

Newtonian Methods in Nonsmooth Optimization via the Lens of Variational Analysis

VO THANH PHAT

Assistant Professor - University of North Dakota, USA

(Based on the joint work with

Pham Duy Khanh, HCMC University of Education, Vietnam)

Boris Mordukhovich, Wayne State University, USA

**26th Midwest Optimization Meeting
Workshop on Large Scale Optimization and Applications
Hosted by the University of Waterloo
November 8-9, 2024**

- Structured Nonconvex Optimization Problems
- Classical Newton's Method and Tools of Variational Analysis
- Coderivative-Based Newton Method for Structured Nonconvex Optimization Problems.
- Applications
- Future Investigation

Structured Nonconvex Optimization Problems

Structured Nonconvex Optimization Problems

Our target optimization problem is

$$\text{minimize } \varphi(x) := f(x) + g(x), \quad (1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is \mathcal{C}^2 -smooth, and $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}} := (-\infty, \infty]$ is a lower semicontinuous, and prox-bounded¹

- Note that both f and g are generally nonconvex and nonsmooth. This makes (1) appropriate for applications to signal and image processing, machine learning, statistics, control, system identification, etc.
- If, in particular, g is the indicator function of a closed set, then (1) becomes a constrained optimization problems with numerous applications.

¹A function is prox-bounded if it is bounded from below by some quadratic function.

Classical Newton's Method and Tools of Variational Analysis

Classical Newton's Method

Consider the gradient system

$$\nabla\varphi(x) = 0, \quad (2)$$

where $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$ be a \mathcal{C}^2 -smooth function.

The classical version of **Newton's method**: $x^0 \in \mathbb{R}^n$ is given, then

$$x^{k+1} = x^k + d^k \quad \text{with} \quad -\nabla\varphi(x^k) = \nabla^2\varphi(x^k)d^k, \quad k = 0, 1, \dots \quad (3)$$

This algorithm is **locally well-defined** and **superlinearly converges** to a solution \bar{x} of (2) if $\nabla^2\varphi(\bar{x})$ is nonsingular.²

²Izmailov, A. F., & Solodov, M. V. (2014). *Newton-type methods for optimization and variational problems (Vol. 1)*. New York: Springer.

Generalized Differentiation

Requirement of generalized derivatives for generalized Newton method:

- Comprehensive calculus rules: sum rules, chain rules, product rules, etc.
- Explicit calculations in a number of settings important for applications.
- Can characterize the convexity, generalized convexity, local optimality, etc.

To develop Newton's method, we use the generalized derivatives including limiting first- and second-order subdifferentials introduced by Mordukhovich, which are presented in the books ³⁴.

³Mordukhovich, B. S. (2018). *Variational analysis and applications* (Vol. 30). Cham: Springer.

⁴Rockafellar, R. T., & Wets, R. J. B. (2009). *Variational analysis* (Vol. 317). Springer Science & Business Media.

Generalized Differentiation

The **normal cone** to $\Omega \subset \mathbb{R}^n$ at $\bar{x} \in \Omega$ from

$$N_{\Omega}(\bar{x}) := \{v \mid \exists x_k \rightarrow \bar{x}, \alpha_k \geq 0, w_k \in \Pi_{\Omega}(x_k), \alpha_k(x_k - w_k) \rightarrow v\}$$

where Π_{Ω} stands for the Euclidean projection. The **coderivative** of $F: \mathbb{R}^n \rightrightarrows \mathbb{R}^m$ at $(\bar{x}, \bar{y}) \in \text{gph } F$

$$D^*F(\bar{x}, \bar{y})(v) := \{u \in \mathbb{R}^n \mid (u, -v) \in N_{\text{gph } F}(\bar{x}, \bar{y})\}, \quad v \in \mathbb{R}^m.$$

When $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is \mathcal{C}^1 -smooth, then

$$D^*F(\bar{x})(v) = \{\nabla F(\bar{x})^* v\}, \quad v \in \mathbb{R}^m,$$

via the adjoint/transpose Jacobian matrix. The (first-order) **subdifferential** of $\varphi: \mathbb{R}^n \rightarrow (-\infty, \infty]$ at $\bar{x} \in \text{dom } \varphi$

$$\partial\varphi(\bar{x}) := \{v \in \mathbb{R}^n \mid (v, -1) \in N_{\text{epi } \varphi}(\bar{x}, \varphi(\bar{x}))\}.$$

Generalized Differentiation - Second-order Subdifferential


Second-order subdifferential/generalized Hessian⁵ of φ at \bar{x} relative to $\bar{v} \in \partial\varphi(\bar{x})$ is

$$\partial^2\varphi(\bar{x}, \bar{v})(u) := (D^*\partial\varphi)(\bar{x}, \bar{v})(u), \quad u \in \mathbb{R}^n$$

If $\varphi \in \mathcal{C}^2$ -smooth around \bar{x} , then

$$\partial^2\varphi(\bar{x}, \bar{v})(u) = \{\nabla^2\varphi(\bar{x})u\}, \quad u \in \mathbb{R}^n$$

It is realized that the generalized Hessian $\partial^2\varphi$ enjoys well-developed **second-order calculus** and can be viewed as an appropriate replacement of the Hessian $\nabla^2\varphi$ for nonsmooth problems. $\partial^2\varphi$ is **fully computed** in terms of the given data for broad classes of problems in optimization and variational analysis.

⁵Mordukhovich, B.S.: *Sensitivity analysis in nonsmooth optimization*. In: Field, D.A., Komkov, V.(eds) *Theoretical Aspects of Industrial Design*, 32–46. *SIAM Proc. Appl. Math.* 58. Philadelphia, PA (1992) 

Goal: Approximate an optimal solution to the following optimization problem

$$\text{minimize } \varphi(x) \quad \text{subject to } x \in \mathbb{R}^n \quad (4)$$

where $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$ is **not necessarily smooth**.

Ideas:

- $\nabla\varphi(x^k) \rightarrow [v^k \in \partial\varphi(x^k)]$.
- $\nabla^2\varphi(x^k) \rightarrow \partial^2\varphi(x^k, v^k)$, for some $v^k \in \partial\varphi(x^k)$.

However, $\partial\varphi(x^k)$ may be empty. So, we will choose the pair $(\hat{x}^k, \hat{v}^k) \in \text{gph } \partial\varphi$ such that $\partial\varphi(\hat{x}^k)$ is **nonempty** in which \hat{x}^k is **not far** from x^k in some senses.

General Framework of Coderivative-Based Newton Method

Suppose that \bar{x} is a stationary point, i.e., $0 \in \partial\varphi(\bar{x})$. Our **generalized coderivative-based Newton method** can be formulated as: $x^0 \in \mathbb{R}^n$ is given, then

$$x^{k+1} = \hat{x}^k + d^k \quad \text{with } -\hat{v}^k \in \partial^2\varphi(\hat{x}^k, \hat{v}^k)(d^k), \hat{v}^k \in \partial\varphi(\hat{x}^k),$$

where (\hat{x}^k, \hat{v}^k) satisfying the following inequality

$$\|(\hat{x}^k, \hat{v}^k) - (x^k, 0)\| \leq \eta \|x^k - \bar{x}\|. \quad (5)$$

The step in (5) is called the **approximate step** in our algorithm.

In fact, we can choose (\hat{x}^k, \hat{v}^k) as an “**approximate projection**” of $(x^k, 0)$ on $\text{gph}\partial\varphi$ in the sense that

$$\|(\hat{x}^k, \hat{v}^k) - (x^k, 0)\| \leq \eta \text{dist}((x^k, 0), \text{gph}\partial\varphi).$$

Convergence Analysis of Coderivative-Based Newtonian Methods

Question: Do we guarantee the **convergence** of the iterative sequence generated by our aforementioned method⁶?

⁶ $x^0 \in \mathbb{R}^n$ is given, then

$$x^{k+1} = \hat{x}^k + d^k \quad \text{with} \quad -\hat{v}^k \in \partial^2 \varphi(\hat{x}^k, \hat{v}^k)(d^k), \quad \hat{v}^k \in \partial \varphi(\hat{x}^k), \quad (6)$$

where $k \in \mathbb{N}$.

Definition

$\varphi: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is **prox-regular**^a at $\bar{x} \in \text{dom } \varphi$ for $\bar{v} \in \partial\varphi(\bar{x})$ if φ is lower semicontinuous around \bar{x} and there are $\varepsilon > 0$ and $\rho \geq 0$ such that for all $x \in \mathbb{B}_\varepsilon(\bar{x})$ with $\varphi(x) \leq \varphi(\bar{x}) + \varepsilon$ we have

$$\varphi(x) \geq \varphi(u) + \langle v, x - u \rangle - \frac{\rho}{2} \|x - u\|^2$$

for all $(u, v) \in (\text{gph } \partial\varphi) \cap \mathbb{B}_\varepsilon(\bar{x}, \bar{v})$.

^aPoliquin, R., & Rockafellar, R. (1996). Prox-regular functions in variational analysis. *Transactions of the American Mathematical Society*, 348(5), 1805-1838.

φ is **subdifferentially continuous** at \bar{x} for \bar{v} if the convergence $(x_k, v_k) \rightarrow (\bar{x}, \bar{v})$ with $v_k \in \partial\varphi(x_k)$ yields $\varphi(x_k) \rightarrow \varphi(\bar{x})$. If both properties hold, φ is **continuously prox-regular**. This is the **major class** in second-order variational analysis.

Definition

^a A mapping $F: \mathbb{R}^n \rightrightarrows \mathbb{R}^m$ is **semismooth*** at $(\bar{x}, \bar{y}) \in \text{gph } F$ if whenever $(u, v) \in \mathbb{R}^n \times \mathbb{R}^m$ we have the equality

$$\langle u^*, u \rangle = \langle v^*, v \rangle \text{ for all } (v^*, u^*) \in \text{gph } D^*F((\bar{x}, \bar{y}); (u, v)).$$

^aGfrerer, H., & Outrata, J. V. (2021). On a semismooth* Newton method for solving generalized equations. *SIAM Journal on Optimization*, 31(1), 489-517.

Example:

- (i) A **continuously differentiable mapping** is semismooth*.
- (ii) A set-valued mapping with the graph represented as a **union of finitely many closed and convex sets** is semismooth*.

Convergence Analysis

- **Local Convergence:** Guarantee the **local convergence** to a **local optimal solution** \bar{x} if φ is **continuously prox-regular** at \bar{x} for 0, $\partial\varphi$ is **semismooth*** at \bar{x} , and \bar{x} satisfies the **second-order sufficient optimality condition** in the sense that

$$0 \in \partial\varphi(\bar{x}) \quad \text{and} \quad \partial^2\varphi(\bar{x}, 0) \text{ is positive definite.}$$

- **Convergence Rate:** **superlinear** in the sense that

$$\lim_{k \rightarrow \infty} \|x^{k+1} - \bar{x}\| / \|x^k - \bar{x}\| = 0.$$

More detail in our work⁷.

⁷Khanh, P.D., Mordukhovich, B.S., Phat, V.T.: *Coderivative-based Newton methods in structured nonconvex and nonsmooth optimization*, arXiv:2403.04262.

Implementation of the Approximate Step

- When φ is a $\mathcal{C}^{1,1}$ -smooth function, we can choose

$$\hat{x}^k := x^k \text{ and } \hat{v}^k := \nabla\varphi(x^k).$$

- When $\text{Prox}_{\lambda\varphi}$ ⁸ can be computed explicitly, we can choose

$$\hat{x}^k := \text{Prox}_{\lambda\varphi}(x^k) \text{ and } \hat{v}^k := \frac{1}{\lambda} \left(x^k - \text{Prox}_{\lambda\varphi}(x^k) \right).$$

The natural question is how to implement the approximate step when $\varphi = f + g$, where f is \mathcal{C}^2 -smooth, g is prox-bounded function? We will discuss in more detail in the next section.

⁸ $\text{Prox}_{\lambda\varphi}(x) = \operatorname{argmin}_{y \in \mathbb{R}^n} \{ \varphi(y) + 1/(2\lambda) \|y - x\|^2 \}$

Coderivative-Based Newton Method for Structured Nonconvex Optimization Problems

Coderivative-Based Newton Method for Structured Nonconvex Optimization Problems

Consider the problem

$$\text{minimize } \varphi(x) := f(x) + g(x), \quad (7)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is \mathcal{C}^2 -smooth, and $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}} := (-\infty, \infty]$ is a lower semicontinuous, and prox-bounded function.

To use our coderivative-based Newton method to find a stationary point \bar{x} to (7), i.e., $0 \in \partial\varphi(\bar{x})$, we need to clarify two following questions:

- How do we implement the **approximate step**?
- How do we guarantee the **global convergence** of the iterative sequence?

Approximate Step in Structured Nonconvex Optimization Problems

When $\varphi = f + g$, where f is \mathcal{C}^2 -smooth, g is prox-bounded, we can choose

$$\hat{x}^k \in \text{Prox}_{\lambda g}(x^k - \lambda \nabla f(x^k))$$

and

$$\hat{v}^k := \nabla f(\hat{x}^k) - \nabla f(x^k) + \frac{1}{\lambda}(x^k - \hat{x}^k).$$

In this case, we have

- $\hat{v}^k \in \partial\varphi(\hat{x}^k)$.
- There is $\eta > 0$ such that

$$\|(\hat{x}^k, \hat{v}^k) - (x^k, 0)\| \leq \eta \|x^k - \bar{x}\|.$$

\implies We can apply our method to guarantee the **locally superlinear convergence** of $\{x^k\}$.

To obtain the **global convergence** of our method for the nonconvex optimization problem, a natural approach is to consider the sequence $\{x^k\}$ as follows:

$$x^{k+1} := \hat{x}^k + \tau_k d^k \quad (8)$$

with an appropriate stepsize selection $\tau_k \in (0, 1]$, with an expectation that the **descent property** holds

$$\varphi(x^{k+1}) = \varphi(\hat{x}^k + \tau_k d^k) < \varphi(x^k), \quad k = 0, 1, \dots$$

However this is impossible to guarantee due to the **discontinuity** of the cost function φ .

Globalization

Fortunately, we can approximate the cost function φ by a differentiable function called **forward-backward envelope**⁹¹⁰¹¹, and we can guarantee the **descent property** of this function, i.e.,

$$\varphi_\lambda(x^{k+1}) = \varphi_\lambda(\hat{x}^k + \tau_k d^k) < \varphi_\lambda(x^k), \quad k = 0, 1, \dots,$$

where φ_λ is defined by

$$\varphi_\lambda(x) := \inf_{y \in \mathbb{R}^n} \left\{ f(x) + \langle \nabla f(x), y - x \rangle + g(y) + \frac{1}{2\lambda} \|y - x\|^2 \right\}, \quad (9)$$

⁹Patrinos, P., & Bemporad, A. (2013, December). Proximal Newton methods for convex composite optimization. In *52nd IEEE Conference on Decision and Control* (pp. 2358-2363). IEEE.

¹⁰Stella, L., Themelis, A., & Patrinos, P. (2017). Forward-backward quasi-Newton methods for nonsmooth optimization problems. *Computational Optimization and Applications*, 67(3), 443-487.

¹¹Themelis, A., Stella, L., & Patrinos, P. (2018). Forward-backward envelope for the sum of two nonconvex functions: Further properties and nonmonotone linesearch algorithms. *SIAM Journal on Optimization*, 28(3), 2274-2303.

Globalized Coderivative-Based Newton Method

$x^0 \in \mathbb{R}^n$ is given, then

$$x^{k+1} = \hat{x}^k + \tau_k d^k \quad \text{with } -\hat{v}^k \in \partial^2 \varphi(\hat{x}^k, \hat{v}^k)(d^k),$$

where

$$\hat{x}^k \in \text{Prox}_{\lambda g}(x^k - \lambda \nabla f(x^k)) \quad \text{and} \quad \hat{v}^k := \nabla f(\hat{x}^k) - \nabla f(x^k) + \frac{1}{\lambda}(x^k - \hat{x}^k)$$

and $\tau_k \in (0, 1]$ satisfying

$$\varphi_\lambda(\hat{x}^k + \tau_k d^k) \leq \varphi_\lambda(x^k) - \sigma \|\hat{v}^k\|^2.$$

Convergence analysis of our method

- **Well-Posedness:** The sequence $\{x^k\}$ is **well-defined**. Both sequences $\{\widehat{v}^k\}$ and $\{\widehat{x}^k - x^k\}$ converge to 0 as $k \rightarrow \infty$. Finally, any **accumulation point** of $\{x^k\}$ is a **stationary point**.
- **Global Convergence:** Guarantee the global convergence of $\{x^k\}$ if g is **continuously prox-regular** at \bar{x} for $-\nabla f(\bar{x})$, $\nabla^2 f$ is **strictly differentiable** at \bar{x} , ∂g is **semismooth*** at \bar{x} , and \bar{x} satisfies the **second-order sufficient optimality condition** in the sense that

$$0 \in \partial\varphi(\bar{x}) \quad \text{and} \quad \partial^2\varphi(\bar{x}, 0) \text{ is positive definite.}$$

- **Convergence Rate:** **superlinear** in the sense that

$$\lim_{k \rightarrow \infty} \|x^{k+1} - \bar{x}\| / \|x^k - \bar{x}\| = 0.$$

Applications

Applications

Given an $m \times n$ matrix A and a vector $b \in \mathbb{R}^m$, the ℓ_0 - ℓ_2 regularized least square regression problem whose importance has been well recognized in applications to practical models of statistics, machine learning, etc¹². This problem is formulated as:

$$\min \varphi(x) := \frac{1}{2} \|Ax - b\|_2^2 + \mu_0 \|x\|_0 + \mu_2 \|x\|_2^2 \text{ subject to } x \in \mathbb{R}^n$$

where μ_0 and μ_2 are positive parameters, and where $\|x\|_0$ is the ℓ_0 norm of x counting the number of nonzero elements of x .

Our numerical experiment in ¹³ shows that our method behaves better than proximal gradient method for solving the above problem.

¹²Hazimeh, H., & Mazumder, R. (2020). Fast best subset selection: Coordinate descent and local combinatorial optimization algorithms. *Operations Research*, 68(5), 1517-1537.

¹³Khanh, P.D., Mordukhovich, B.S., Phat, V.T.: Coderivative-based Newton methods in structured nonconvex and nonsmooth optimization, *arXiv:2403.04262*.

Future Investigation

Future Investigation

- Establish the generalized Newton method for solving difference programming.
- Establish the generalized Newton method for solving multiobjective optimization problems.
- Establish the generalized Newton method for solving bilevel optimization problems.
- Establish the coderivative-based stochastic Newton method for solving nonsmooth and nonconvex optimization problems with high dimension.
- Applications to other important classes of models in data science, machine learning, statistic, and related disciplines.

- Khanh, P.D., Mordukhovich, B.S., Phat, V.T.: *A generalized Newton method for subgradient systems*. **Math. Oper. Res.** 48, 1811–1845 (2022).
- Khanh, P.D., Mordukhovich, B.S., Phat, V.T., Tran, D. B.: *Generalized damped Newton algorithms in nonsmooth optimization via second-order subdifferentials*. **J. Global Optim.** 86, 93–122 (2023).
- Khanh, P.D., Mordukhovich, B.S., Phat, V.T., Tran, D. B.: *Globally convergent coderivative-based generalized Newton methods in nonsmooth optimization*. **Math. Program.** 205, 373–429 (2024).
- Khanh, P.D., Mordukhovich, B.S., Phat, V.T.: *Coderivative-based Newton methods in structured nonconvex and nonsmooth optimization*, arXiv:2403.04262.

THANK YOU FOR YOUR ATTENTION