

Newton-MR: Newton's Method Without Smoothness or Convexity

Michael W. Mahoney

ICSI and Department of Statistics
University of California at Berkeley

Joint work with Fred Roosta, Yang Liu, and Peng Xu

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

Newton's Method

Classical Newton's Method

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \underbrace{\alpha_k}_{\text{step-size}} \underbrace{[\nabla^2 f(\mathbf{x}^{(k)})]^{-1} \nabla f(\mathbf{x}^{(k)})}_{\text{Newton Direction}}$$



First Order Methods

Classical Gradient Descent

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$$



Machine Learning  First Order Methods...



But why 1st?

Q: But Why 1st Order Methods?

- Cheap Iterations
- Easy To Implement
- “Good” Worst-Case Complexities
- Good Generalization

But why **Not** 2nd?

Q: But Why **Not** 2nd Order Methods?

- ~~Cheap~~ Expensive Iterations
- ~~Easy~~ Hard To Implement
- ~~“Good”~~ “Bad” Worst-Case Complexities
- ~~Good~~ Bad Generalization

Our Goal...

Goal: Improve 2nd Order Methods...

- Cheap ~~Expensive~~ Iterations
- Easy ~~Hard~~ To Use
- “Good” ~~“Bad”~~ Average(?) -Case Complexities
- Good ~~Bad~~ Generalization

Our Goal..

Any Other Advantages?

- Effective Iterations \Rightarrow Less Iterations \Rightarrow Less Communications
- Saddle Points For Non-Convex Problems
- Less Sensitive to Parameter Tuning
- Less Sensitive to Initialization

Achilles' heel for most 2nd-order methods is...



Achilles' heel: Solving the **Sub-problems!!!**

Sub-Problems

- Trust Region:

$$\mathbf{s}^{(k)} = \arg \min_{\|\mathbf{s}\| \leq \Delta_k} \left\langle \mathbf{s}, \nabla f(\mathbf{x}^{(k)}) \right\rangle + \frac{1}{2} \left\langle \mathbf{s}, \nabla^2 f(\mathbf{x}^{(k)}) \mathbf{s} \right\rangle$$

- Cubic Regularization:

$$\mathbf{s}^{(k)} = \arg \min_{\mathbf{s} \in \mathbb{R}^d} \left\langle \mathbf{s}, \nabla f(\mathbf{x}^{(k)}) \right\rangle + \frac{1}{2} \left\langle \mathbf{s}, \nabla^2 f(\mathbf{x}^{(k)}) \mathbf{s} \right\rangle + \frac{\sigma_k}{3} \|\mathbf{s}\|^3$$

Newton's Method

Recall: Classical **Newton's method**

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \underbrace{[\nabla^2 f(\mathbf{x}^{(k)})]^{-1} \nabla f(\mathbf{x}^{(k)})}_{\text{Linear System}}$$



$$\nabla^2 f(\mathbf{x}^{(k)}) \mathbf{p} = -\nabla f(\mathbf{x}^{(k)})$$

We know how to solve "**Ax = b**" **very well!**

Newton-CG

f : **Strongly** Convex \implies Newton-**CG** $\implies \nabla^2 f(\mathbf{x}^{(k)})\mathbf{p} \approx -\nabla f(\mathbf{x}^{(k)})$

$$\mathbf{p} \approx \operatorname{argmin}_{\mathbf{p} \in \mathbb{R}^d} \left\langle \mathbf{p}, \nabla f(\mathbf{x}^{(k)}) \right\rangle + \frac{1}{2} \left\langle \mathbf{p}, \nabla^2 f(\mathbf{x}^{(k)})\mathbf{p} \right\rangle$$

Why CG?

- f is **strongly** convex $\implies \nabla^2 f(\mathbf{x}^{(k)})$ is SPD
- More subtly...

$$\mathbf{p}^{(t)} = \operatorname{argmin}_{\mathbf{p} \in \mathcal{K}_t} \left\langle \mathbf{p}, \nabla f(\mathbf{x}^{(k)}) \right\rangle + \frac{1}{2} \left\langle \mathbf{p}, \nabla^2 f(\mathbf{x}^{(k)}) \mathbf{p} \right\rangle$$



$$\left\langle \mathbf{p}, \nabla f(\mathbf{x}^{(k)}) \right\rangle \leq -\frac{1}{2} \left\langle \mathbf{p}, \nabla^2 f(\mathbf{x}^{(k)}) \mathbf{p} \right\rangle < 0$$



$\mathbf{p}^{(t)}$ is a **descent** direction for f for all t !

Classical Newton's Method

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

- **Indefinite** Hessian \implies Unbounded sub-problem
- **Singular** Hessian and $\nabla f(\mathbf{x}) \notin \text{Range}(\nabla^2 f(\mathbf{x})) \implies$ Unbounded sub-problem
 - $\nabla^2 f(\mathbf{x})\mathbf{p} = -\nabla f(\mathbf{x})$ has no solution

strong convexity \implies linear system sub-problems

~~strong convexity~~ \implies ~~linear system~~ sub-problems

$$\underbrace{\mathbf{Ax} = \mathbf{b}}_{\text{Linear System}} \implies \underbrace{\|\mathbf{Ax} - \mathbf{b}\|}_{\text{Least Squares}}$$

$$\min_{\mathbf{p} \in \mathbb{R}^d} \left\| \overbrace{\nabla^2 f(\mathbf{x}_k)}^{\mathbf{A}} \overbrace{\mathbf{p}}^{\mathbf{x}} + \overbrace{\nabla f(\mathbf{x}_k)}^{-\mathbf{b}} \right\|$$

The underlying matrix in OLS is

- symmetric
- (possibly) indefinite
- (possibly) singular
- (possibly) ill-conditioned

MINRES-type OLS Solvers \implies MINRES-QLP [Choi et al., 2011]

Sub-problems of MINRES:

$$\mathbf{p}^{(t)} = \operatorname{argmin}_{\mathbf{p} \in \mathcal{K}_t} \frac{1}{2} \|\nabla^2 f(\mathbf{x}_k) \mathbf{p} + \nabla f(\mathbf{x}_k)\|^2$$

- There is always a solution (sometimes infinitely many)
-

$$\mathbf{p}^{(t)} = \operatorname{argmin}_{\mathbf{p} \in \mathcal{K}_t} \frac{1}{2} \|\nabla^2 f(\mathbf{x}_k) \mathbf{p} + \nabla f(\mathbf{x}_k)\|^2$$



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



$\mathbf{p}^{(t)}$ is a **descent** direction for $\|\nabla f(\mathbf{x})\|^2$ for all t !

Newton-MR vs. Newton-CG

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}_k$$

Newton-CG:

$$\mathbf{p}_k \approx \operatorname{argmin}_{\mathbf{p} \in \mathbb{R}^d} \langle \mathbf{g}_k, \mathbf{p} \rangle + \frac{1}{2} \langle \mathbf{p}, \mathbf{H}_k \mathbf{p} \rangle = -[\mathbf{H}_k]^{-1} \mathbf{g}_k$$

$$\alpha_k : f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \leq f(\mathbf{x}_k) + \alpha_k \beta \langle \mathbf{p}_k, \mathbf{g}_k \rangle$$

Newton-MR:

$$\mathbf{p}_k \approx \operatorname{argmin}_{\mathbf{p} \in \mathbb{R}^d} \|\mathbf{H}_k \mathbf{p} + \mathbf{g}_k\|^2 = -[\mathbf{H}_k]^\dagger \mathbf{g}_k$$

$$\alpha_k : \|\mathbf{g}(\mathbf{x}_k + \alpha_k \mathbf{p}_k)\|^2 \leq \|\mathbf{g}_k\|^2 + 2\alpha_k \beta \langle \mathbf{p}_k, \mathbf{H}_k \mathbf{g}_k \rangle$$

Newton-MR vs. Newton-CG

	Newton-CG	Newton-MR
Sub-problems	$\min_{\mathbf{p} \in \mathcal{K}_t} \frac{1}{2} \langle \mathbf{p}, \mathbf{H}\mathbf{p} \rangle + \langle \mathbf{p}, \mathbf{g} \rangle$	$\min_{\mathbf{p} \in \mathcal{K}_t} \ \mathbf{H}\mathbf{p} + \mathbf{g}\ ^2$
Line Search	$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \alpha \rho \langle \mathbf{p}_k, \mathbf{g}_k \rangle$	$\ \mathbf{g}_{k+1}\ ^2 \leq \ \mathbf{g}_k\ ^2 + 2\alpha \rho \langle \mathbf{p}_k, \mathbf{H}_k \mathbf{g}_k \rangle$

Invexity

Invexity

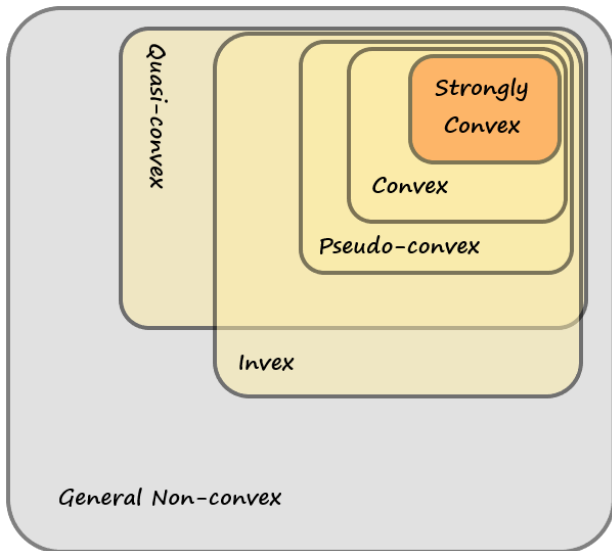
$$f(\mathbf{y}) - f(\mathbf{x}) \geq \langle \phi(\mathbf{y}, \mathbf{x}), \nabla f(\mathbf{x}) \rangle$$

- Necessary and sufficient for optimality: $\nabla f(\mathbf{x}) = 0$
- E.g.: Convex $\implies \phi(\mathbf{y}, \mathbf{x}) = \mathbf{y} - \mathbf{x}$
-

$$\left. \begin{array}{l} g : \mathbb{R}^p \rightarrow \mathbb{R} \text{ is differentiable and convex} \\ \mathbf{h} : \mathbb{R}^d \rightarrow \mathbb{R}^p \text{ has full-rank Jacobian } (p \leq d) \end{array} \right\} \implies g \circ \mathbf{h} \text{ is invex}$$

- “Global optimality” of stationary points in deep residual networks [Bartlett et al., 2018]

Strong Convexity \subsetneq Invexity



Newton-MR vs. Newton-CG

	Newton-CG	Newton-MR
Sub-problems	$\min_{\mathbf{p} \in \mathcal{K}_t} \frac{1}{2} \langle \mathbf{p}, \mathbf{H}\mathbf{p} \rangle + \langle \mathbf{p}, \mathbf{g} \rangle$	$\min_{\mathbf{p} \in \mathcal{K}_t} \ \mathbf{H}\mathbf{p} + \mathbf{g}\ ^2$
Line Search	$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \alpha\rho \langle \mathbf{p}_k, \mathbf{g}_k \rangle$	$\ \mathbf{g}_{k+1}\ ^2 \leq \ \mathbf{g}_k\ ^2 + 2\alpha\rho \langle \mathbf{p}_k, \mathbf{H}_k \mathbf{g}_k \rangle$
Problem class	Strongly Convex	Invex

Moral Smoothness

(Recall) Typical Smoothness Assumptions:

Lipschitz Gradient: $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L_g \|\mathbf{x} - \mathbf{y}\|$

Lipschitz Hessian: $\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L_H \|\mathbf{x} - \mathbf{y}\|$

These smoothness assumptions
are *stronger* than
what is required for first-order methods.

Moral Smoothness

Moral-Smoothness

Let $\mathcal{X}_0 \triangleq \{\mathbf{x} \in \mathbb{R}^d \mid \|\nabla f(\mathbf{x})\| \leq \|\nabla f(\mathbf{x}_0)\|\}$. For any $\mathbf{x}_0 \in \mathbb{R}^d$, there is a constant $0 < L(\mathbf{x}_0) < \infty$, such that $\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X}_0 \times \mathbb{R}^d$, we have

$$\|\nabla^2 f(\mathbf{y})\nabla f(\mathbf{y}) - \nabla^2 f(\mathbf{x})\nabla f(\mathbf{x})\| \leq L(\mathbf{x}_0) \|\mathbf{y} - \mathbf{x}\|.$$

Smoothness \subsetneq Moral Smoothness

Moral Smoothness

Hessian of the quadratically smoothed hinge-loss is **not continuous**.

$$f(\mathbf{x}) = \frac{1}{2} \max \left\{ 0, b \langle \mathbf{a}, \mathbf{x} \rangle \right\}^2$$

But it satisfies moral-smoothness with $L = b^4 \|\mathbf{a}\|^4$.

Newton-MR vs. Newton-CG

	Newton-CG	Newton-MR
Sub-problems	$\min_{\mathbf{p} \in \mathcal{K}_t} \frac{1}{2} \langle \mathbf{p}, \mathbf{H}\mathbf{p} \rangle + \langle \mathbf{p}, \mathbf{g} \rangle$	$\min_{\mathbf{p} \in \mathcal{K}_t} \ \mathbf{H}\mathbf{p} + \mathbf{g}\ ^2$
Line Search	$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \alpha\rho \langle \mathbf{p}_k, \mathbf{g}_k \rangle$	$\ \mathbf{g}_{k+1}\ ^2 \leq \ \mathbf{g}_k\ ^2 + 2\alpha\rho \langle \mathbf{p}_k, \mathbf{H}_k \mathbf{g}_k \rangle$
Problem class	Strongly Convex	Invex
Smoothness	H & g	Hg

Null-Space Property

For any $\mathbf{x} \in \mathbb{R}^d$, let

- \mathbf{U}_x be an orthogonal basis for $\text{Range}(\nabla^2 f(\mathbf{x}))$
- \mathbf{U}_x^\perp be its orthogonal complement

Gradient-Hessian Null-Space Property

$$\left\| \left(\mathbf{U}_x^\perp \right)^T \nabla f(\mathbf{x}) \right\|^2 \leq \left(\frac{1-\nu}{\nu} \right) \left\| \mathbf{U}_x^T \nabla f(\mathbf{x}) \right\|^2, \quad \forall \mathbf{x} \in \mathbb{R}^d, \quad 0 < \nu \leq 1$$

- Strictly convex $f(\mathbf{x})$: $\nu = 1$
- Non-convex $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{a}_i^T \mathbf{x})$: $\nu = 1$
- Some fractional programming: $\nu = 8/9$
- Some non-linear composition of functions $f(\mathbf{x}) = g(\mathbf{h}(\mathbf{x}))$

Inexactness

- **Newton-CG** [Roosta and Mahoney, Mathematical Programming, 2018]

$$\|\mathbf{H}_k \mathbf{p}_k + \mathbf{g}_k\| \leq \theta \|\mathbf{g}_k\| \implies \theta \leq 1/\sqrt{\kappa}$$

- **Newton-MR** [Roosta, Liu, Xu and Mahoney, arXiv, 2019]

$$\langle \mathbf{H}_k \mathbf{p}_k, \mathbf{g}_k \rangle \leq -(1 - \theta) \|\mathbf{g}_k\|^2 \implies 1 - \nu \leq \theta < 1$$

Examples of Convergence Results

Global Linear Rate in “ $\|\mathbf{g}\|$ ”

$$\|\mathbf{g}^{(k+1)}\|^2 \leq \left(1 - \frac{4\rho(1-\rho)\gamma^2(1-\theta)^2}{L(\mathbf{x}_0)}\right) \|\mathbf{g}_k\|^2.$$

Global Linear Rate in “ $f(\mathbf{x}) - f^*$ ” Under Polyak-Łojasiewicz

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

Error Recursion with $\alpha_k = 1$ Under Error Bound

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

Newton-MR vs. Newton-CG

	Newton-CG	Newton-MR
Sub-problems	$\min_{\mathbf{p} \in \mathcal{K}_t} \frac{1}{2} \langle \mathbf{p}, \mathbf{H}\mathbf{p} \rangle + \langle \mathbf{p}, \mathbf{g} \rangle$	$\min_{\mathbf{p} \in \mathcal{K}_t} \ \mathbf{H}\mathbf{p} + \mathbf{g}\ ^2$
Line Search	$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \alpha\rho \langle \mathbf{p}_k, \mathbf{g}_k \rangle$	$\ \mathbf{g}_{k+1}\ ^2 \leq \ \mathbf{g}_k\ ^2 + 2\alpha\rho \langle \mathbf{p}_k, \mathbf{H}_k \mathbf{g}_k \rangle$
Problem class	Strongly Convex	Invex
Smoothness	H & g	Hg
Inexactness	$\ \mathbf{H}\mathbf{p} + \mathbf{g}\ \leq \theta \ \mathbf{g}\ $ $\theta < 1/\sqrt{\kappa}$	$\langle \mathbf{p}, \mathbf{H}\mathbf{g} \rangle \leq -(1 - \theta) \ \mathbf{g}\ $ $\theta < 1$
Metric / Rate	$\ \mathbf{g}\ $: R-linear $f(\mathbf{x}) - f^*$: Q-Linear	$\ \mathbf{g}\ $: Q-linear $f(\mathbf{x}) - f^*$: R-Linear (GPL)

Newton-MR vs. Newton-CG for $\min f(\mathbf{x})$

&

MINRES vs. CG for $\mathbf{Ax} = \mathbf{b}$

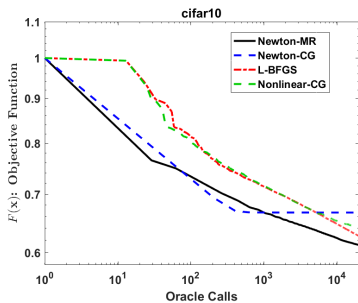
Newton-MR vs. Newton-CG

	$\min f(\mathbf{x})$	
	Newton-CG	Newton-MR
Sub-problems	$\min_{\mathbf{p} \in \mathcal{K}_t} \frac{1}{2} \langle \mathbf{p}, \mathbf{H}\mathbf{p} \rangle + \langle \mathbf{p}, \mathbf{g} \rangle$	$\min_{\mathbf{p} \in \mathcal{K}_t} \ \mathbf{H}\mathbf{p} + \mathbf{g}\ ^2$
Problem class	Strongly Convex	Invex
Metric / Rate	$\ \mathbf{g}\ $: R-linear $f(\mathbf{x}) - f^*$: Q-Linear	$\ \mathbf{g}\ $: Q-linear $f(\mathbf{x}) - f^*$: R-Linear (GPL)

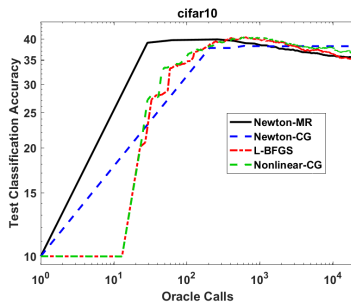
MINRES vs. CG

	$Ax = b$	
	CG	MINRES
Sub-problems	$\min_{x \in \mathcal{K}_t} \frac{1}{2} \langle x, Ax \rangle + \langle x, b \rangle$	$\min_{x \in \mathcal{K}_t} \ Ax - b\ ^2$
Problem class	Symmetric Positive Definite	Symmetric
Metric / Rate	$\ Ax - b\ $: R-linear $\ x - x^*\ _A$: Q-Linear	$\ Ax - b\ $: Q-linear $\ x - x^*\ _A$: R-Linear (SPD)

Weakly-Convex ($n = 50,000, d = 27,648$): Softmax-Cross Entropy

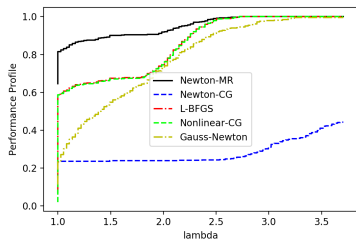
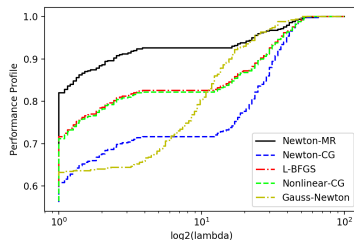


(a) $f(\mathbf{x}_k)$



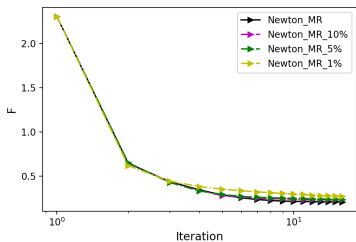
(b) Test Accuracy

Non-Convex: GMM

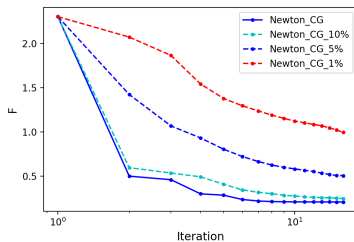
(c) $f(\mathbf{x}_k)$ 

(d) Estimation error

Weakly-Convex ($n = 50,000, d = 7,056$): Softmax-Cross Entropy

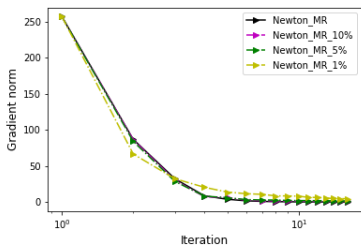


(e) $f(\mathbf{x}_k)$

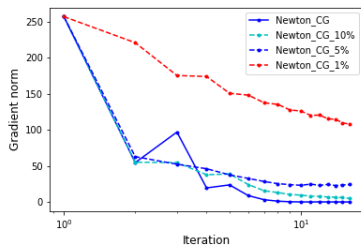


(f) $f(\mathbf{x}_k)$

Weakly-Convex ($n = 50,000, d = 7,056$): Softmax-Cross Entropy

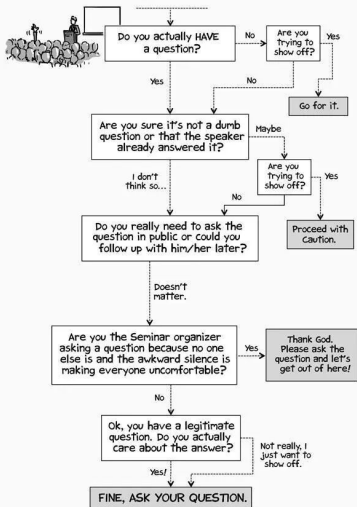


(g) $\|\nabla f(\mathbf{x}_k)\|$



(h) $\|\nabla f(\mathbf{x}_k)\|$

Should you ask a Question during Seminar?



THANK YOU!