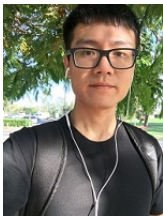# A Newton-MR Algorithm with Complexity Guarantee for Non-Convex Problems

Fred Roosta

School of Mathematics and Physics
University of Queensland

Yang Liu (UQ)

Michael Mahoney
(Berkeley)

Peng Xu (Stanford)

## Problem

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

$f : \mathbb{R}^d \to \mathbb{R}$ is

- twice differentiable
- (potentially) non-convex

---

**Algorithm** Generic 2$^{nd}$-order Method

---

Start from $\boldsymbol{x}_0$

---

---

**Algorithm** Generic 2$^{\text{nd}}$-order Method

---

Start from $\boldsymbol{x}_0$

**for** $k = 1, 2, \ldots$ **do**

**end for**

---

---

**Algorithm** Generic 2$^{nd}$-order Method

---

Start from $\boldsymbol{x}_0$

**for** $k = 1, 2, \ldots$ **do**

$$\boldsymbol{p}_k = \begin{cases} \alpha_k \boldsymbol{p} & \text{where} \quad \boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k \qquad \text{(Line Search)} \\ \\ \\ \end{cases}$$

**end for**

---

---

**Algorithm** Generic 2$^{\text{nd}}$-order Method

---

Start from $\boldsymbol{x}_0$

**for** $k = 1, 2, \ldots$ **do**

$$
\boldsymbol{p}_k = 
\begin{cases}
\alpha_k \boldsymbol{p} \quad \text{where} \quad \boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k & \text{(Line Search)} \\[2em]
\underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle / 2 & \text{(Trust Region)}
\end{cases}
$$

**end for**

---

**Algorithm** Generic 2$^{\text{nd}}$-order Method

Start from $\boldsymbol{x}_0$

**for** $k = 1, 2, \ldots$ **do**

$$
\boldsymbol{p}_k =
\begin{cases}
\alpha_k \boldsymbol{p} & \text{where} \quad \boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k & \text{(Line Search)} \\[2em]
\underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle / 2 & \text{(Trust Region)}
\end{cases}
$$

$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{p}_k$

**end for**

Article    Talk

# Conjugate gradient method

From Wikipedia, the free encyclopedia
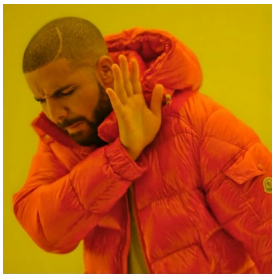
In mathematics, the **conjugate gradient method** is an algorithm for the numerical soluti implemented as an iterative algorithm, applicable to sparse systems that are too large to numerically solving partial differential equations or optimization problems.

The conjugate gradient method can also be used to solve unconstrained optimization pr and extensively researched.[4][5]
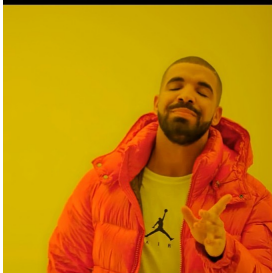
The biconjugate gradient method provides a generalization to non-symmetric matrices. \

Contents [hide]

1 Description of the problem addressed by conjugate gradients
2 Derivation as a direct method

CG

Minres

**Q: Why CG?**

# Why CG?

## Why CG?

- Let's consider the simple case of $H_k p \approx -g_k$

# Why CG?

- Let's consider the simple case of $H_k p \approx -g_k$

- When $f$ is strongly convex

## Why CG?

- Let's consider the simple case of $H_k p \approx -g_k$

- When $f$ is strongly convex $\implies H_k$ is SPD

## Why CG?

- Let's consider the simple case of $H_k p \approx -g_k$

- When $f$ is strongly convex $\implies H_k$ is SPD

- More subtly...

$$p^{(t)} = \underset{p \in \mathcal{K}_t}{\arg \min} \ \langle p, g_k \rangle + \frac{1}{2} \langle p, H_k p \rangle$$

# Why CG?

- Let's consider the simple case of $\boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k$

- When $f$ is strongly convex $\implies \boldsymbol{H}_k$ is SPD

- More subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle$$

## Why CG?

- Let's consider the simple case of $\boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k$

- When $f$ is strongly convex $\implies \boldsymbol{H}_k$ is SPD

- More subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle$$

$$\langle \boldsymbol{p}, \boldsymbol{g}_k \rangle \leq -\frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle$$

## Why CG?

- Let's consider the simple case of $\boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k$

- When $f$ is strongly convex $\implies \boldsymbol{H}_k$ is SPD

- More subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle$$

$$\langle \boldsymbol{p}, \boldsymbol{g}_k \rangle \leq -\frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle < 0$$

# Why CG?

- Let's consider the simple case of $\boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k$

- When $f$ is strongly convex $\Longrightarrow \boldsymbol{H}_k$ is SPD

- More subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle$$

$$\langle \boldsymbol{p}, \boldsymbol{g}_k \rangle \leq -\frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle < 0$$

# Why CG?

- Let's consider the simple case of $\boldsymbol{H}_k\boldsymbol{p} \approx -\boldsymbol{g}_k$

- When $f$ is strongly convex $\implies \boldsymbol{H}_k$ is SPD

- More subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p}\in\mathcal{K}_t}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \frac{1}{2}\langle \boldsymbol{p}, \boldsymbol{H}_k\boldsymbol{p}\rangle$$

$$\langle \boldsymbol{p}, \boldsymbol{g}_k \rangle \leq -\frac{1}{2}\langle \boldsymbol{p}, \boldsymbol{H}_k\boldsymbol{p}\rangle < 0$$

$\boldsymbol{p}^{(t)}$ is a descent direction for $f(\boldsymbol{x})$ <u>for all $t$</u>!

## Why CG?

- Let's consider the simple case of $\boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k$

- When $f$ is strongly convex $\implies \boldsymbol{H}_k$ is SPD

- More subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle$$

$$\langle \boldsymbol{p}, \boldsymbol{g}_k \rangle \leq -\frac{1}{2} \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle < 0$$

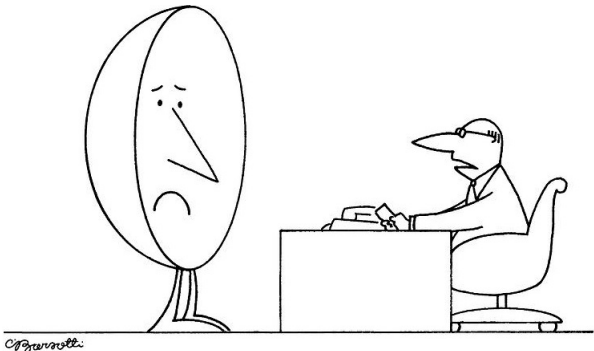$\boldsymbol{p}^{(t)}$ is a descent direction for $f(\boldsymbol{x})$ <u>for all $t$</u>!

## Classical Newton's Method



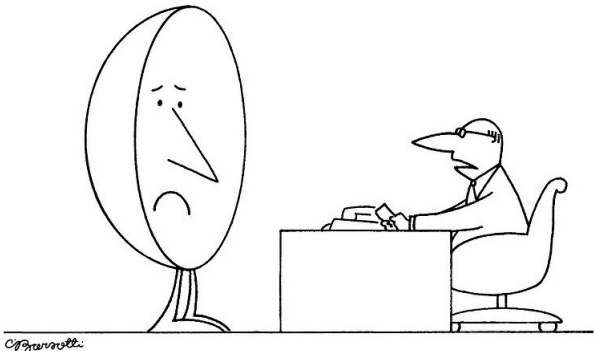*"Actually, the job calls for someone who is convex."*

## Classical Newton's Method



*"Actually, the job calls for someone who is convex."*

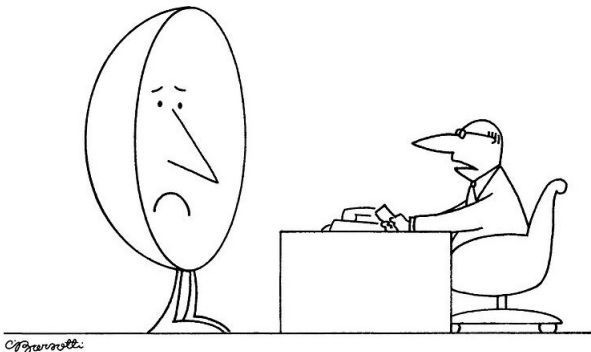But...what if the Hessian is indefinite and/or singular?

## Classical Newton's Method



*"Actually, the job calls for someone who is convex."*

But...what if the Hessian is indefinite and/or singular?

- Indefinite Hessian $\implies$ Unbounded sub-problem

## Classical Newton's Method



*"Actually, the job calls for someone who is convex."*

But...what if the Hessian is indefinite and/or singular?

- Indefinite Hessian $\implies$ Unbounded sub-problem
- $\boldsymbol{g} \notin \text{Range}(\boldsymbol{H}) \implies$ Unbounded sub-problem

$$\min_{\boldsymbol{p} \in \mathbb{R}^d} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|$$

The underlying matrix in OLS is

$$\min_{\boldsymbol{p} \in \mathbb{R}^d} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|$$

The underlying matrix in OLS is

- symmetric

$$\min_{\boldsymbol{p} \in \mathbb{R}^d} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|$$

The underlying matrix in OLS is

- symmetric

- (possibly) indefinite

$$\min_{\boldsymbol{p} \in \mathbb{R}^d} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|$$

The underlying matrix in OLS is

- symmetric

- (possibly) indefinite

- (possibly) singular

$$\min_{\boldsymbol{p} \in \mathbb{R}^d} \| \boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k \|$$

The underlying matrix in OLS is

- symmetric

- (possibly) indefinite

- (possibly) singular

- (possibly) ill-conditioned

$$\min_{\boldsymbol{p} \in \mathbb{R}^d} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|$$

The underlying matrix in OLS is

- symmetric

- (possibly) indefinite

- (possibly) singular

- (possibly) ill-conditioned

MINRES-type OLS Solvers

## Newton-MR-type Algorithms

A class of Newton-type algorithms with MINRES as sub-problem solver

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p} \in \mathcal{K}_t} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \operatorname*{arg\,min}_{\boldsymbol{p} \in \mathcal{K}_t} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p} \in \mathcal{K}_t} \frac{1}{2}\|\boldsymbol{H}_k\boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p} \in \mathcal{K}_t} \frac{1}{2}\|\boldsymbol{H}_k\boldsymbol{p} + \boldsymbol{g}_k\|^2$$

$$\left\langle \boldsymbol{p}^{(t)}, \boldsymbol{H}_k\boldsymbol{g}_k \right\rangle \leq -\frac{1}{2}\|\nabla^2 f(\boldsymbol{x}_k)\boldsymbol{p}^{(t)}\|^2$$

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

$$\left\langle \boldsymbol{p}^{(t)}, \boldsymbol{H}_k \boldsymbol{g}_k \right\rangle \leq -\frac{1}{2} \|\nabla^2 f(\boldsymbol{x}_k) \boldsymbol{p}^{(t)}\|^2 < 0$$

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p} \in \mathcal{K}_t} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p} \in \mathcal{K}_t} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

$$\left\langle \boldsymbol{p}^{(t)}, \boldsymbol{H}_k \boldsymbol{g}_k \right\rangle \leq -\frac{1}{2} \|\nabla^2 f(\boldsymbol{x}_k) \boldsymbol{p}^{(t)}\|^2 < 0$$

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \ \frac{1}{2} \|\boldsymbol{H}_k \boldsymbol{p} + \boldsymbol{g}_k\|^2$$

$$\Big\langle \boldsymbol{p}^{(t)}, \boldsymbol{H}_k \boldsymbol{g}_k \Big\rangle \leq -\frac{1}{2} \|\nabla^2 f(\boldsymbol{x}_k)\boldsymbol{p}^{(t)}\|^2 < 0$$

$\boldsymbol{p}^{(t)}$ is a descent direction for $\|\boldsymbol{g}\|^2$ <u>for all $t$!</u>

Sub-problems of MINRES:

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p}\in\mathcal{K}_t} \frac{1}{2}\|\boldsymbol{H}_k\boldsymbol{p} + \boldsymbol{g}_k\|^2$$

- There is always a solution (sometimes infinitely many)
- But more subtly...

$$\boldsymbol{p}^{(t)} = \arg\min_{\boldsymbol{p}\in\mathcal{K}_t} \frac{1}{2}\|\boldsymbol{H}_k\boldsymbol{p} + \boldsymbol{g}_k\|^2$$

$$\left\langle \boldsymbol{p}^{(t)}, \boldsymbol{H}_k\boldsymbol{g}_k \right\rangle \leq -\frac{1}{2}\|\nabla^2 f(\boldsymbol{x}_k)\boldsymbol{p}^{(t)}\|^2 < 0$$

$\boldsymbol{p}^{(t)}$ is a descent direction for $\|\boldsymbol{g}\|^2$ <u>for all $t$!</u>

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

Introduction      Newton-MR (Invex)      Newton-MR (Non-convex)   References

0000000      ○●○○○○○○○○○○○○○○○○○○○○○○○     ○○○○○○○○○○○○○○○○○○○○○○○○○

$$\min_{\boldsymbol{x} \in \mathbb{R}^d} \ \|\boldsymbol{g}\|$$

Convex

Non-Convex

Introduction
0000000

Newton-MR (Invex)
0000000000000000000000000

Newton-MR (Non-convex)
00000000000000000000000

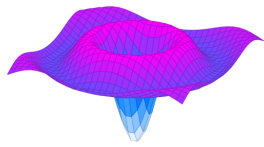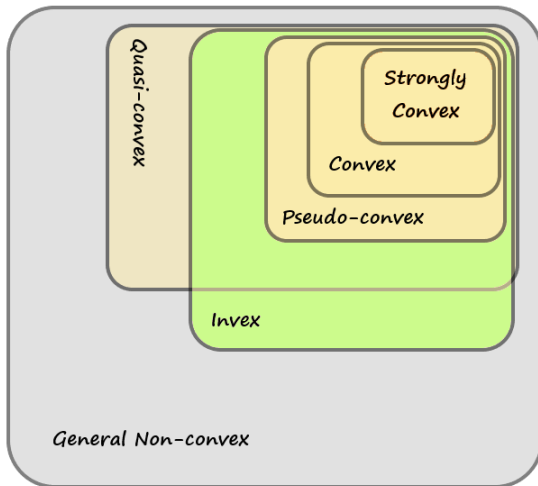References



Convex          Invex          Non-Convex

---

**Algorithm** Newton-MR (Invex)

1: **Input:** $\boldsymbol{x}_0$, $0 < \tau < 1$, $0 < \rho < 1$

2: **for** $k = 0, 1, 2, \ldots$ until $\|\boldsymbol{g}_k\| \leq \tau$ **do**

3: $\quad \boldsymbol{p}_k \approx -\boldsymbol{H}_k^\dagger \boldsymbol{g}_k$

4: $\quad$ Find $\alpha_k$ such that $\|\boldsymbol{g}_{k+1}\|^2 \leq \|\boldsymbol{g}_k\|^2 + 2\rho\alpha_k \langle \boldsymbol{p}_k, \boldsymbol{H}_k \boldsymbol{g}_k \rangle$

5: $\quad$ Update $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k$

6: **end for**

7: **Output:** $\boldsymbol{x}$ for which $\|\boldsymbol{g}_k\| \leq \tau$

---

## Examples of Convergence Results

### Global Linear Rate in "$\|\mathbf{g}\|$"

$$\left\|\mathbf{g}^{(k+1)}\right\|^2 \leq (1-\eta)\left\|\mathbf{g}_k\right\|^2, \quad 0 < \eta \leq 1.$$

### Global Linear Rate in "$f(\mathbf{x}) - \min_{\mathbf{x}} f$" Under Polyak-Łojasiewicz

$$f(\mathbf{x}_k) - \min_{\mathbf{x}} f \leq C\zeta^k, \quad 0 < \zeta \leq 1.$$

### Error Recursion with $\alpha_k = 1$

$$\min_{\mathbf{y} \in \mathcal{X}^\star} \|\mathbf{x}_{k+1} - \mathbf{y}\| \leq c_1 \min_{\mathbf{y} \in \mathcal{X}^\star} \|\mathbf{x}_k - \mathbf{y}\|^2 + \sqrt{(1-\nu)}c_2 \min_{\mathbf{y} \in \mathcal{X}^\star} \|\mathbf{x}_k - \mathbf{y}\|.$$

## Inexact Hessian

$$\tilde{H} \approx H$$

## Inexact Hessian

$$\left\| \tilde{\boldsymbol{H}} - \boldsymbol{H} \right\| \leq \epsilon$$

### Finite-sum Optimization

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\mathbf{x})$$

### Finite-sum Optimization

$$f(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{x})$$

$$\boldsymbol{H} = \frac{1}{n} \sum_{i=1}^{n} \nabla^2 f_i(\boldsymbol{x}).$$

## Finite-sum Optimization

$$f(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{x})$$

$$\tilde{\boldsymbol{H}} = \frac{1}{|\mathcal{S}|} \sum_{j \in \mathcal{S}} \nabla^2 f_j(\boldsymbol{x}),$$

### Finite-sum Optimization

$$f(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{x})$$

$$|\mathcal{S}| \in \mathcal{O}\left(\epsilon^{-2} \log\left(\frac{2d}{\delta}\right)\right)$$

## Finite-sum Optimization

$$f(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{x})$$

$$\mathbb{P}\left(\left\|\tilde{\boldsymbol{H}} - \boldsymbol{H}\right\| \le \epsilon\right) \ge 1 - \delta$$

## Newton-MR with Inexact Hessian

---

**Algorithm** Newton-MR With Inexact Hessian Information

---

1: **Input:** $\boldsymbol{x}_0$, $0 < \tau < 1$, $0 < \rho < 1$

2: **for** $k = 0, 1, 2, \ldots$ until $\|\boldsymbol{g}_k\| \leq \tau$ **do**

3:     $\boldsymbol{p}_k \approx -\tilde{\boldsymbol{H}}_k^\dagger \boldsymbol{g}_k$

4:     Find $\alpha_k$ such that $\|\boldsymbol{g}_{k+1}\|^2 \leq \|\boldsymbol{g}_k\|^2 + 2\rho\alpha_k \left\langle \boldsymbol{p}_k, \tilde{\boldsymbol{H}}_k \boldsymbol{g}_k \right\rangle$

5:     Update $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k$

6: **end for**

7: **Output:** $\boldsymbol{x}$ for which $\|\boldsymbol{g}_k\| \leq \tau$

---

**Recall: Newton's Method**

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \boldsymbol{H}_k^{-1} \boldsymbol{g}_k$$

**Recall: Newton's Method w. Inexact Hessian**

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \tilde{\boldsymbol{H}}_k^{-1} \boldsymbol{g}_k$$

## Recall: Newton's Method w. Inexact Hessian

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \tilde{\boldsymbol{H}}_k^{-1} \boldsymbol{g}_k$$

$$(1 - \tilde{\epsilon}_1)\boldsymbol{H} \preceq \tilde{\boldsymbol{H}} \preceq (1 + \tilde{\epsilon}_1)\boldsymbol{H}$$

## Recall: Newton's Method w. Inexact Hessian

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \tilde{\boldsymbol{H}}_k^{-1} \boldsymbol{g}_k$$

$$(1 - \tilde{\epsilon}_2)\boldsymbol{H}^{-1} \preceq \tilde{\boldsymbol{H}}^{-1} \preceq (1 + \tilde{\epsilon}_1)\boldsymbol{H}^{-1}$$

**Recall: Newton's Method w. Inexact Hessian**

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \tilde{\boldsymbol{H}}_k^{-1} \boldsymbol{g}_k$$

$$\left\| \tilde{\boldsymbol{H}}^{-1} - \boldsymbol{H}^{-1} \right\| \leq \tilde{\epsilon}_3$$

**Newton-MR Method**

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{H}_k^\dagger \mathbf{g}_k$$

## Newton-MR Method w. Inexact Hessian

$$x^{(k+1)} = x^{(k)} - \alpha_k \tilde{H}_k^\dagger g_k$$

## Newton-MR Method w. Inexact Hessian

$$x^{(k+1)} = x^{(k)} - \alpha_k \tilde{H}_k^\dagger g_k$$

$$\left\| \tilde{H}^\dagger - H^\dagger \right\| \leq \tilde{\epsilon}_3 \quad (?)$$

$$\lim_{\epsilon \to 0} \tilde{\boldsymbol{H}}^{\dagger} = \boldsymbol{H}^{\dagger}$$

$$\lim_{\epsilon \to 0} \tilde{\boldsymbol{H}}^{\dagger} = \boldsymbol{H}^{\dagger} \iff \mathrm{Rank}(\tilde{\boldsymbol{H}}) = \mathrm{Rank}(\boldsymbol{H})$$

[Matrix Perturbation Theory, Gilbert W. Stewart and Ji-guang Sun]

$$\lim_{\epsilon \to 0} \tilde{\boldsymbol{H}}^{\dagger} = \boldsymbol{H}^{\dagger} \iff \mathrm{Rank}(\tilde{\boldsymbol{H}}) = \mathrm{Rank}(\boldsymbol{H})$$

$$\left\| \tilde{\boldsymbol{H}}^{\dagger} - \boldsymbol{H}^{\dagger} \right\| \leq \tilde{\epsilon}_3 \quad \text{✗}$$

[Matrix Perturbation Theory, Gilbert W. Stewart and Ji-guang Sun]

$$\left\| \tilde{\boldsymbol{H}}^{\dagger} - \boldsymbol{H}^{\dagger} \right\| \leq \left( \frac{1 + \sqrt{5}}{2} \right) \max \left\{ \left\| \boldsymbol{H}^{\dagger} \right\|^{2}, \left\| \tilde{\boldsymbol{H}}^{\dagger} \right\|^{2} \right\} \epsilon$$

[Matrix Perturbation Theory, Gilbert W. Stewart and Ji-guang Sun]

$$\left\| \tilde{\boldsymbol{H}}^\dagger - \boldsymbol{H}^\dagger \right\| \le \left( \frac{1 + \sqrt{5}}{2} \right) \max \left\{ \left\| \boldsymbol{H}^\dagger \right\|^2, \left\| \tilde{\boldsymbol{H}}^\dagger \right\|^2 \right\} \epsilon$$

$$\left\| \tilde{\boldsymbol{H}}^\dagger \right\| \in o \left( \frac{1}{\sqrt{\epsilon}} \right)$$

[Matrix Perturbation Theory, Gilbert W. Stewart and Ji-guang Sun]

$$\left\|\tilde{\boldsymbol{H}}^{\dagger} - \boldsymbol{H}^{\dagger}\right\| \leq \left(\frac{1 + \sqrt{5}}{2}\right) \max\left\{\left\|\boldsymbol{H}^{\dagger}\right\|^{2}, \left\|\tilde{\boldsymbol{H}}^{\dagger}\right\|^{2}\right\} \epsilon$$

$$\left\|\tilde{\boldsymbol{H}}^{\dagger}\right\| \in o\left(\frac{1}{\sqrt{\epsilon}}\right) \quad \textcolor{red}{\text{✗}}$$

[Matrix Perturbation Theory, Gilbert W. Stewart and Ji-guang Sun]

$$\left\|\tilde{\boldsymbol{H}}^{\dagger} - \boldsymbol{H}^{\dagger}\right\| \leq \left(\frac{1 + \sqrt{5}}{2}\right) \max\left\{\left\|\boldsymbol{H}^{\dagger}\right\|^2, \left\|\tilde{\boldsymbol{H}}^{\dagger}\right\|^2\right\} \epsilon$$

$$\|\tilde{\boldsymbol{H}}^{\dagger}\| \in \mathcal{O}\left(\frac{1}{\epsilon}\right)$$

[Matrix Perturbation Theory, Gilbert W. Stewart and Ji-guang Sun]

$$\langle \boldsymbol{Hg}, \boldsymbol{p} \rangle$$

$$\langle \boldsymbol{Hg}, \boldsymbol{p} \rangle = - \left\langle \boldsymbol{Hg}, \boldsymbol{H}^{\dagger} \boldsymbol{g} \right\rangle$$

Introduction                    Newton-MR (Invex)                    Newton-MR (Non-convex)                    References
0000000                  000000000000000●00000000                  00000000000000000000000

24 / 53

$$\langle \boldsymbol{Hg}, \boldsymbol{p} \rangle = - \left\langle \boldsymbol{Hg}, \boldsymbol{H}^{\dagger}\boldsymbol{g} \right\rangle = - \left\| \boldsymbol{UU}^{\mathsf{T}}\boldsymbol{g} \right\|$$

$$\langle \boldsymbol{Hg}, \boldsymbol{p} \rangle = -\left\langle \boldsymbol{Hg}, \boldsymbol{H}^\dagger \boldsymbol{g} \right\rangle = -\left\| \boldsymbol{UU}^\mathsf{T} \boldsymbol{g} \right\|$$

$$\left\| \textcolor{red}{\tilde{\boldsymbol{U}} \tilde{\boldsymbol{U}}^\mathsf{T}} - \boldsymbol{UU}^\mathsf{T} \right\| \leq \tilde{\epsilon}_3 \quad \textbf{(?)}$$

$$\langle \boldsymbol{H}\boldsymbol{g}, \boldsymbol{p} \rangle = - \left\langle \boldsymbol{H}\boldsymbol{g}, \boldsymbol{H}^{\dagger}\boldsymbol{g} \right\rangle = - \left\| \boldsymbol{U}\boldsymbol{U}^{\mathsf{T}}\boldsymbol{g} \right\|$$

$$\left\| \tilde{\boldsymbol{U}}\tilde{\boldsymbol{U}}^{\mathsf{T}} - \boldsymbol{U}\boldsymbol{U}^{\mathsf{T}} \right\| \le \tilde{\epsilon}_3 \quad \textbf{(?)}$$

$$\mathrm{Rank}(\tilde{\boldsymbol{H}}) \ne \mathrm{Rank}(\boldsymbol{H})$$

$$\langle \boldsymbol{Hg}, \boldsymbol{p} \rangle = - \left\langle \boldsymbol{Hg}, \boldsymbol{H}^{\dagger}\boldsymbol{g} \right\rangle = - \left\| \boldsymbol{UU}^{\mathsf{T}}\boldsymbol{g} \right\|$$

$$\left\| \tilde{\boldsymbol{U}}\tilde{\boldsymbol{U}}^{\mathsf{T}} - \boldsymbol{UU}^{\mathsf{T}} \right\| \leq \tilde{\epsilon}_3 \quad \textbf{(?)}$$

$$\text{Rank}(\tilde{\boldsymbol{H}}) \neq \text{Rank}(\boldsymbol{H}) \Longrightarrow \left\| \tilde{\boldsymbol{U}}\tilde{\boldsymbol{U}}^{\mathsf{T}} - \boldsymbol{UU}^{\mathsf{T}} \right\| = 1$$

$$\langle \boldsymbol{Hg}, \boldsymbol{p} \rangle = - \left\langle \boldsymbol{Hg}, \boldsymbol{H}^{\dagger}\boldsymbol{g} \right\rangle = - \left\| \boldsymbol{UU}^{\mathsf{T}}\boldsymbol{g} \right\|$$

$$\left\| \tilde{\boldsymbol{U}}\tilde{\boldsymbol{U}}^{\mathsf{T}} - \boldsymbol{UU}^{\mathsf{T}} \right\| \leq \tilde{\epsilon}_3 \quad \textbf{✗}$$

$$\mathsf{Rank}(\tilde{\boldsymbol{H}}) \neq \mathsf{Rank}(\boldsymbol{H}) \Longrightarrow \left\| \tilde{\boldsymbol{U}}\tilde{\boldsymbol{U}}^{\mathsf{T}} - \boldsymbol{UU}^{\mathsf{T}} \right\| = 1$$

Instead of

$$\left\| \tilde{\boldsymbol{H}}\tilde{\boldsymbol{H}}^{\dagger} - \boldsymbol{H}\boldsymbol{H}^{\dagger} \right\| \le \tilde{\epsilon}_3$$

which implies

$$\left\| \left( \tilde{\boldsymbol{H}} \tilde{\boldsymbol{H}}^{\dagger} - \boldsymbol{H} \boldsymbol{H}^{\dagger} \right) \boldsymbol{v} \right\| \leq \tilde{\epsilon}_3 \left\| \boldsymbol{v} \right\|, \quad \text{for all } \boldsymbol{v}$$

we only need

$$\left\| \left( \tilde{\boldsymbol{H}} \tilde{\boldsymbol{H}}^\dagger - \boldsymbol{H} \boldsymbol{H}^\dagger \right) \boldsymbol{g} \right\| \le \tilde{\epsilon}_3 \left\| \boldsymbol{g} \right\|$$

$$\left\|\left(\tilde{\boldsymbol{H}}\tilde{\boldsymbol{H}}^{\dagger} - \boldsymbol{H}\boldsymbol{H}^{\dagger}\right)\boldsymbol{g}\right\| \le \left(\mathcal{O}(\epsilon) + \sqrt{1-\nu}\right)\|\boldsymbol{g}\|$$

$$\boldsymbol{p}_k^{(t)} \approx \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \left\| \tilde{\boldsymbol{H}}_k \boldsymbol{p} + \boldsymbol{g}_k \right\|$$

$$\boldsymbol{p}_k^{(t)} \approx \underset{\boldsymbol{p} \in \mathcal{K}_t}{\arg\min} \left\| \tilde{\boldsymbol{H}}_k \boldsymbol{p} + \boldsymbol{g}_k \right\|$$

$$\left\| \boldsymbol{p}_k^{(t)} \right\| \leq \left( \mathcal{O}(1) + \frac{\sqrt{1-\nu}}{\epsilon} \right) \|\boldsymbol{g}_k\|, \quad t = 1, 2, \ldots, \mathsf{Rank}(\tilde{\boldsymbol{H}}_k)$$

### Global Convergence: Inherent Stability

$$\|\boldsymbol{g}_{k+1}\|^2 \leq (1 - \eta + \mathcal{O}(\epsilon)) \|\boldsymbol{g}_k\|^2$$

### Local Convergence: Inherent Stability

$$\|\boldsymbol{g}(\boldsymbol{x}_{k+1})\| \leq c_1 \|\boldsymbol{g}(\boldsymbol{x}_k)\|^2 + (c_2 + \mathcal{O}(\epsilon)) \|\boldsymbol{g}(\boldsymbol{x}_k)\|$$

# Softmax-Cross Entropy: HAPT



$f(\boldsymbol{x}_k)$ vs. Iterations



$\|\boldsymbol{g}_k\|$ vs. Iterations



$f(\boldsymbol{x}_k)$ vs. Iterations



$\|\boldsymbol{g}_k\|$ vs. Iterations

DenseNet-201 with SoftPlus activation and `CIFAR100` dataset

The factors involving $1 - \nu$ have real effect!

$$f(x_1, x_2) = \frac{ax_1^2}{b - x_2}, \quad x_1 \in \mathbb{R}, \ x_2 \in (-\infty, b) \cup (b, \infty)$$

$$f(x_1, x_2) = \frac{ax_1^2}{b - x_2}, \quad x_1 \in \mathbb{R}, \; x_2 \in (-\infty, b) \cup (b, \infty)$$

$$\nu = \frac{8}{9}$$

$f(\mathbf{x}_k)$ vs. Iterations

$\|\mathbf{g}_k\|$ vs. Iterations



Step-size vs. Iterations

## Recall...

**Algorithm** Generic 2$^{nd}$-order Method

Start from $\boldsymbol{x}_0$

**for** $k = 1, 2, \ldots$ **do**

$$\boldsymbol{p}_k = \begin{cases} \alpha_k \boldsymbol{p} \quad \text{where} \quad \boldsymbol{H}_k \boldsymbol{p} \approx -\boldsymbol{g}_k & \text{(Line Search)} \\[2em] \underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g}_k \rangle + \langle \boldsymbol{p}, \boldsymbol{H}_k \boldsymbol{p} \rangle / 2 & \text{(Trust Region)} \end{cases}$$

$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{p}_k$

**end for**

Introduction
0000000

Newton-MR (Invex)
00000000000000000000000

Newton-MR (Non-convex)
00000000000000000000000

References

major iteration, we define a tolerance $\epsilon_k$ that specifies the required accuracy of t
solution. For concreteness, we choose the forcing sequence to be $\eta_k = \min(0$
to obtain a superlinear convergence rate, but other choices are possible.

**Algorithm 7.1** (Line Search Newton–CG).

Given initial point $x_0$;
**for** $k = 0, 1, 2, \ldots$
    Define tolerance $\epsilon_k = \min(0.5, \sqrt{\|\nabla f_k\|})\|\nabla f_k\|$;
    Set $z_0 = 0, r_0 = \nabla f_k, d_0 = -r_0 = -\nabla f_k$;
    **for** $j = 0, 1, 2, \ldots$
        **if** $d_j^T B_k d_j \leq 0$
            **if** $j = 0$
                **return** $p_k = -\nabla f_k$;
            **else**
                **return** $p_k = z_j$;
        Set $\alpha_j = r_j^T r_j / d_j^T B_k d_j$;
        Set $z_{j+1} = z_j + \alpha_j d_j$;
        Set $r_{j+1} = r_j + \alpha_j B_k d_j$;
        **if** $\|r_{j+1}\| < \epsilon_k$
            **return** $p_k = z_{j+1}$;
        Set $\beta_{j+1} = r_{j+1}^T r_{j+1} / r_j^T r_j$;
        Set $d_{j+1} = -r_{j+1} + \beta_{j+1} d_j$;
    **end (for)**
    Set $x_{k+1} = x_k + \alpha_k p_k$, where $\alpha_k$ satisfies the Wolfe, Goldstein, or
        Armijo backtracking conditions (using $\alpha_k = 1$ if possible);
**end**

Springer Series in
Operations Research

Jorge Nocedal
Stephen J. Wright

Numerical
Optimization
Second Edition

Springer

where $B_k = \nabla^2 f_k$. As in Algorithm 7.1, we use $d_j$ to denote the search directic
modified CG iteration and $z_j$ to d

**Algorithm 7.2** (CG–Steihaug).
  Given tolerance $\epsilon_k > 0$;
  Set $z_0 = 0, r_0 = \nabla f_k, d_0 = -r_0$
  **if** $\|r_0\| < \epsilon_k$
        return $p_k = z_0 = 0$;
  **for** $j = 0, 1, 2, \ldots$
        **if** $d_j^T B_k d_j \leq 0$
              Find $\tau$ such that $p_k$
                    and satisfies
              **return** $p_k$;
        Set $\alpha_j = r_j^T r_j / d_j^T B_k d_j$;
        Set $z_{j+1} = z_j + \alpha_j d_j$;
        **if** $\|z_{j+1}\| \geq \Delta_k$
              Find $\tau \geq 0$ such tha... $p_k = z_j + \tau d_j$ satisfies $\|p_k\| = \Delta_k$;
              **return** $p_k$;
        Set $r_{j+1} = r_j + \alpha_j B_k d_j$;
        **if** $\|r_{j+1}\| < \epsilon_k$
              **return** $p_k = z_{j+1}$;
        Set $\beta_{j+1} = r_{j+1}^T r_{j+1} / r_j^T r_j$;
        Set $d_{j+1} = -r_{j+1} + \beta_{j+1} d_j$;
  **end** (**for**).

Springer Series in
Operations Research

Jorge Nocedal
Stephen J. Wright

Numerical
Optimization

Second Edition

② Springer

# A Newton-CG algorithm with complexity guarantees for smooth unconstrained optimization

Clément W. Royer[1] · Michael O'Neill[2] · Stephen J. Wright[2]

**Abstract**
We consider minimization of a smooth nonconvex algorithm based on Newton's method and the linear explicit detection and use of negative curvature directive function. The algorithm tracks Newton-conjugate in the 1980s closely, but includes enhancements results ...
second ...
results ...

# TRUST-REGION NEWTON-CG WITH STRO COMPLEXITY GUARANTEES FOR OPTIMIZATION*

FRANK E. CURTIS†, DANIEL P. ROBINSON‡, J. WRIGHT

# A log-barrier Newton-CG method for bound constrained optimization with complexity guarantees

MICHAEL O'NEILL* AND STEPHEN J. WRIGHT
Department of Computer Sciences, University of Wisconsin-Madison, 1210 W. Dayton Street,
Madison, WI 53706, USA
*Corresponding author: moneill@cs.wisc.edu

[Received on 06 April 2019; revised on 03 December 2019]

We describe an algorithm based on a logarithmic barrier function, Newton's method and linear conjugate gradients that seeks an approximate minimizer of a smooth function over the non-negative orthant. We develop a bound on the complexity of the approach, stated in terms of the required accuracy and the cost of a single gradient evaluation of the objective function and/or a matrix-vector multiplication involving the Hessian of the objective. The approach can be implemented without explicit calculation or storage of the Hessian.

*Keywords:* nonconvex optimization; log-barrier methods; worst-case complexity; bound constraints.

## 1. Introduction

We consider the following constrained optimization problem:

$$\min f(x) \quad \text{subject to } x \geq 0, \tag{1.1}$$

where $f : \mathbb{R}^n \to \mathbb{R}$ is a nonconvex function, twice uniformly Lipschitz continuously differentiable in the interior of the non-negative orthant. We assume that explicit storage of the Hessian $\nabla^2 f(x)$ for $x \geq 0$ is undesirable, but that Hessian-vector products of the form $\nabla^2 f(x)v$ can be computed at any $x \geq 0$ for ...

# Complexity of Projected Newton Methods for Bound-constrained Optimization

Yue Xie · Stephen J. Wright

**Abstract** We analyze the iteration complexity of two methods based on the projected gradient and Newton methods for solving bound-constrained optimization problems. The first method is a scaled variant of Bertsekas's two-metric projection method [2], which can be shown to output an $\epsilon$-approximate first-order point in $\mathcal{O}(\epsilon^{-2})$ iterations. The second is a projected Newton-Conjugate Gradient (CG) method, which locates an $\epsilon$-approximate second-order point with high probability in $\mathcal{O}(\epsilon^{-3/2})$ iterations, at a cost of $\mathcal{O}(\epsilon^{-7/4})$ gradient evaluations or Hessian-vector products (omitting logarithmic factors). Besides having good complexity properties, both methods are appealing from a practical point of view, as we show using some illustrative numerical results.

**Keywords** Nonconvex Bound-constrained Optimization · Global Complexity Guarantees · Two-Metric Projection Method · Projected Newton Method

**Mathematics Subject Classification (2010)** 49M15 · 68Q25 · 90C06 ·

# Complexity of Proximal Augmented Lagrangian for Nonconvex Optimization with Nonlinear Equality Constraints

Yue Xie[1] · Stephen J. Wright[2]

**Abstract**
We analyze worst-case complexity of a Proximal augmented Lagrangian (Proximal AL) framework for nonconvex optimization with nonlinear equality constraints. When an approximate first-order (second-order) optimal point is obtained in the subproblem, an $\epsilon$ first-order (second-order) optimal point for the original problem can be guaranteed within $\mathcal{O}(1/\epsilon^{2-\eta})$ outer iterations (where $\eta$ is a user-defined parameter with $\eta \in [0, 2)$ for the first-order result and $\eta \in [1, 2]$ for the second-order result) when the proximal term coefficient $\beta$ and penalty parameter $\rho$ satisfy $\beta = \mathcal{O}(\epsilon^\eta)$ and $\rho = \Omega(1/\epsilon^\eta)$, respectively. We also investigate the total iteration complexity and operation complexity when a Newton-conjugate-gradient algorithm is used to solve the subproblems. Finally, we discuss an adaptive scheme for determining a value of the parameter $\rho$ that satisfies the requirements of the analysis.

Conjugate Gradient

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \ \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle \ /2$$

Conjugate Gradient

$$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg \min} \; \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle \, / 2$$

- **Useful for trust region:**

$$\boxed{\begin{array}{c} \fbox{Conjugate Gradient} \\[2mm] \boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \ \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle \ /2 \end{array}}$$

- **Useful for trust region:**
  - Similarity to TR's sub-problem: $\underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle \ /2$

$$\boxed{\begin{array}{c} \fbox{Conjugate Gradient} \\[1em] \boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \; \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle /2 \end{array}}$$

- **Useful for trust region:**
  - Similarity to TR's sub-problem: $\underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \; \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle /2$
  - $\left\| \boldsymbol{p}^{(t)} \right\|$ increasing with $t$

> ### Conjugate Gradient
>
> $$\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \ \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle \ /2$$

- **Useful for trust region:**
  - Similarity to TR's sub-problem: $\underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle \ /2$
  - $\left\| \boldsymbol{p}^{(t)} \right\|$ increasing with $t$

- **Useful for Newton-CG and trust-region:**

$$
\boxed{
\begin{array}{c}
\text{Conjugate Gradient} \\[2mm]
\boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \; \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle /2
\end{array}
}
$$

- **Useful for trust region:**
    - Similarity to TR's sub-problem: $\underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \; \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle /2$
    - $\left\| \boldsymbol{p}^{(t)} \right\|$ increasing with $t$

- **Useful for Newton-CG and trust-region:**
    - Negative curvature direction

$$\boxed{\begin{array}{c} \fbox{Conjugate Gradient} \\[2mm] \boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{Hp} \rangle /2 \end{array}}$$

- **Useful for trust region:**
  - Similarity to TR's sub-problem: $\underset{\|\boldsymbol{p}\| \leq \Delta}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{g} \rangle + \langle \boldsymbol{p}, \boldsymbol{Hp} \rangle /2$
  - $\left\| \boldsymbol{p}^{(t)} \right\|$ increasing with $t$

- **Useful for Newton-CG and trust-region:**
  - Negative curvature direction

**Algorithm 7.1** (Line Search Newton–CG).

Given initial point $x_0$;

**for** $k = 0, 1, 2, \ldots$

Define tolerance $\epsilon_k = \min(0.5, \sqrt{\|\nabla f_k\|}) \|\nabla f_k\|$;

Set $z_0 = 0$, $r_0 = \nabla f_k$, $d_0 = -r_0 = -\nabla f_k$;

**for** $j = 0, 1, 2, \ldots$

**if** $d_j^T B_k d_j \leq 0$

**if** $j = 0$

**return** $p_k = -\nabla f_k$;

**else**

**return** $p_k = z_j$;

Set $\alpha_j = r_j^T r_j / d_j^T B_k d_j$;

Set $z_{j+1} = z_j + \alpha_j d_j$;

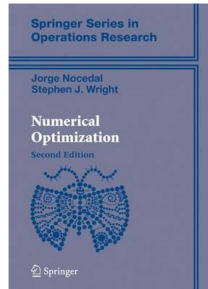Set $r_{j+1} = r_j + \alpha_j B_k d_j$;

**if** $\|r_{j+1}\| < \epsilon_k$

**return** $p_k = z_{j+1}$;

Set $\beta_{j+1} = r_{j+1}^T r_{j+1} / r_j^T r_j$;

Set $d_{j+1} = -r_{j+1} + \beta_{j+1} d_j$;

**end (for)**

Set $x_{k+1} = x_k + \alpha_k p_k$, where $\alpha_k$ satisfies the Wolfe, Goldstein, or

Armijo backtracking conditions (using $\alpha_k = 1$ if possible);

**end**

Introduction
OOOOOOOO

Newton-MR (Invex)
OOOOOOOOOOOOOOOOOOOOOOOOOOO

Newton-MR (Non-convex)
OOOOO●OOOOOOOOOOOOOOOOOO

References

**Algorithm 7.2** (CG–Steihaug).

Given tolerance $\epsilon_k > 0$;

Set $z_0 = 0, r_0 = \nabla f_k, d_0 = -r_0 = -\nabla f_k$;

**if** $\|r_0\| < \epsilon_k$

       **return** $p_k = z_0 = 0$;

**for** $j = 0, 1, 2, \ldots$

    **if** $d_j^T B_k d_j \leq 0$

        Find $\tau$ such that $p_k = z_j + \tau d_j$ minimizes $m_k(p_k)$ in (4.5)

        and satisfies $\|p_k\| = \Delta_k$;

       **return** $p_k$;

    Set $\alpha_j = r_j^T r_j / d_j^T B_k d_j$;

    Set $z_{j+1} = z_j + \alpha_j d_j$;

    **if** $\|z_{j+1}\| \geq \Delta_k$

       Find $\tau \geq 0$ such that $p_k = z_j + \tau d_j$ satisfies $\|p_k\| = \Delta_k$;

       **return** $p_k$;

    Set $r_{j+1} = r_j + \alpha_j B_k d_j$;

    **if** $\|r_{j+1}\| < \epsilon_k$

       **return** $p_k = z_{j+1}$;

    Set $\beta_{j+1} = r_{j+1}^T r_{j+1} / r_j^T r_j$;

    Set $d_{j+1} = -r_{j+1} + \beta_{j+1} d_j$;

**end** (**for**).

Springer Series in
Operations Research

Jorge Nocedal
Stephen J. Wright

Numerical
Optimization
Second Edition

Springer

Introduction
oooooooo

Newton-MR (Invex)
oooooooooooooooooooooooo

Newton-MR (Non-convex)
ooooooo●ooooooooooooooo

References

38 / 53

# CG VERSUS MINRES: AN EMPIRICAL COMPARISON[*]

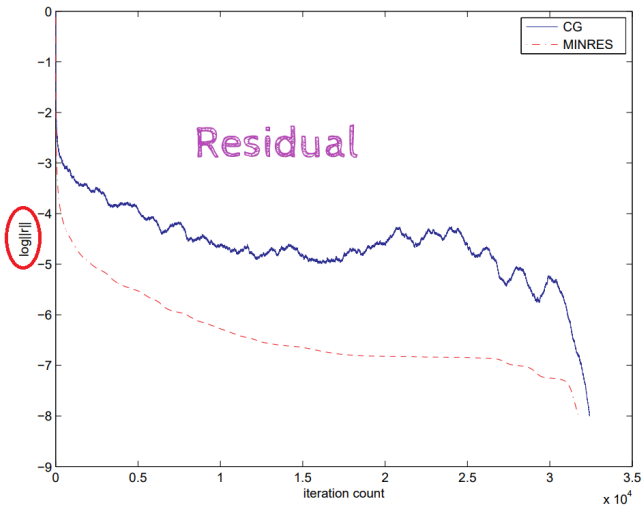### DAVID CHIN-LUNG FONG[†] AND MICHAEL SAUNDERS[‡]

**Abstract.** For iterative solution of symmetric systems $Ax = b$, the conjugate gradient method (CG) is commonly used when $A$ is positive definite, while the minimum residual method (MINRES) is typically reserved for indefinite systems. We investigate the sequence of approximate solutions $x_k$ generated by each method and suggest that even if $A$ is positive definite, MINRES may be preferable to CG if iterations are to be terminated early. In particular, we show for MINRES that the solution norms $\|x_k\|$ are monotonically increasing when $A$ is positive definite (as was already known for CG), and the solution errors $\|x^* - x_k\|$ are monotonically decreasing. We also show that the backward errors for the MINRES iterates $x_k$ are monotonically decreasing.

**Key words.** conjugate gradient method, minimum residual method, iterative method, sparse matrix, linear equations, CG, CR, MINRES, Krylov subspace method, trust-region method

**1. Introduction.** The conjugate gradient method (CG) [11] and the minimum residual method (MINRES) [18] are both Krylov subspace methods for the iterative solution of symmetric linear equations $Ax = b$. CG is commonly used when the matrix $A$ is positive definite, while MINRES is generally reserved for indefinite systems [27, p85]. We reexamine this wisdom from the point of view of early termination on positive-definite systems.

We assume that the system $Ax = b$ is real with $A$ symmetric positive definite (spd) and of dimension $n \times n$. The Lanczos process [13] with starting vector $b$ may be used to generate the $n \times k$ matrix $V_k \equiv (v_1 \quad v_2 \quad \ldots \quad v_k)$ and the $(k+1) \times k$

Fong, D.C., & Saunders, M. (2012). CG Versus MINRES: An Empirical Comparison. Sultan Qaboos University Journal for
Science, 17, 44-62.

$$\boxed{\begin{array}{c} \fbox{Minimum Residual} \\ \boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \; \left\| \boldsymbol{g} + \boldsymbol{H}\boldsymbol{p} \right\|^2 / 2 \end{array}}$$

$$\boxed{\begin{array}{c} \text{Minimum Residual} \\[4pt] \boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \ \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{g} \rangle + \left\langle \boldsymbol{p}, \boldsymbol{H}^2 \boldsymbol{p} \right\rangle / 2 \end{array}}$$

$$\boxed{\begin{array}{c} \fbox{Minimum Residual} \\ \boldsymbol{p}^{(t)} = \underset{\boldsymbol{p} \in \mathcal{K}_t(\boldsymbol{H}, -\boldsymbol{g})}{\arg\min} \; \langle \boldsymbol{p}, \boldsymbol{Hg} \rangle + \underbrace{\left\langle \boldsymbol{p}, \boldsymbol{H}^2 \boldsymbol{p} \right\rangle}_{\ddot{\frown}} / 2 \end{array}}$$

**MINRES**:

**MINRES**:

- **Negative Curvature** or **PSD Certificate**?
  (<u>without any additional work</u>)

**MINRES**:

- **Negative Curvature** or **PSD Certificate**?
  (<u>without any additional work</u>)

- **Monotonicity Properties**?

$$A \longleftarrow H, \quad b \longleftarrow -g$$

Starting from $\boldsymbol{x}_0 = \boldsymbol{0}$, we have

$$\boldsymbol{x}_t = \operatorname*{arg\,min}_{\boldsymbol{x} \in \mathcal{K}_t(\boldsymbol{A}, \boldsymbol{b})} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|$$

### Lemma (Liu and Roosta, 2021)

*As part of MINRES iterations, we can <u>readily</u> compute*

$$\frac{\langle \boldsymbol{r}_{t-1}, \boldsymbol{A}\boldsymbol{r}_{t-1}\rangle}{\langle \boldsymbol{r}_{t-1}, \boldsymbol{r}_{t-1}\rangle} = \spadesuit_{t-1} \times \clubsuit_t$$

*where* $\boldsymbol{r}_{t-1} = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}_{t-1}$

Since $r_{t-1} \in \mathcal{K}_t \left( A, b \right)$

Since $\boldsymbol{r}_{t-1} \in \mathcal{K}_t\left(\boldsymbol{A}, \boldsymbol{b}\right)$, from Lanczos Process, we get

$$\boldsymbol{A}\boldsymbol{V}_t = \boldsymbol{V}_{t+1}\widetilde{\boldsymbol{T}}_t, \quad \widetilde{\boldsymbol{T}}_t = \begin{bmatrix} \boldsymbol{T}_t \\ \beta_{t+1}\boldsymbol{e}_t^{\mathsf{T}} \end{bmatrix}$$

Since $\boldsymbol{r}_{t-1} \in \mathcal{K}_t(\boldsymbol{A}, \boldsymbol{b})$, from Lanczos Process, we get

$$\boldsymbol{A}\boldsymbol{V}_t = \boldsymbol{V}_{t+1}\widetilde{\boldsymbol{T}}_t, \quad \widetilde{\boldsymbol{T}}_t = \begin{bmatrix} \boldsymbol{T}_t \\ \beta_{t+1}\boldsymbol{e}_t^{\mathsf{T}} \end{bmatrix} \implies \boldsymbol{V}_t^{\mathsf{T}}\boldsymbol{A}\boldsymbol{V}_t = \boldsymbol{T}_t$$

Since $\boldsymbol{r}_{t-1} \in \mathcal{K}_t(\boldsymbol{A}, \boldsymbol{b})$, from Lanczos Process, we get

$$\boldsymbol{A}\boldsymbol{V}_t = \boldsymbol{V}_{t+1}\widetilde{\boldsymbol{T}}_t, \quad \widetilde{\boldsymbol{T}}_t = \begin{bmatrix} \boldsymbol{T}_t \\ \beta_{t+1}\boldsymbol{e}_t^{\mathsf{T}} \end{bmatrix} \implies \boldsymbol{V}_t^{\mathsf{T}}\boldsymbol{A}\boldsymbol{V}_t = \boldsymbol{T}_t$$

Range($\boldsymbol{V}_t$) = $\mathcal{K}_t(\boldsymbol{A}, \boldsymbol{b})$

Since $r_{t-1} \in \mathcal{K}_t(A, b)$, from Lanczos Process, we get

$$AV_t = V_{t+1}\widetilde{T}_t, \quad \widetilde{T}_t = \begin{bmatrix} T_t \\ \beta_{t+1}e_t^{\mathsf{T}} \end{bmatrix} \implies V_t^{\mathsf{T}}AV_t = T_t$$

$$\text{Range}(V_t) = \mathcal{K}_t(A, b) \implies r_{t-1} = V_t z$$

Since $r_{t-1} \in \mathcal{K}_t\left(A, b\right)$, from Lanczos Process, we get

$$AV_t = V_{t+1}\widetilde{T}_t, \quad \widetilde{T}_t = \begin{bmatrix} T_t \\ \beta_{t+1}e_t^\mathsf{T} \end{bmatrix} \implies V_t^\mathsf{T}AV_t = T_t$$

$$\text{Range}(V_t) = \mathcal{K}_t\left(A, b\right) \implies r_{t-1} = V_t z$$
$$\implies r_{t-1}^\mathsf{T}Ar_{t-1}$$

Since $r_{t-1} \in \mathcal{K}_t(A, b)$, from Lanczos Process, we get

$$AV_t = V_{t+1} \widetilde{T}_t, \quad \widetilde{T}_t = \begin{bmatrix} T_t \\ \beta_{t+1} e_t^\mathsf{T} \end{bmatrix} \implies V_t^\mathsf{T} A V_t = T_t$$

$$\text{Range}(V_t) = \mathcal{K}_t(A, b) \implies r_{t-1} = V_t z$$
$$\implies r_{t-1}^\mathsf{T} A r_{t-1} = z^\mathsf{T} V_t^\mathsf{T} A V_t z$$

Since $r_{t-1} \in \mathcal{K}_t(A, b)$, from Lanczos Process, we get

$$AV_t = V_{t+1}\widetilde{T}_t, \quad \widetilde{T}_t = \begin{bmatrix} T_t \\ \beta_{t+1}e_t^\mathsf{T} \end{bmatrix} \implies V_t^\mathsf{T} AV_t = T_t$$

$$\begin{aligned}
\text{Range}(V_t) = \mathcal{K}_t(A, b) &\implies r_{t-1} = V_t z \\
&\implies r_{t-1}^\mathsf{T} Ar_{t-1} = z^\mathsf{T} V_t^\mathsf{T} AV_t z = z^\mathsf{T} T_t z
\end{aligned}$$

Since $r_{t-1} \in \mathcal{K}_t(A, b)$, from Lanczos Process, we get

$$AV_t = V_{t+1}\widetilde{T}_t, \quad \widetilde{T}_t = \begin{bmatrix} T_t \\ \beta_{t+1}e_t^\mathsf{T} \end{bmatrix} \implies V_t^\mathsf{T} A V_t = T_t$$

$$\begin{aligned} \text{Range}(V_t) = \mathcal{K}_t(A, b) &\implies r_{t-1} = V_t z \\ &\implies r_{t-1}^\mathsf{T} A r_{t-1} = z^\mathsf{T} V_t^\mathsf{T} A V_t z = z^\mathsf{T} T_t z \end{aligned}$$

$$T_t \succ 0 \implies r_{i-1}^\mathsf{T} A r_{i-1} > 0, \ 1 \le i \le t$$

Since $r_{t-1} \in \mathcal{K}_t(A, b)$, from Lanczos Process, we get

$$AV_t = V_{t+1}\widetilde{T}_t, \quad \widetilde{T}_t = \begin{bmatrix} T_t \\ \beta_{t+1}e_t^\mathsf{T} \end{bmatrix} \implies V_t^\mathsf{T}AV_t = T_t$$

$$\text{Range}(V_t) = \mathcal{K}_t(A, b) \implies r_{t-1} = V_t z$$
$$\implies r_{t-1}^\mathsf{T}Ar_{t-1} = z^\mathsf{T}V_t^\mathsf{T}AV_t z = z^\mathsf{T}T_t z$$

$$T_t \succ 0 \implies r_{i-1}^\mathsf{T}Ar_{i-1} > 0, \ 1 \le i \le t$$

**How about the converse?**

## MINRES: Non-positive Curvature Detection and Monotonicity Properties

$$r_{i-1}^{\mathsf{T}} A r_{i-1} > 0, \ 1 \leq i \leq t \implies \begin{cases} T_i \succ 0, & 1 \leq i \leq t \\ \\ \\ \\ \\ \\ \\ \\ \\ \end{cases}$$

## MINRES: Non-positive Curvature Detection and Monotonicity Properties

$$r_{i-1}^\mathsf{T} A r_{i-1} > 0,\ 1 \leq i \leq t \implies \begin{cases} T_i \succ 0, \quad 1 \leq i \leq t \\[2em] A \succeq 0, \quad \text{if} \quad t = g(A, b) \geq \text{Rank}(A) \\[8em] \end{cases}$$

## MINRES: Non-positive Curvature Detection and Monotonicity Properties

$$
r_{i-1}^{\mathsf{T}} \boldsymbol{A} r_{i-1} > 0, \ 1 \leq i \leq t \implies
\begin{cases}
\boldsymbol{T}_i \succ \boldsymbol{0}, \quad 1 \leq i \leq t \\[2em]
\boldsymbol{A} \succeq \boldsymbol{0}, \quad \text{if} \quad t = g(\boldsymbol{A}, \boldsymbol{b}) \geq \mathrm{Rank}(\boldsymbol{A})
\end{cases}
$$

**E.g.:** Picking $\boldsymbol{b}$ uniformly at random from unit sphere guarantees w.p.1 that $t = g(\boldsymbol{A}, \boldsymbol{b}) \geq \mathrm{Rank}(\boldsymbol{A})$

## MINRES: Non-positive Curvature Detection and Monotonicity Properties

$$r_{i-1}^{\mathsf{T}} \boldsymbol{A} r_{i-1} > 0, \ 1 \leq i \leq t \implies \begin{cases} \boldsymbol{T}_i \succ \boldsymbol{0}, \quad 1 \leq i \leq t \\[2ex] \boldsymbol{A} \succeq \boldsymbol{0}, \quad \text{if} \quad t = g(\boldsymbol{A}, \boldsymbol{b}) \geq \mathrm{Rank}(\boldsymbol{A}) \\[2ex] \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{b} > \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{A} \boldsymbol{x}_t, \quad 1 \leq i \leq t \end{cases}$$

**E.g.:** Picking $\boldsymbol{b}$ uniformly at random from unit sphere guarantees w.p.1 that $t = g(\boldsymbol{A}, \boldsymbol{b}) \geq \mathrm{Rank}(\boldsymbol{A})$

## MINRES: Non-positive Curvature Detection and Monotonicity Properties

$$r_{i-1}^\mathsf{T} A r_{i-1} > 0,\ 1 \le i \le t \implies \begin{cases} T_i \succ 0, \quad 1 \le i \le t \\[2mm] A \succeq 0, \quad \text{if} \quad t = g(A, b) \ge \text{Rank}(A) \\[2mm] x_i^\mathsf{T} b > x_i^\mathsf{T} A x_t, \quad 1 \le i \le t \\[2mm] \langle x_i, A x_i \rangle / 2 - \langle b, x_i \rangle \Big\downarrow, \quad 1 \le i \le t \end{cases}$$

**E.g.:** Picking $b$ uniformly at random from unit sphere guarantees w.p.1 that $t = g(A, b) \ge \text{Rank}(A)$

## MINRES: Non-positive Curvature Detection and Monotonicity Properties

$$r_{i-1}^{\mathsf{T}} \boldsymbol{A} r_{i-1} > 0, \ 1 \leq i \leq t \implies \begin{cases} \boldsymbol{T}_i \succ \boldsymbol{0}, \quad 1 \leq i \leq t \\[2ex] \boldsymbol{A} \succeq \boldsymbol{0}, \quad \text{if} \quad t = g(\boldsymbol{A}, \boldsymbol{b}) \geq \text{Rank}(\boldsymbol{A}) \\[2ex] \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{b} > \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{A} \boldsymbol{x}_t, \quad 1 \leq i \leq t \\[2ex] \langle \boldsymbol{x}_i, \boldsymbol{A} \boldsymbol{x}_i \rangle / 2 - \langle \boldsymbol{b}, \boldsymbol{x}_i \rangle \ \downarrow, \quad 1 \leq i \leq t \\[2ex] \|\boldsymbol{x}_i\| \ \uparrow, \quad 1 \leq i \leq t \end{cases}$$

**E.g.:** Picking $\boldsymbol{b}$ uniformly at random from unit sphere guarantees w.p.1 that $t = g(\boldsymbol{A}, \boldsymbol{b}) \geq \text{Rank}(\boldsymbol{A})$

**Approximate Optimality Conditions:**

**Approximate Optimality Conditions:**

- **First-order**

$$\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$$

**Approximate Optimality Conditions:**

- **First-order**

$$\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$$

- **Second-order**

$$\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}, \quad \text{and} \quad \lambda_{\min}\left(\boldsymbol{H}_k\right) \geq -\epsilon_{\boldsymbol{H}}\boldsymbol{I}$$

We use the perturbation approach by Royer, O'Neill, and Wright, 2020

$$\boldsymbol{H} \Longleftarrow \boldsymbol{H} + \epsilon_{\boldsymbol{H}} \boldsymbol{I}$$

FULL LENGTH PAPER

Series A

**A Newton-CG algorithm with complexity guarantees for smooth unconstrained optimization**

Clément W. Royer[1] · Michael O'Neill[2] · Stephen J. Wright[2]

**Abstract**
We consider minimization of a smooth nonconvex objective function using an iterative algorithm based on Newton's method and the linear conjugate gradient algorithm, with explicit detection and use of negative curvature directions for the Hessian of the objective function. The algorithm tracks Newton-conjugate gradient procedures developed in the 1980s closely, but includes enhancements that allow worst-case complexity results to be proved for convergence to points that satisfy approximate first-order and second-order optimality conditions. The complexity results match the best known results in the literature for second-order methods.

**Algorithm** Newton-MR (Non-convex)

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

   **end for**

**Algorithm** Newton-MR (Non-convex)

  **for** $k = 1, 2, \ldots$ **do**
    **if** $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ **then**

  **end for**

**Algorithm** Newton-MR (Non-convex)

**for** $k = 1, 2, \ldots$ **do**
  **if** $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ **then**
    $\boldsymbol{g}_k \sim \mathcal{B}(\boldsymbol{0}, 1)$
  **end if**

**end for**

**Algorithm** Newton-MR (Non-convex)

**for** $k = 1, 2, \ldots$ **do**

    **if** $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ **then**

        $\boldsymbol{g}_k \sim \mathcal{B}(\boldsymbol{0}, 1)$

    **end if**

    Run MINRES to obtain $\boldsymbol{p}_k \approx \underset{\boldsymbol{p} \in \mathbb{R}^d}{\arg \min} \|(\boldsymbol{H}_k + \epsilon_{\boldsymbol{H}} \boldsymbol{I}) \boldsymbol{p} + \boldsymbol{g}_k\|$

**end for**

**Algorithm** Newton-MR (Non-convex)

**for** $k = 1, 2, \ldots$ **do**

    **if** $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ **then**

      $\boldsymbol{g}_k \sim \mathcal{B}(\boldsymbol{0}, 1)$

    **end if**

    Run MINRES to obtain $\boldsymbol{p}_k \approx \underset{\boldsymbol{p} \in \mathbb{R}^d}{\arg\min} \|(\boldsymbol{H}_k + \epsilon_{\boldsymbol{H}}\boldsymbol{I})\,\boldsymbol{p} + \boldsymbol{g}_k\|$

    **if** $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ and MINRES certifies $\boldsymbol{H}_k \geq -\epsilon_{\boldsymbol{H}}\boldsymbol{I}$ **then**

      Terminate

    **end if**

**end for**

**Algorithm** Newton-MR (Non-convex)

for $k = 1, 2, \ldots$ do

    if $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ then

        $\boldsymbol{g}_k \sim \mathcal{B}(\boldsymbol{0}, 1)$

    end if

    Run MINRES to obtain $\boldsymbol{p}_k \approx \underset{\boldsymbol{p} \in \mathbb{R}^d}{\arg\min} \|(\boldsymbol{H}_k + \epsilon_{\boldsymbol{H}} \boldsymbol{I}) \boldsymbol{p} + \boldsymbol{g}_k\|$

    if $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ and MINRES certifies $\boldsymbol{H}_k \geq -\epsilon_{\boldsymbol{H}} \boldsymbol{I}$ then

        Terminate

    end if

    Find $\alpha_k$, with the initial trial step-size $\alpha_k = 2$, such that

$$f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) < f(\boldsymbol{x}_k) - \rho \alpha_k^3 \|\boldsymbol{p}_k\|^3$$

end for

**Algorithm** Newton-MR (Non-convex)

for $k = 1, 2, \ldots$ do
   if $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ then
      $\boldsymbol{g}_k \sim \mathcal{B}(\boldsymbol{0}, 1)$
   end if
   Run MINRES to obtain $\boldsymbol{p}_k \approx \underset{\boldsymbol{p} \in \mathbb{R}^d}{\arg \min} \|(\boldsymbol{H}_k + \epsilon_{\boldsymbol{H}}\boldsymbol{I})\,\boldsymbol{p} + \boldsymbol{g}_k\|$
   if $\|\boldsymbol{g}_k\| \leq \epsilon_{\boldsymbol{g}}$ and MINRES certifies $\boldsymbol{H}_k \geq -\epsilon_{\boldsymbol{H}}\boldsymbol{I}$ then
      Terminate
   end if
   Find $\alpha_k$, with the initial trial step-size $\alpha_k = 2$, such that

$$f(\boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k) < f(\boldsymbol{x}_k) - \rho \alpha_k^3 \|\boldsymbol{p}_k\|^3$$

   $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \alpha_k \boldsymbol{p}_k$
end for

$$\overbrace{\langle \boldsymbol{p}, \boldsymbol{g} \rangle \leq - \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle /2}^{\text{CG}} \qquad\qquad \overbrace{\langle \boldsymbol{p}, \boldsymbol{g} \rangle \leq - \langle \boldsymbol{p}, \boldsymbol{H}\boldsymbol{p} \rangle}^{\text{MINRES}}$$

$$\overbrace{\langle \boldsymbol{p}, \boldsymbol{g} \rangle \leq - \langle \boldsymbol{p}, \boldsymbol{Hp} \rangle /2}^{\text{CG}}$$

$$\overbrace{\langle \boldsymbol{p}, \boldsymbol{g} \rangle \leq - \langle \boldsymbol{p}, \boldsymbol{Hp} \rangle}^{\text{MINRES}}$$

$$\alpha = 1$$

$$\alpha = 2$$

## Operation Complexity

## Operation Complexity

- **First-order**

$$\tilde{\mathcal{O}}\left(\max\left\{\epsilon_{\boldsymbol{g}}^{-3}\epsilon_{\boldsymbol{H}}^{5/2}, \epsilon_{\boldsymbol{H}}^{-7/2}\right\}\right)$$

## Operation Complexity

- **First-order**

$$\tilde{\mathcal{O}}\left(\max\left\{\epsilon_{\boldsymbol{g}}^{-3}\epsilon_{\boldsymbol{H}}^{5/2}, \epsilon_{\boldsymbol{H}}^{-7/2}\right\}\right) \stackrel{\epsilon_{\boldsymbol{H}}^{2}=\epsilon_{\boldsymbol{g}}=\epsilon}{\Longrightarrow} \tilde{\mathcal{O}}\left(\epsilon^{-7/4}\right)$$

## Operation Complexity

- **First-order**

$$\tilde{\mathcal{O}}\left(\max\left\{\epsilon_{\boldsymbol{g}}^{-3}\epsilon_{\boldsymbol{H}}^{5/2}, \epsilon_{\boldsymbol{H}}^{-7/2}\right\}\right) \stackrel{\epsilon_{\boldsymbol{H}}^{2}=\epsilon_{\boldsymbol{g}}=\epsilon}{\Longrightarrow} \tilde{\mathcal{O}}\left(\epsilon^{-7/4}\right)$$

- **Second-order**

$$\tilde{\mathcal{O}}\left(\max\{\epsilon_{\boldsymbol{H}}^{-7/2}, \epsilon_{\boldsymbol{H}}^{-1/2}\epsilon_{\boldsymbol{g}}^{-3/2}\}\right)$$

## Operation Complexity

- **First-order**

$$\tilde{\mathcal{O}} \left( \max \left\{ \epsilon_{\boldsymbol{g}}^{-3} \epsilon_{\boldsymbol{H}}^{5/2}, \epsilon_{\boldsymbol{H}}^{-7/2} \right\} \right) \overset{\epsilon_{\boldsymbol{H}}^2 = \epsilon_{\boldsymbol{g}} = \epsilon}{\Longrightarrow} \tilde{\mathcal{O}} \left( \epsilon^{-7/4} \right)$$

- **Second-order**

$$\tilde{\mathcal{O}} \left( \max \{ \epsilon_{\boldsymbol{H}}^{-7/2}, \epsilon_{\boldsymbol{H}}^{-1/2} \epsilon_{\boldsymbol{g}}^{-3/2} \} \right) \overset{\epsilon_{\boldsymbol{H}}^2 = \epsilon_{\boldsymbol{g}} = \epsilon}{\Longrightarrow} \tilde{\mathcal{O}} \left( \epsilon^{-7/4} \right)$$

Figure: None-linear Least-square problem with CIFAR10 dataset.

Figure: None-linear Least-square problem with CIFAR10 dataset.

Figure: Auto-encoder with CIFAR10 dataset.

Figure: Auto-encoder with CIFAR10 dataset.

Figure: Performance Profile on 252 CUTEst Problems

📄 Roosta, Fred et al. (2018). "Newton-MR: Inexact Newton Method With Minimum Residual Sub-problem Solver". In: *arXiv preprint arXiv:1810.00303*.

📄 Royer, Clément W, Michael O'Neill, and Stephen J Wright (2020). "A Newton-CG algorithm with complexity guarantees for smooth unconstrained optimization". In: *Mathematical Programming* 180.1, pp. 451–488.

📄 Liu, Yang and Fred Roosta (2021a). "A Newton-MR Algorithm With Complexity Guarantee Non-convex Optimization". In: In preparation.

📄 — (2021b). "Convergence of Newton-MR under Inexact Hessian Information". In: *SIAM Journal on Optimization* 31.1, pp. 59–90.

📄 — (2021c). "MINRES: From Negative Curvature Detection to Monotonicity Properties". In: Submitted.