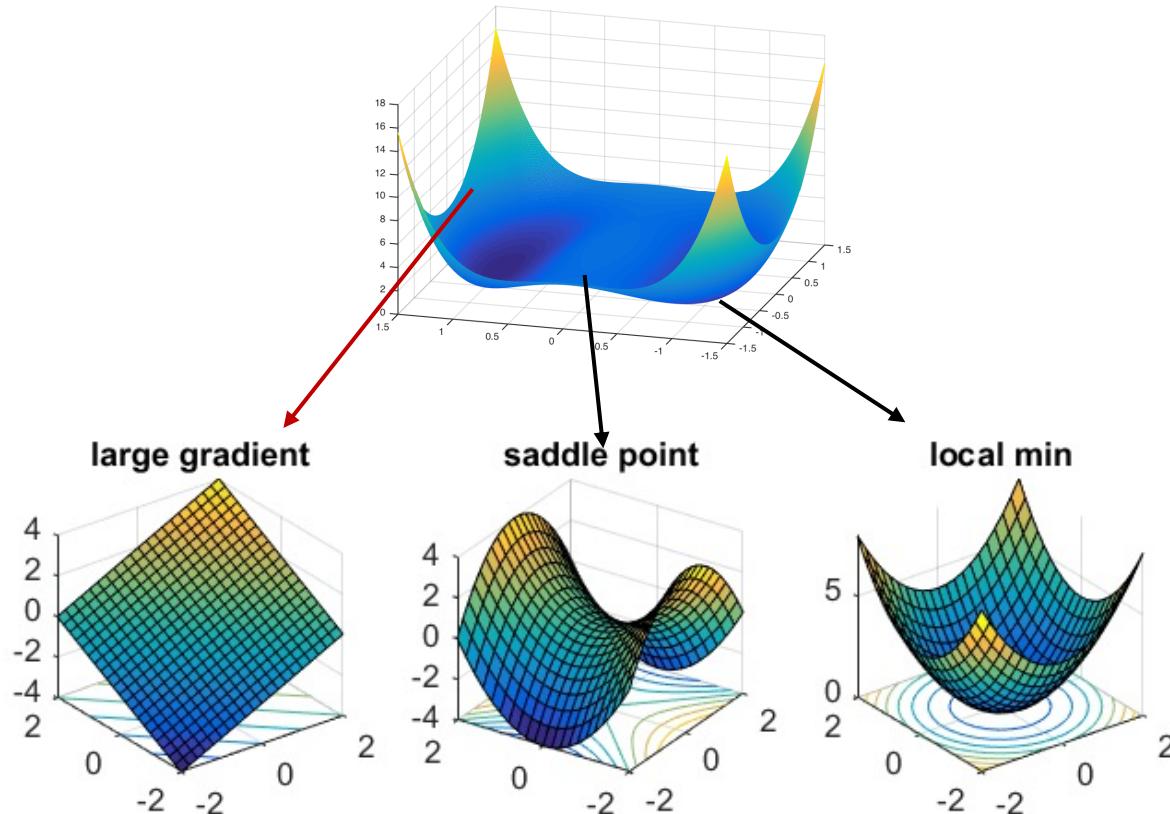
Strict Saddle Function in \mathbb{R}^n

Definition. (Strict Saddle Function in \mathbb{R}^n , Ge et al.'15)

A function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is $(\alpha, \beta, \gamma, \delta)$ -strict saddle, if $\forall x \in \mathbb{R}^n$ obeys *at least one* of the following:

- **[Large gradient]** $\|\nabla f(x)\|_2 \geq \beta$;
- **[Negative curvature]** $\exists v \in \mathbb{S}^{n-1}$, such that
$$v^\top \nabla^2 f(x) v \leq -\alpha;$$
- **[Strong convexity around minimizers]**
 $\exists x_\star$ such that $\|x - x_\star\|_2 \leq \delta$, and for all $y \in \mathcal{B}(x_\star, 2\delta)$, we have $\nabla^2 f(y) \succeq \gamma I$.

Strict Saddle Function in \mathbb{R}^n



Convergence to 2nd Order Critical Points

- For strict saddle function, we show that the iterates produced by TRM converges to a point \boldsymbol{x}_\star that satisfies

$$\nabla f(\boldsymbol{x}_\star) = \mathbf{0}, \quad \nabla^2 f(\boldsymbol{x}_\star) \succeq \mathbf{0}.$$

- For analysis, we divide \mathbb{R}^n into three regions

$$\mathcal{R}_g := \{\boldsymbol{x} \in \mathbb{R}^n \mid \|\nabla f(\boldsymbol{x})\|_2 \geq \beta\}$$

$$\mathcal{R}_n := \{\boldsymbol{x} \in \mathbb{R}^n \mid \lambda_{\min}(\nabla^2 f(\boldsymbol{x})) \leq -\alpha\}$$

$$\mathcal{R}_s := \{\boldsymbol{x} \in \mathbb{R}^n \mid \nabla^2 f(\boldsymbol{x}) \succeq \gamma \mathbf{I}\}$$

Constant Decrease in Gradient Region

Proposition 1 (Constant Decrease in \mathcal{R}_g)

Suppose $\mathbf{x}_k \in \mathcal{R}_g := \{\mathbf{x} \in \mathbb{R}^n \mid \|\nabla f(\mathbf{x})\|_2 \geq \beta\}$,
and the radius of TRM subproblem satisfies

$$\Delta \leq \min \left\{ \frac{\beta}{L_1}, \sqrt{\frac{3\beta}{4L_2}} \right\}.$$

Then for $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$, we have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{\beta\Delta}{4}.$$

Decrease in Negative Curvature Region

Proposition 2 (Constant Decrease in \mathcal{R}_n)

Suppose $\mathbf{x}_k \in \mathcal{R}_n := \{\mathbf{x} \in \mathbb{R}^n \mid \lambda_{\min}(\nabla^2 f(\mathbf{x})) \leq -\alpha\}$,
and the radius of TRM subproblem satisfies

$$\Delta \leq \frac{3}{8} \frac{\alpha}{L_2}.$$

Then for $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$, we have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{\alpha \Delta^2}{8}.$$

Quadratic Convergence Near Local Minimizer

Proposition 3 (Local Quadratic Convergence in \mathcal{R}_s)

Suppose $x_k \in \mathcal{R}_s := \{x \in \mathbb{R}^n \mid \nabla^2 f(x) \geq \gamma I\}$, and the radius of TRM subproblem is sufficiently small.

Suppose there exists a local minimizer $x_\star \in \mathcal{R}_s$, then for $x_{k+1} = x_k + d_k$, we have

$$\|\nabla f(x_{k+1})\|_2 \leq c_1 \|\nabla f(x_k)\|_2^2$$

$$\|x_k - x_\star\|_2 \leq c_2 2^{-2^k}$$

Further Readings

- *Numerical Optimization*. Jorge Nocedal and Stephen J. Wright.
(Chapter 4)
- *Trust Region Methods*. Andrew R. Conn, Nicholas I.M. Gould, and Philippe L. Toint. MOS-SIAM Series on Optimization, 2000.

Lecture Agenda

- Trust-Region Method
 - Algorithmic Introduction
 - Solving the TRM Subproblem
 - Convergence for Strict Saddle Function
- **Cubic Regularization of Newton's Method**

Assumptions

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x})$$

- The *nonlinear* function $f(\boldsymbol{x}) : \mathbb{R}^n \mapsto \mathbb{R}$ is twice continuously differentiable;
- The gradient is L_1 -Lipschitz continuous

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}')\|_2 \leq L_1 \|\boldsymbol{x} - \boldsymbol{x}'\|_2, \quad \forall \boldsymbol{x}, \boldsymbol{x}' \in \mathbb{R}^n.$$

- The Hessian is L_2 -Lipschitz continuous

$$\|\nabla^2 f(\boldsymbol{x}) - \nabla^2 f(\boldsymbol{x}')\|_2 \leq L_2 \|\boldsymbol{x} - \boldsymbol{x}'\|_2, \quad \forall \boldsymbol{x}, \boldsymbol{x}' \in \mathbb{R}^n.$$

Recall from Proximal Method

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x})$$

Minimize the quadratic upper bound:

$$\begin{aligned}\widehat{F}_\mu(\boldsymbol{x}, \boldsymbol{x}_k) &= \underbrace{f(\boldsymbol{x}_k) + \langle \nabla f(\boldsymbol{x}_k), \boldsymbol{x} - \boldsymbol{x}_k \rangle + g(\boldsymbol{x})}_{=: \ell_F(\boldsymbol{x}, \boldsymbol{x}_k)} + \frac{1}{2\mu} \|\boldsymbol{x} - \boldsymbol{x}_k\|_2^2\end{aligned}$$

that for any $\mu \in (0, 1/L_1)$,

$$F(\boldsymbol{x}) \leq \widehat{F}_\mu(\boldsymbol{x}, \boldsymbol{x}_k)$$

Recall from Proximal Method

$$\begin{aligned}\widehat{F}_\mu(\mathbf{x}, \mathbf{x}_k) &= g(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{x} - (\mathbf{x}_k - \mu \nabla f(\mathbf{x}_k))\|_2^2 \\ &\quad + f(\mathbf{x}_k) - \frac{\mu}{2} \|\nabla f(\mathbf{x}_k)\|_2^2\end{aligned}$$

Let $\mathbf{w}_k = \mathbf{x}_k - \frac{1}{\mu} \nabla f(\mathbf{x}_k)$,

$$\begin{aligned}\mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} \left\{ g(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{w}_k\|_2^2 \right\} \\ &= \text{prox}_{\mu g}(\mathbf{w}_k).\end{aligned}$$

Recall from Proximal Method

$$G_\mu(\mathbf{x}) := \frac{1}{\mu} (\mathbf{x} - \text{prox}_{\mu g}(\mathbf{x} - \mu \nabla f(\mathbf{x}))), \text{ then}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu \cdot \underbrace{G_\mu(\mathbf{x}_k)}_{\text{proximal gradient}}.$$

- However, minimizing quadratic upper bound cannot help us escape saddle point.
- It is built upon a linear approximation, with no Hessian (negative curvature) information utilized.
- To escape strict saddle, we need second-order approximation.

Idea: Minimize the Cubic Upper Bound?

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

From Lecture 3, when $\nabla^2 f$ is L_2 -Lipschitz, we have

$$f(\mathbf{x} + \mathbf{d}) \leq \underbrace{f_Q(\mathbf{d}; \mathbf{x}) + \frac{L_2}{6} \|\mathbf{d}\|_2^3}_{\widehat{f}_Q(\mathbf{d}; \mathbf{x})},$$

with that

$$f_Q(\mathbf{d}; \mathbf{x}) := f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}) \mathbf{d}.$$

Minimize the Cubic Upper Bound?

Idea: instead of minimizing the quadratic upper bound,
minimize the *cubic upper bound*:

$$\mathbf{d}_k = \arg \min_{\mathbf{d}} \hat{f}_Q(\mathbf{d}; \mathbf{x}_k)$$

With

$$\hat{f}_Q(\mathbf{d}; \mathbf{x}) := f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}) \mathbf{d} + \frac{L_2}{6} \|\mathbf{d}\|_2^3.$$

The resulting method is known as the *Cubic Regularized Newton's Method*.

Cubic Regularized Newton's Method

Problem Class:

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^n,$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is nonconvex, and is twice continuously differentiable, with both gradient and Hessian Lipschitz continuous. We have access to the second order oracle: $\nabla f(\mathbf{x}) \in \mathbb{R}^n$ and $\nabla^2 f(\mathbf{x}) \in \mathbb{R}^{n \times n}$.

Setup: Let $\hat{f}(\mathbf{y}, \mathbf{x})$ be defined similarly as in (9.2.4):

$$\hat{f}(\mathbf{y}, \mathbf{x}) \doteq \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} (\mathbf{y} - \mathbf{x})^* \nabla^2 f(\mathbf{x}) (\mathbf{y} - \mathbf{x}) + \frac{L_2}{6} \|\mathbf{y} - \mathbf{x}\|_2^3.$$

Initialization: Set $\mathbf{x}_0 \in \mathbb{R}^n$,

Iteration: For $k = 0, 1, 2, \dots$

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{y}} \hat{f}(\mathbf{y}, \mathbf{x}_k).$$

Convergence Guarantee: \mathbf{x}_k converges with $\lim_{k \rightarrow \infty} \mu(\mathbf{x}_k) = 0$.

Relationship with Trust-Region Method

- Subproblem of cubic regularized Newton:

$$\min_{\mathbf{d} \in \mathbb{R}^n} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d} + \frac{L_2}{6} \|\mathbf{d}\|_2^3.$$

- Subproblem of trust-region method:

$$\min_{\mathbf{d} \in \mathbb{R}^n} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d}, \text{ s.t. } \|\mathbf{d}\|_2 \leq \Delta_k,$$

The cubic regularization method can be viewed as replacing the constraint $\|\mathbf{d}\|_2 \leq \Delta_k$ in the trust region method by a cubic regularization $\frac{L_2}{6} \|\mathbf{d}\|_2^3$.

Relationship with Trust-Region Method

- Subproblem of cubic regularized Newton:

$$\min_{\mathbf{d} \in \mathbb{R}^n} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d} + \frac{L_2}{6} \|\mathbf{d}\|_2^3.$$

- Subproblem of trust-region method:

$$\min_{\mathbf{d} \in \mathbb{R}^n} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d}, \text{ s.t. } \|\mathbf{d}\|_2 \leq \Delta_k,$$

If L_2 is not known, we can estimate it using backtracking linesearch similar to the proximal method

Solving the Subproblem of Cubic Regularization

- Subproblem of cubic regularized Newton:

$$\min_{\mathbf{d} \in \mathbb{R}^n} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d} + \frac{L_2}{6} \|\mathbf{d}\|_2^3.$$

- When $\nabla^2 f(\mathbf{x}_k)$ is not PSD (i.e., $\nabla^2 f(\mathbf{x}_k) \not\succeq \mathbf{0}$), the subproblem is still *nonconvex*.
- The subproblem shares a lot of similarities to TRM subproblem, where we might potentially run into “hard cases”, similar to TRM subproblem.

Solving the Subproblem of Cubic Regularization

- **Subproblem of cubic regularized Newton:**

$$\min_{\mathbf{d} \in \mathbb{R}^n} \psi(\mathbf{d}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d} + \frac{L_2}{6} \|\mathbf{d}\|_2^3.$$

➤ Nonetheless, we can still run gradient descent, with

$$\nabla \psi(\mathbf{d}) = \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) \mathbf{d} + \frac{L_2}{6} \nabla \|\mathbf{d}\|_2^3$$

$$\nabla^2 f(\mathbf{x}_k) \mathbf{d} \approx \frac{\nabla f(\mathbf{x}_k + t\mathbf{d}) - \nabla f(\mathbf{x}_k)}{t}$$

➤ Recent work showed that gradient descent with small noise can globally optimize the problem with $O(\varepsilon^{-1} \log(1/\varepsilon))$.

Convergence to 2nd Order Critical Points

Similar to TRM, the cubic regularized Newton's method converges to a 2nd-order critical point \boldsymbol{x}_\star with

$$\nabla f(\boldsymbol{x}_\star) = \mathbf{0}, \quad \nabla^2 f(\boldsymbol{x}_\star) \succeq \mathbf{0}.$$

- We measure our progress by the following quantity

$$\mu(\boldsymbol{x}) := \max \left\{ \sqrt{\frac{1}{L_2} \|\nabla f(\boldsymbol{x})\|_2}, -\frac{2}{3L_2} \lambda_{\min} (\nabla^2 f(\boldsymbol{x})) \right\}$$

- If $\mu(\boldsymbol{x}) \rightarrow 0$, then $\boldsymbol{x}_k \rightarrow \boldsymbol{x}_\star$.

Convergence to 2nd Order Critical Points

Theorem (Convergence Rate of Cubic Newton's Method)

Suppose $f(\mathbf{x})$ is bounded from below. The sequence $\{\mathbf{x}_k\}$ that

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{d}} \hat{f}_Q(\mathbf{d}; \mathbf{x}_k)$$

converges to a non-empty set of limit points \mathcal{X} . Let $\mathbf{x}_* \in \mathcal{X}$, then we further have $\lim_{k \rightarrow \infty} \mu(\mathbf{x}_k) = 0$ and for any $k \geq 1$,

$$\min_{1 \leq i \leq k} \mu(\mathbf{x}_i) \leq C \left(\frac{f(\mathbf{x}_0) - f(\mathbf{x}_*)}{k \cdot L_2} \right)^{1/3}$$

for some constant $C > 0$.

Convergence to 2nd Order Critical Points

- The fact $\mu(\mathbf{x}_k) \rightarrow 0$ implies that $\{\mathbf{x}_k\}_{k \geq 1}$ converges to a 2nd order critical point with

$$\nabla f(\mathbf{x}_*) = \mathbf{0}, \quad \nabla^2 f(\mathbf{x}_*) \succeq \mathbf{0}$$

- The bound on $\mu(\mathbf{x}_k)$ implies that

$$\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|_2 \leq O(k^{-2/3})$$

which improves over $O(k^{-1/2})$ for 1st-order gradient descent method, and it is *tight* for methods with access to 2nd-order oracle.

Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications.* John Wright, Yi Ma. **(Chapter 9.2)**
- *Cubic Regularization of Newton Method and Its Global Performance.* Yurii Nesterov, and B.T. Polyak. Mathematical Programming, 2006.
- *Gradient Descent Finds the Cubic-regularized Nonconvex Newton Step.* Yair Carmon and John Duchi, SIAM Journal on Optimization, 2019.