

Second-Order Optimization with Lazy Hessians

Nikita Doikov

El Mahdi Chayti

Martin Jaggi



Problem and Motivation

We want to solve **unconstrained minimization** problem:

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

- f is differentiable and can be **non-convex**
- First-order gradient methods: **cheap** to implement, but **slow** rates
- Second-order methods (Newton's Method): **fast** rates but **expensive**

This work: we propose to use a previously seen Hessian for several iterations (*lazy Hessian updates*):

- **Provable improvement** of the total arithmetic complexity

Newton's Method with Cubic Regularization

Assume that the **Hessian is Lipschitz Continuous**:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^d$$

\Rightarrow **global upper model** of the objective, for $H \geq L$:

$$f(y) \leq \Omega(x; y) + \frac{H}{6}\|y - x\|^3, \quad \forall x, y \in \mathbb{R}^d,$$

where

$$\Omega(x; y) \stackrel{\text{def}}{=} f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \langle \nabla^2 f(x)(y - x), y - x \rangle$$

Cubic Newton Method [1]. Iterate, $k \geq 0$:

$$x_{k+1} = \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ M_H(x_k; y) \equiv \Omega(x_k; y) + \frac{H}{6}\|y - x_k\|^3 \right\}$$

Theorem. Let $H := L$. Then, to find $\|\nabla f(\bar{x}_k)\| \leq \varepsilon$, the Cubic Newton needs

$$K = \mathcal{O}\left(\frac{\sqrt{L}(f(x_0) - f^*)}{\varepsilon^{3/2}}\right).$$

- For the Gradient Method, we need $\mathcal{O}(1/\varepsilon^2)$ iterations
- We also can prove convergence to a **second-order stationary point** for the Cubic Newton: $\nabla^2 f(\bar{x}_k) \succeq -\sqrt{L}\varepsilon I$
- **Adaptive strategy** for H : ensure $f(x_{k+1}) \leq M_H(x_k; x_{k+1})$ [1, 2] — the method becomes **universal**, adapting automatically to the most appropriate problem class [3, 4]

Solving the Cubic Subproblem

How to compute one step? $h^+ = \operatorname{argmin}_{h \in \mathbb{R}^d} \left\{ \langle g, h \rangle + \frac{1}{2} \langle Ah, h \rangle + \frac{H}{6} \|h\|^3 \right\}$

- **Step 1:** compute **factorization** of $A = U\Lambda U^\top$, where U is orthonormal basis $UU^\top = I$, and Λ is **diagonal** or **tridiagonal** — $\mathcal{O}(d^3)$ arithmetic operations (the most expensive part)
- **Step 2:** solve the **dual problem** (concave univariate maximization):

$$\max_{\tau \in \mathbb{R}: \tau > [-\lambda_{\min}]_+} \left\{ -\frac{1}{2} \langle (\Lambda + \tau I)^{-1} \bar{g}, \bar{g} \rangle - \frac{\tau^4}{3H^2} \tau^3 \right\} \Rightarrow h^+ = -(\Lambda + \tau^* I)^{-1} \bar{g}$$

Lazy Hessian Updates

Main Idea: use the same Hessian for $m \geq 1$ iterations

Lazy Hessian Updates: compute new Hessian once per m iterations.

Hessians:	$\nabla^2 f(x_0)$	reuse Hessian \rightarrow		$\nabla^2 f(x_m)$	reuse Hessian \rightarrow	
Gradients:	$\nabla f(x_0)$	$\nabla f(x_1)$...	$\nabla f(x_{m-1})$	$\nabla f(x_m)$	$\nabla f(x_{m+1})$...

- Appeared first in [5]

Algorithm: Cubic Newton with Lazy Hessians

Define step of the method with Hessian at some previous point z :

$$T_H(x, z) \stackrel{\text{def}}{=} \operatorname{argmin}_{y \in \mathbb{R}^d} \left\{ \langle \nabla f(x), y - x \rangle + \frac{1}{2} \langle \nabla^2 f(z)(y - x), y - x \rangle + \frac{H}{6} \|y - x\|^3 \right\}$$

Define $\pi(k) \stackrel{\text{def}}{=} k - k \bmod m$

Cubic Newton with Lazy Hessians

Iterate, $k \geq 0$:

1. Set last snapshot point $z_k = x_{\pi(k)}$
2. Compute lazy cubic step $x_{k+1} = T_H(x_k, z_k)$

Theory: Convergence Rate

Theorem. Let $H := 6mL$. Then, to find $\|\nabla f(\bar{x})\| \leq \varepsilon$, the method needs

$$K = \mathcal{O}\left(\frac{\sqrt{mL}(f(x_0) - f^*)}{\varepsilon^{3/2}}\right)$$

- Worse than the full Cubic Newton by the factor \sqrt{m}

Note: the total number of **Hessian updates** during these steps is

$$\frac{K}{m} = \mathcal{O}\left(\frac{\sqrt{L}(f(x_0) - f^*)}{\sqrt{m}\varepsilon^{3/2}}\right)$$

Choice of m ?

\gg **Optimize the total cost:**

$$\text{Arithmetic complexity} = K \times \text{GradCost} + \frac{K}{m} \times \text{HessCost}$$

In many problems: $\boxed{\text{HessCost} = d \times \text{GradCost}} \Rightarrow m := d$

- Generalized Linear Models (**Logistic Regression**)
- Log-sum-exp (**Soft Max**)
- **Neural Networks:** computing $\nabla^2 f(x)h$ is the same cost as $\nabla f(x)$, for any x, h , by using backpropagation. Then

$$\nabla^2 f(x) = [\nabla^2 f(x)e_1 \mid \dots \mid \nabla^2 f(x)e_d]$$

Total Arithmetic Complexity

• Gradient Method:

$$\mathcal{O}\left(\frac{L_1(f(x_0) - f^*)}{\varepsilon^2}\right) \times \text{GradCost}$$

• Full Cubic Newton:

$$\mathcal{O}\left(\frac{\sqrt{L_2}(f(x_0) - f^*)}{\varepsilon^{3/2}}\right) \times \text{GradCost} \times d$$

• Lazy Cubic Newton ($m := d$):

$$\mathcal{O}\left(\frac{\sqrt{L_2}(f(x_0) - f^*)}{\varepsilon^{3/2}}\right) \times \text{GradCost} \times \sqrt{d}$$

This provably improves the total arithmetic complexity of the Cubic Newton by factor \sqrt{d}

Locally, we also have superlinear convergence

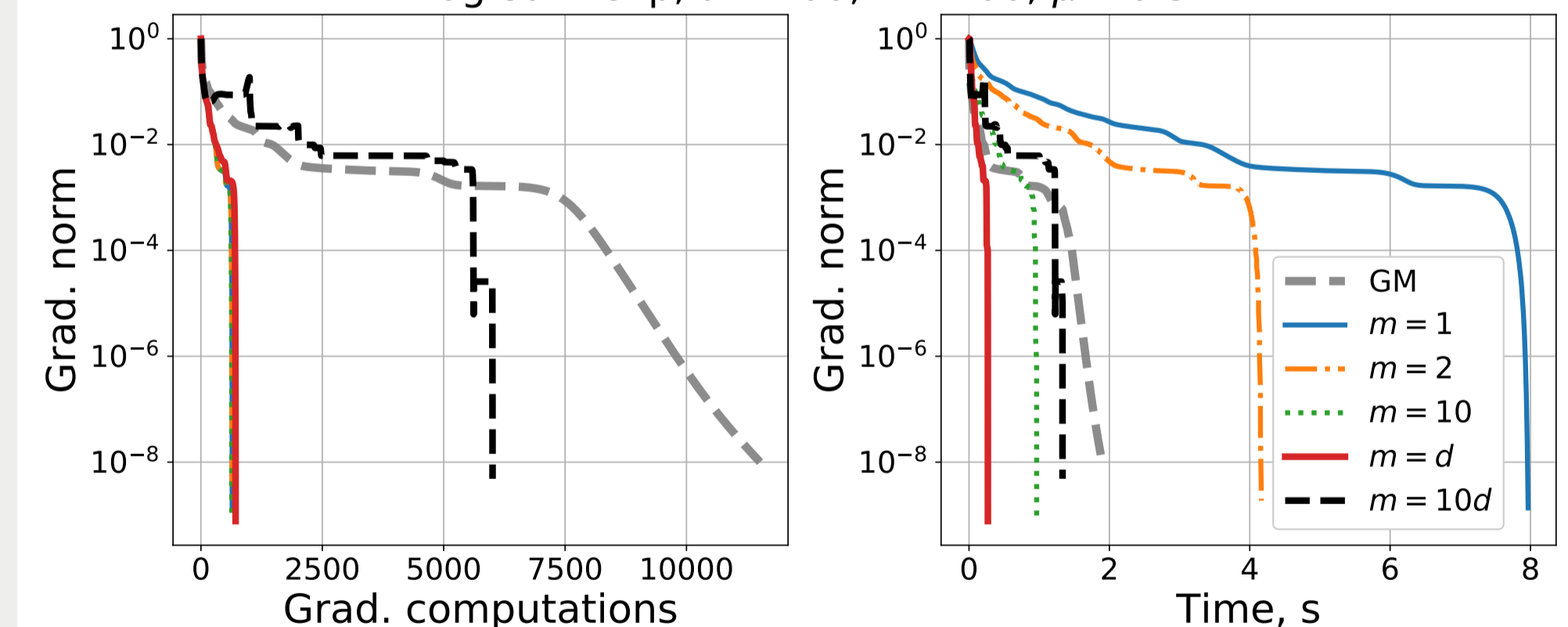
For **convex problems**, we can use the Gradient Regularization technique with lazy Hessian updates, achieving the same global rates:

$$T_H(x, z) = x - \left(\nabla^2 f(z) + \sqrt{H} \|\nabla f(x)\| I \right)^{-1} \nabla f(x)$$

Experiment: Soft Max

$$\min_{x \in \mathbb{R}^d} f(x) = \mu \ln \left(\sum_{i=1}^n \exp\left(\frac{\langle a_i, x \rangle - b_i}{\mu}\right) \right) \approx \max_{1 \leq i \leq n} [\langle a_i, x \rangle - b_i]$$

Log-sum-exp, $d = 100$, $n = 100$, $\mu = 0.5$



References

- [1] Yurii Nesterov and Boris Polyak. "Cubic regularization of Newton's method and its global performance". In: *Mathematical Programming* 108.1 (2006), pp. 177–205
- [2] Coralia Cartis, Nicholas IM Gould, and Philippe L Toint. "Adaptive cubic regularisation methods for unconstrained optimization. Part I: motivation, convergence and numerical results". In: *Mathematical Programming* 127.2 (2011), pp. 245–295
- [3] Geovani N Grapiglia and Yurii Nesterov. "Regularized Newton Methods for Minimizing Functions with Hölder Continuous Hessians". In: *SIAM Journal on Optimization* 27.1 (2017), pp. 478–506
- [4] Nikita Doikov and Yurii Nesterov. "Minimizing uniformly convex functions by cubic regularization of Newton method". In: *Journal of Optimization Theory and Applications* (2021), pp. 1–23
- [5] VE Shamanskii. "A modification of Newton's method". In: *Ukrainian Mathematical Journal* 19.1 (1967), pp. 118–122