

Complexity of Cubically Regularized Newton Method for Minimizing Uniformly Convex Functions

Nikita Doikov

Yurii Nesterov

Catholic University of Louvain, Belgium

FGS-19, Nice
September 18, 2019

1. Introduction: first-order nondegeneracy
2. Regularization of Newton method
3. Universal Cubic Newton
4. Numerical examples

1. Introduction: first-order nondegeneracy
2. Regularization of Newton method
3. Universal Cubic Newton
4. Numerical examples

Optimization problem

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

- ▶ Problem class: $\mu I \preceq \nabla^2 f(x) \preceq LI, \quad \forall x \in \mathbb{R}^n$
- ▶ Gradient Method:
$$x_{k+1} = x_k - \frac{1}{\alpha_k} \nabla f(x_k)$$
$$= \operatorname{argmin}_{y \in \mathbb{R}^n} \left\{ f(x_k) + \langle \nabla f(x_k), y - x_k \rangle + \frac{\alpha_k}{2} \|y - x_k\|^2 \right\}$$
- ▶ Linear convergence rate. Set $\alpha_k = L$. Then
$$f(x_k) - f(x_{k+1}) \geq \frac{1}{2L} \|\nabla f(x_k)\|^2 \geq \frac{\mu}{L} (f(x_k) - f^*).$$

Review: First order Nondegeneracy

- ▶ Problem class: $\mu I \preceq \nabla^2 f(x) \preceq LI, \quad \forall x \in \mathbb{R}^n.$

The **condition number** of f :

$$\kappa \equiv \frac{L}{\mu} \leq 1.$$

- ▶ Iteration complexity. $f(x_N) - f^* \leq \varepsilon$

Gradient Method: $N = O\left(\frac{1}{\kappa} \log \frac{f(x_0) - f^*}{\varepsilon}\right).$

Fast Gradient Method: $N = O\left(\sqrt{\frac{1}{\kappa}} \log \frac{f(x_0) - f^*}{\varepsilon}\right).$

Second Order Methods: ?

Damped Newton method

$$\begin{aligned}x_{k+1} &= x_k - \frac{1}{\alpha_k} (\nabla^2 f(x_k))^{-1} \nabla f(x_k) \\ &= \operatorname{argmin}_{y \in \mathbb{R}^n} \left\{ f(x_k) + \langle \nabla f(x_k), y - x_k \rangle + \frac{\alpha_k}{2} \langle \nabla^2 f(x_k)(y - x_k), y - x_k \rangle \right\}\end{aligned}$$

- ▶ Global convergence. Set $\alpha_k = L/\mu$.

$$\begin{aligned}f(x_k) - f(x_{k+1}) &\geq \frac{\mu}{2L} \langle (\nabla^2 f(x_k))^{-1} \nabla f(x_k), \nabla f(x_k) \rangle \\ &\geq \frac{\mu}{2L^2} \|\nabla f(x_k)\|^2 \geq \frac{\mu^2}{L^2} (f(x_k) - f^*).\end{aligned}$$

- ▶ Complexity: $N = O\left(\frac{1}{\kappa^2} \log \frac{f(x_0) - f^*}{\varepsilon}\right)$. It is worse than GM!

Our work: we use **Regularization of Newton method** [Nesterov & Polyak 06; Grapiglia & Nesterov 17] and show global linear rate of convergence which is **faster** than for Gradient Methods.

1. Introduction: first-order nondegeneracy
2. Regularization of Newton method
3. Universal Cubic Newton
4. Numerical examples

Functions with Hölder continuous Hessian

For $\nu \in [0, 1]$ Hessian of f is Hölder continuous of degree ν on a convex set $C \subseteq \text{dom } f$ iff

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq \mathcal{H}_f(\nu) \|y - x\|^\nu, \quad \forall x, y \in C.$$

- ▶ $\nu = 1$. Lipschitz continuous Hessian:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq \mathcal{H}_f(1) \|y - x\|$$

- ▶ $\nu = 0$. Hessian with bounded variation:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq \mathcal{H}_f(0)$$

Note: for $\mu I \preceq \nabla^2 f(x) \preceq LI$ we have

$$\mathcal{H}_f(0) \leq L - \mu.$$

Regularization of Newton iterations

For $\nu \in [0, 1]$:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq \mathcal{H}_f(\nu) \|y - x\|^\nu, \quad \forall x, y \in \mathcal{C}.$$

Bound for Quadratic approximation:

► $|f(y) - Q(x; y)| \leq \frac{\mathcal{H}_f(\nu)}{(1+\nu)(2+\nu)} \|y - x\|^{2+\nu}$, where

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

Regularized Newton step:

$$T_{\nu, H}(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{y \in \mathcal{C}} \left\{ Q(x; y) + \frac{H}{(1+\nu)(2+\nu)} \|y - x\|^{2+\nu} \right\}$$

Iterations: $x_{k+1} = T_{\nu, H_k}(x_k)$, $k \geq 0$.

[Nesterov & Polyak 06; Grapiglia & Nesterov 17]

Uniformly convex functions

$f : \text{dom } f \rightarrow \mathbb{R}$ is called **uniformly convex** of degree $p \geq 2$ on a convex set $C \subseteq \text{dom } f$ iff

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\sigma_f(p)}{p} \|x - y\|^p, \quad \forall x, y \in C.$$

$\sigma_f(p)$ – constant of uniform convexity.

- ▶ Strongly convex functions: $p = 2$, $\mu = \sigma_f(2)$.
- ▶ **Example:**

$$f(x) = \frac{1}{p} \|x - x_0\|^p, \quad p \geq 2$$

is uniformly convex of degree p with constant $\sigma_f(p) = 2^{2-p}$.

- ▶ Sum of convex and uniformly convex functions gives uniformly convex.

How to measure nondegeneracy?

$$\min_x f(x)$$

- ▶ Hölder Hessian of degree $\nu \in [0, 1]$:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq \mathcal{H}_f(\nu) \|x - y\|^\nu, \quad \forall x, y \in \mathcal{C}.$$

- ▶ Uniformly convex of degree $p \geq 2$:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\sigma_f(p)}{p} \|x - y\|^p, \quad \forall x, y \in \mathcal{C}.$$

Which ν and p to use?

$$O\left((\mathcal{H}_f(\nu))^\alpha (\sigma_f(p))^\beta \log \frac{f(x_0) - f^*}{\varepsilon}\right) \quad \# \text{ of iterations}$$

Physical dimensions: $[f]$ and $[x]$. Then

$$[\mathcal{H}_f(\nu)] = \frac{[f]}{[x]^{2+\nu}}, \quad [\sigma_f(p)] = \frac{[f]}{[x]^p}, \quad \left[(\mathcal{H}_f(\nu))^\alpha (\sigma_f(p))^\beta\right] = \frac{[f]^{\alpha+\beta}}{[x]^{\alpha(2+\nu)+\beta p}}.$$

$$\alpha + \beta = 0 \text{ and } \alpha(2 + \nu) + \beta p = 0 \quad \Rightarrow \quad \boxed{p = 2 + \nu}$$

Linear rate of regularized Newton

Fix $\nu \in [0, 1]$. Model of the objective:

$$M_{\nu, H}(x; y) \stackrel{\text{def}}{=} Q(x; y) + \frac{H}{(1+\nu)(2+\nu)} \|y - x\|^{2+\nu}$$

Regularized Newton step: $T_{\nu, H}(x) \stackrel{\text{def}}{=} \underset{y \in C}{\operatorname{argmin}} M_{\nu, H}(x; y)$.

Iterations: $x_{k+1} = T_{\nu, H_k}(x_k), \quad k \geq 0$

New result: Theorem 1

- ▶ Let $0 < \mathcal{H}_f(\nu) < +\infty$.
- ▶ Let $\sigma_f(2 + \nu) > 0$.
- ▶ Let $0 \leq H_k \leq \beta \mathcal{H}_f(\nu)$ and $f(x_{k+1}) \leq M_{\nu, H_k}(x_k; x_{k+1}), \forall k$.

Then: $f(x_K) - f^* \leq \varepsilon$ for

$$K = O\left(\max\left\{1, \left(\frac{\mathcal{H}_f(\nu)}{\sigma_f(2+\nu)}\right)^{\frac{1}{1+\nu}}\right\} \cdot \log \frac{f(x_0) - f^*}{\varepsilon}\right).$$

Condition number of degree ν .

Denote

$$\gamma_f(\nu) \stackrel{\text{def}}{=} \frac{\sigma_f(2+\nu)}{\mathcal{H}_f(\nu)}$$

second-order **condition number** of degree $\nu \in [0, 1]$.

Properties (unbounded case, $\text{diam } C = +\infty$):

1. $\gamma_f(\nu) \leq \frac{1}{1+\nu}$, $\nu \in (0, 1]$.
2. Let for some $\nu > 0$ we have $\gamma_f(\nu) > 0$. Then

$$\gamma_f(\nu') = 0 \quad \forall \nu' \in [0, 1] \setminus \{\nu\}.$$

In practice we may not know exact value for ν . Universal methods?

1. Introduction: first-order nondegeneracy
2. Regularization of Newton method
3. Universal Cubic Newton
4. Numerical examples

Composite optimization problem

$$F^* = \min_{x \in \text{dom } F} F(x) \stackrel{\text{def}}{=} f(x) + h(x).$$

- ▶ f is a smooth part, $\gamma_f(\nu) > 0$ for some $\nu \in [0, 1]$.
- ▶ $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a composite part (possibly nonsmooth but *simple*).

Cubically regularized (composite) Newton:

$$x_{k+1} = T_{H_k}(x_k), \quad k \geq 0$$

for

$$T_H(x) \stackrel{\text{def}}{=} \underset{y \in \text{dom } F}{\text{argmin}} M_H(x; y)$$

where $M_H(x; y) \stackrel{\text{def}}{=} Q(x; y) + \frac{H}{6} \|y - x\|^3 + h(y)$.

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

$$M_H^*(x) \stackrel{\text{def}}{=} M_H(x; T_H(x))$$

Algorithm 1 Adaptive cubically regularized Newton method

Initialization: Choose $x_0 \in \text{dom } F$, $H_0 > 0$.

Iterations: $k \geq 0$.

- 1: Find minimum integer $i_k \geq 0$ s.t. it holds

$$F(T_{H_k 2^{i_k}}(x_k)) \leq M_{H_k 2^{i_k}}^*(x_k).$$

- 2: Make the Cubic step $x_{k+1} := T_{H_k 2^{i_k}}(x_k)$.
- 3: Set $H_{k+1} := H_k 2^{i_k - 1}$.

[Nesterov-Polyak 06; Cartis-Gould-Toint 11; Grapiglia-Nesterov 17]

New result: Theorem 2

- ▶ Let $\gamma_f(\nu) > 0$ for some $\nu \in [0, 1]$.

Then $F(x_K) - F^* \leq \varepsilon$ for

$$K = O\left(\max\{1, (\gamma_f(\nu))^{\frac{-1}{1+\nu}}\} \cdot \log \frac{F(x_0) - F^*}{\varepsilon}\right).$$

Total number of oracle calls N_K during K iterations is bounded as:

$$N_K \leq 2K + \log_2 \frac{\kappa_f(\nu)}{\varepsilon^{(1-\nu)/(2+\nu)} H_0}$$

The algorithm **does not** use real $\nu \in [0, 1]$ and adapt to **the best possible bound**:

$$O\left(\max\{1, \inf_{\nu \in [0, 1]} (\gamma_f(\nu))^{\frac{-1}{1+\nu}}\} \cdot \frac{F(x_0) - F^*}{\varepsilon}\right).$$

Local convergence

Classic Newton method achieves **local quadratic convergence** under the conditions

1. Strong convexity: $\sigma_f(2) > 0$.
2. Hessian is Lipschitz continuous: $\mathcal{H}_f(1) < +\infty$.
3. Initial point x_0 *close* to the optimum x^* .

For **Adaptive cubically regularized Newton** we have:

$$F(x_{k+1}) - F^* \leq \frac{9\mathcal{H}_f^2(1)}{2\sigma_f^3(2)} (F(x_k) - F^*)^2, \quad \forall k.$$

And the **region of quadratic convergence** is

$$\mathcal{Q} = \left\{ x \in \text{dom } F : F(x) - F^* \leq \frac{2\sigma_f^3(2)}{9\mathcal{H}_f^2(1)} \right\}.$$

Number of iterations to find ϵ -solution:

$$O\left(\log \log \frac{\sigma_f^3(2)}{\mathcal{H}_f^2(1)\epsilon}\right).$$

The same order as for Classic Newton.

Cubic Newton vs. Gradient Method

Composite minimization problem: $\min_{x \in \mathbb{R}^n} F(x) \stackrel{\text{def}}{=} f(x) + h(x)$.

- ▶ Problem class: $\mu I \preceq \nabla^2 f(x) \preceq LI$. **Gradient Method** needs $O\left(\frac{L}{\mu} \log \frac{F(x_0) - F^*}{\varepsilon}\right)$ iterations to find ε -solution.
- ▶ Problem class: for some $\nu \in [0, 1]$ it holds $\gamma_f(\nu) > 0$.

Adaptive Cubic Newton needs

$O\left(\max\{1, \inf_{\nu \in [0,1]} (\gamma_f(\nu))^{\frac{-1}{1+\nu}}\} \log \frac{F(x_0) - F^*}{\varepsilon}\right)$ iterations to find ε -solution, **without knowledge of ν** . (Until entering \mathcal{Q}).

Note:

$$(\gamma_f(0))^{-1} \leq \frac{L-\mu}{\mu} < \frac{L}{\mu}.$$

Consider $\tilde{f}(x) = f(x) + \frac{1}{2}\langle Ax, x \rangle$. Then

$$\tilde{L} = L + \lambda_{\max}(A) \quad \text{but} \quad (\gamma_{\tilde{f}}(\nu))^{-1} \leq (\gamma_f(\nu))^{-1}.$$

\Rightarrow Cubic Newton is not affected by any quadratic parts of F .

Cubic Newton vs. Classic Newton

Composite minimization problem: $\min_{x \in \mathbb{R}^n} F(x) \stackrel{\text{def}}{=} f(x) + h(x)$.

- ▶ Arithmetical complexity of every iteration: $O(n^3)$.
- ▶ Strongly convex functions with Lipschitz Hessian:
local quadratic convergence.
- ▶ Strongly convex functions with bounded Hessian,

$$\mu I \preceq \nabla^2 f(x) \preceq LI :$$

$O\left(\left(\frac{L}{\mu}\right)^2 \log \frac{F(x_0) - F^*}{\varepsilon}\right)$ for **Damped Newton method**,

$O\left(\frac{L - \mu}{\mu} \log \frac{F(x_0) - F^*}{\varepsilon}\right)$ for **Adaptive Cubic Newton**.

1. Introduction: first-order nondegeneracy
2. Regularization of Newton method
3. Universal Cubic Newton
4. Numerical examples

$$\min_{x \in \mathbb{R}^n} f(x) = \log \left(\sum_{i=1}^m e^{\langle a_i, x \rangle} \right) + \frac{\mu}{2} \langle x, x \rangle.$$

- ▶ $a_1, \dots, a_m \in \mathbb{R}^n$ — given data.
- ▶ Denote $B \equiv \sum_{i=1}^m a_i a_i^T \succeq 0$, and use $\|x\| \equiv \langle Bx, x \rangle^{1/2}$.

- ▶ We have

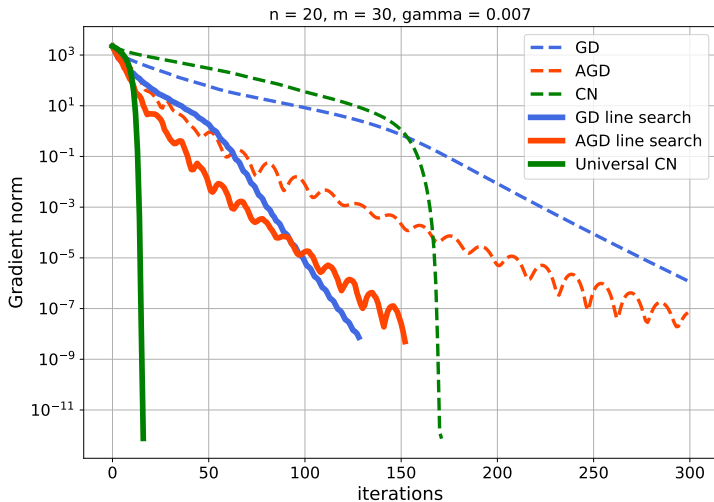
$$\mathcal{H}_f(0) \leq 1, \quad \mathcal{H}_f(1) \leq 2.$$

- ▶ For any $\nu \in [0, 1]$:

$$\mathcal{H}_f(\nu) \leq (\mathcal{H}_f(0))^{1-\nu} \cdot (\mathcal{H}_f(1))^\nu \leq 2^\nu.$$

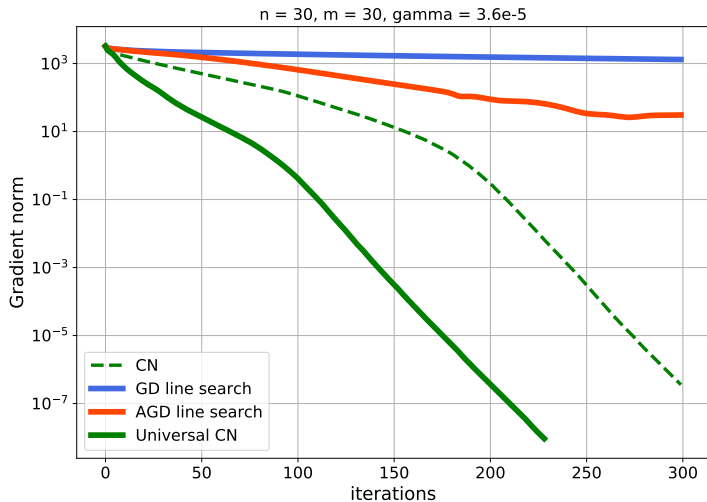
Log-sum-exp convergence

$$\min_{x \in \mathbb{R}^n} f(x) = \log \left(\sum_{i=1}^m e^{\langle a_i, x \rangle} \right) + \frac{\mu}{2} \langle x, x \rangle.$$



Log-sum-exp, ill-conditioned problem

$$\min_{x \in \mathbb{R}^n} f(x) = \log \left(\sum_{i=1}^m e^{\langle a_i, x \rangle} \right) + \frac{\mu}{2} \langle x, x \rangle.$$



Conclusion

- ▶ Second-order condition number $\gamma_f(\nu)$ of degree $\nu \in [0, 1]$.
- ▶ Adaptive cubically regularized Newton achieves linear rate of convergence when $\gamma_f(\nu) > 0$ for some $\nu \in [0, 1]$. The main complexity factor is $(\gamma_f(\nu))^{\frac{-1}{1+\nu}}$.
- ▶ The algorithm is **universal**: it does not need to know right $\nu \in [0, 1]$ and any other parameters of the class.
- ▶ For the class $f \in \mathcal{S}_{\mu, L}^{2,1}$: $\mu l \preceq \nabla^2 f(x) \preceq Ll$ provided complexity bounds for Cubic Newton are always better than for Gradient Method.

Thank you for your attention!