

Affine-invariant contracting-point methods for Convex Optimization

Nikita Doikov

Joint work with Yurii Nesterov

UCLouvain, Belgium

Symposium on Numerical Analysis and Optimization, UFPR
March 4, 2021

1. Introduction
2. Contracting-Point Methods
3. Inexact and Stochastic Algorithms
4. Conclusions

1. Introduction
2. Contracting-Point Methods
3. Inexact and Stochastic Algorithms
4. Conclusions

Composite Optimization Problem

$$\min_x \left\{ F(x) \stackrel{\text{def}}{=} f(x) + \psi(x) \right\}$$

- ▶ f is convex and several times differentiable (the *difficult part*).
- ▶ $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a *simple* convex function.
- ▶ We assume that the domain of ψ ,

$$\text{dom } \psi \stackrel{\text{def}}{=} \left\{ x \in \mathbb{R}^n : \psi(x) < +\infty \right\},$$

is **bounded**.

Example: Indicator of a Set

1. Let $Q \subset \mathbb{R}^n$ be a simple bounded convex set. We can use

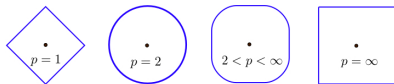
$$\psi(x) = \text{Ind}_Q(x) := \begin{cases} 0, & x \in Q \\ +\infty, & \text{otherwise.} \end{cases}$$

\Rightarrow Then our problem is

$$\min_{x \in Q} f(x)$$

Examples of sets:

► The Ball: $Q = \left\{ x \in \mathbb{R}^n : \|x\|_p := \left(\sum_{i=1}^n |x^{(i)}|^p \right)^{1/p} \leq \frac{D}{2} \right\}$.



► The standard Simplex: $Q = \mathbb{S}_n := \left\{ x \in \mathbb{R}_+^n : \sum_{i=1}^n x^{(i)} = 1 \right\}$.

0-simplex

1-simplex

2-simplex

3-simplex

Example: ℓ_1 -Regularization

2. Let

$$\psi(x) = \begin{cases} \lambda \|x\|_1, & x \in Q \\ +\infty, & \text{otherwise.} \end{cases}$$

\Rightarrow Adding ℓ_1 -Regularizer to the problem:

$$\min_{x \in Q} f(x) + \lambda \|x\|_1.$$

Enforce solutions to be **sparse**.

- ▶ Machine Learning, Statistics. Empirical Risk Minimization:

$$f(x) = \frac{1}{M} \sum_{i=1}^M f_i(x).$$

Generalized Linear Models: $f_i(x) = \phi(\langle a_i, x \rangle - b_i)$, where $\{a_i, b_i\}$ is given data, ϕ is a convex loss.

Logistic Regression: $\phi(x) = \log(1 + e^x)$.

- ▶ Smooth Approximation of a Nondifferentiable Function.

$$f(x) \approx \bar{f}(x) = \max_{i=1}^m \{\langle a_i, x \rangle - b_i\} \quad (\text{pointwise maximum}).$$

Log-sum-exp (SoftMax): $f(x) = \mu \log\left(\sum_{i=1}^m e^{(\langle a_i, x \rangle - b_i)/\mu}\right)$.

The Gradient Method:

$$x_{k+1} = \operatorname{argmin}_y \left\{ \langle \nabla f(x_k), y - x_k \rangle + \frac{L}{2} \|y - x_k\|^2 + \psi(y) \right\}.$$

Note: when $\psi(x) = \operatorname{Ind}_Q(x)$ and the norm is Euclidean, the iterations can be rewritten in a canonical form:

$$x_{k+1} = \operatorname{proj}_Q(x_k - \frac{1}{L} \nabla f(x_k)).$$

- ▶ Global convergence: $F(x_K) - F^* \leq \varepsilon$ for $K = \mathcal{O}(\frac{1}{\varepsilon})$, if the gradient is Lipschitz continuous.
- ▶ The method depends on the **norm** $\|\cdot\|$.
- ▶ Cheap iterations, but the rate of convergence is **slow**.

The Newton's Method:

$$x_{k+1} = \operatorname{argmin}_y \left\{ \langle \nabla f(x_k), y - x_k \rangle + \frac{1}{2} \langle \nabla^2 f(x_k)(y - x_k), y - x_k \rangle + \psi(y) \right\}.$$

If $\psi(x) \equiv 0$, then

$$x_{k+1} = x_k - (\nabla^2 f(x_k))^{-1} \nabla f(x_k).$$

- ▶ Quadratic convergence $\mathcal{O}(\log \log \frac{1}{\varepsilon})$, if $\nabla^2 f(x^*) \succ 0$ and x_0 close to x^* .
- ▶ No global convergence. A heuristic: use line-search in practice.
- ▶ The method is affine-invariant (it does not use any norms).

The goal: to develop second- and high-order algorithms with global convergence guarantees.

- ▶ The rate of second-order methods should be **better** than that of first-order methods.

We propose a general framework of **Contracting-Point Methods**.

- ▶ New **affine-invariant** algorithms of different order $p \geq 1$.
- ▶ We prove: $F(x_k) - F^* \leq \mathcal{O}(1/k^p)$.

1. Introduction
2. Contracting-Point Methods
3. Inexact and Stochastic Algorithms
4. Conclusions

Conceptual Contracting-Point Scheme

$g_k(x) = f(\gamma_k x + (1 - \gamma_k)x_k)$, $\gamma_k \in (0, 1]$. Then

$$Dg_k(x) = \gamma_k Df(\gamma_k x + (1 - \gamma_k)x_k),$$

$$D^2 g_k(x) = \gamma_k^2 D^2 f(\gamma_k x + (1 - \gamma_k)x_k), \dots$$

- ▶ Smoothness properties of $g_k(\cdot)$ can be better than $f(\cdot)$.

Conceptual Contracting-Point Method:

$$\begin{aligned} v_{k+1} &\in \underset{x}{\operatorname{Argmin}} \left\{ f(\gamma_k x + (1 - \gamma_k)x_k) + \gamma_k \psi(x) \right\}, \\ x_{k+1} &= \gamma_k v_{k+1} + (1 - \gamma_k)x_k. \end{aligned}$$

Lemma: Inexact Contracting Step

The subproblem: $\min_x \left\{ F_k(x) \stackrel{\text{def}}{=} f(\gamma_k x + (1 - \gamma_k)x_k) + \gamma_k \psi(x) \right\}$.

Lemma.

Let v_{k+1} be an **approximate** solution: $F_k(v_{k+1}) - F_k^* \leq \delta_{k+1}$. Then

$$F(x_{k+1}) \leq (1 - \gamma_k)F(x_k) + \gamma_k F^* + \delta_{k+1}.$$

Proof: For any $v \in \text{dom } \psi$, it holds

$$\begin{aligned} & (1 - \gamma_k)F(x_k) + \gamma_k F(v) \\ &= (1 - \gamma_k)f(x_k) + \gamma_k f(v) + (1 - \gamma_k)\psi(x_k) + \gamma_k \psi(v) \end{aligned}$$

$$\stackrel{\text{convexity}}{\geq} F_k(v) + (1 - \gamma_k)\psi(x_k)$$

$$\stackrel{\text{def. of } v_{k+1}}{\geq} F_k(v_{k+1}) + (1 - \gamma_k)\psi(x_k) - \delta_{k+1}$$

$$\stackrel{\text{convexity}}{\geq} F(\gamma_k v_{k+1} + (1 - \gamma_k)x_k) - \delta_{k+1} = F(x_{k+1}) - \delta_{k+1}. \quad \square$$

Corollary: Global Convergence

The Lemma provides us with the following guarantee of one inexact Contracting-Point step:

$$F(x_{k+1}) \leq (1 - \gamma_k)F(x_k) + \gamma_k F^* + \delta_{k+1} \quad (*)$$

Let us take an arbitrary increasing sequence $\{A_k\}_{k \geq 1}$, $A_0 := 0$, and denote:

$$a_{k+1} := A_{k+1} - A_k > 0, \quad \gamma_k := \frac{a_{k+1}}{A_{k+1}}.$$

Substituting this into $(*)$, we obtain

$$A_{k+1}F(x_{k+1}) \leq A_k F(x_k) + a_{k+1}F^* + A_{k+1}\delta_{k+1}.$$

Corollary: $F(x_k) - F^* \leq \frac{1}{A_k} \sum_{i=1}^k A_i \delta_i.$

- ▶ If A_k grow \nearrow sufficiently fast, and δ_i are small, we have global convergence.

Affine-Invariant Smoothness Condition

Fix $p \geq 1$. For a bounded convex set Q , denote

$$\mathcal{V}_Q^{(p+1)}(f) \stackrel{\text{def}}{=} \sup_{x, y, v \in Q} |D^{p+1}f(y)[v - x]^{p+1}|.$$

Note: for a fixed norm, we have $\mathcal{V}_Q^{(p+1)}(f) \leq L_p(\text{diam } Q)^{p+1}$, where L_p is the Lipschitz constant for p th derivative.

It holds, $\forall x, x_k \in Q$ and $\forall \gamma_k \in (0, 1]$:

$$\begin{aligned} & \left| f(\gamma_k x + (1 - \gamma_k)x_k) - f(x_k) - \sum_{i=1}^p \frac{\gamma_k^i}{i!} D^i f(x_k)[x - x_k]^i \right| \\ & \leq \frac{\gamma_k^{p+1}}{(p+1)!} \mathcal{V}_Q^{(p+1)}(f). \quad (\text{Taylor's Theorem}). \end{aligned}$$

Contracting-Point Tensor Method:

$$\begin{aligned} v_{k+1} &\in \underset{x}{\operatorname{Argmin}} \left\{ \sum_{i=1}^p \frac{\gamma_k^i}{i!} D^i f(x_k) [x - x_k]^i + \gamma_k \psi(x) \right\}, \\ x_{k+1} &= \gamma_k v_{k+1} + (1 - \gamma_k) x_k. \end{aligned}$$

Since $\operatorname{dom} \psi$ is bounded, the subproblem is well-defined.

- ▶ $p = 1$: The Conditional Gradient Method [Frank-Wolfe, 1956].
- ▶ $p = 2$: Contracting Newton. (new)

Global Convergence Rate

Theorem. Set $\gamma_k := \frac{\rho+1}{k+\rho+1}$. Then $F(x_k) - F^* \leq \mathcal{O}\left(\frac{\mathcal{V}_{\text{dom } \psi}^{(\rho+1)}(f)}{k^\rho}\right)$.

Proof: Let us consider one Contracting-Point Tensor step. We have, for any $x \in \text{dom } \psi$

$$F_k(x) \equiv f(\gamma_k x + (1 - \gamma_k)x_k) + \gamma_k \psi(x)$$

$$\stackrel{\text{Taylor's}}{\geq} f(x_k) + \sum_{i=1}^{\rho} \frac{\gamma_k^i}{i!} D^i f(x_k)[x - x_k]^i + \gamma_k \psi(x) - \frac{\gamma_k^{\rho+1} \mathcal{V}_{\text{dom } \psi}^{(\rho+1)}(f)}{(\rho+1)!}$$

$$\stackrel{\text{Step}}{\geq} f(x_k) + \sum_{i=1}^{\rho} \frac{\gamma_k^i}{i!} D^i f(x_k)[v_{k+1} - x_k]^i + \gamma_k \psi(v_{k+1}) - \frac{\gamma_k^{\rho+1} \mathcal{V}_{\text{dom } \psi}^{(\rho+1)}(f)}{(\rho+1)!}$$

$$\stackrel{\text{Taylor's}}{\geq} f(\gamma_k v_{k+1} + (1 - \gamma_k)x_k) + \gamma_k \psi(v_{k+1}) - \frac{2\gamma_k^{\rho+1} \mathcal{V}_{\text{dom } \psi}^{(\rho+1)}(f)}{(\rho+1)!}$$

$$\equiv F_k(v_{k+1}) - \frac{2\gamma_k^{\rho+1} \mathcal{V}_{\text{dom } \psi}^{(\rho+1)}(f)}{(\rho+1)!}.$$

Global Convergence Rate – Continue with the Proof

Hence, one step of the Tensor method approximates a step of the conceptual Contracting-Point scheme:

$$F_k(v_{k+1}) - F_k^* \leq \delta_{k+1} := \frac{2\gamma_k^{p+1} \mathcal{V}_{\text{dom } \psi}^{(p+1)}(f)}{(p+1)!}.$$

The Corollary implies the global convergence:

$$F(x_k) - F^* \leq \frac{1}{A_k} \sum_{i=1}^k A_i \delta_i = \frac{2\mathcal{V}_{\text{dom } \psi}^{(p+1)}(f)}{(p+1)!} \cdot \frac{1}{A_k} \sum_{i=1}^k \frac{a_i^{p+1}}{A_i^p}.$$

Note: $A_k \approx k^{p+1}$, then $a_k = A_k - A_{k-1} \approx k^p$, $\frac{a_i^{p+1}}{A_i^p} \approx \text{const}$, and

$$\frac{1}{A_k} \sum_{i=1}^k \frac{a_i^{p+1}}{A_i^p} \approx \frac{1}{k^p}.$$

An appropriate choice

$$A_k := k \cdot (k+1) \cdot \dots \cdot (k+p) \Rightarrow \boxed{\gamma_k = \frac{a_{k+1}}{A_{k+1}} = \frac{p+1}{k+p+1}}$$



1. Introduction
2. Contracting-Point Methods
3. Inexact and Stochastic Algorithms
4. Conclusions

$p = 1$: Conditional Gradient Method:

$$\begin{aligned} v_{k+1} &\in \underset{x}{\operatorname{Argmin}} \left\{ \langle \nabla f(x_k), x - x_k \rangle + \psi(x) \right\}, \\ x_{k+1} &= \gamma_k v_{k+1} + (1 - \gamma_k) x_k. \end{aligned}$$

► $F(x_k) - F^* \leq \mathcal{O}(1/k)$. The same rate as of the GM.

$p = 2$: Contracting Newton Method:

$$\begin{aligned} v_{k+1} &\in \underset{x}{\operatorname{Argmin}} \left\{ \langle \nabla f(x_k), x - x_k \rangle + \frac{\gamma_k}{2} \langle \nabla^2 f(x_k)(x - x_k), x - x_k \rangle \right. \\ &\quad \left. + \psi(x) \right\}, \\ x_{k+1} &= \gamma_k v_{k+1} + (1 - \gamma_k) x_k. \end{aligned}$$

► $F(x_k) - F^* \leq \mathcal{O}(1/k^2)$.

Contracting Newton Method (reformulation):

$$x_{k+1} = \operatorname{argmin}_y \left\{ \langle \nabla f(x_k), y - x_k \rangle + \frac{1}{2} \langle \nabla^2 f(x_k)(y - x_k), y - x_k \rangle + \gamma_k \psi(x_k + \frac{1}{\gamma_k}(y - x_k)) \right\}.$$

Regularization of quadratic model by the asymmetric **trust region**.

» How to solve the subproblem?

If $\psi(x) = \operatorname{Ind}_Q(x)$, where $Q = \{x \in \mathbb{R}^n : \|x\| \leq \frac{D}{2}\}$ is the ball, we can use techniques developed for Trust-Region methods.

See books: [Conn-Gould-Toint, 2000], [Nocedal-Wright, 2006].

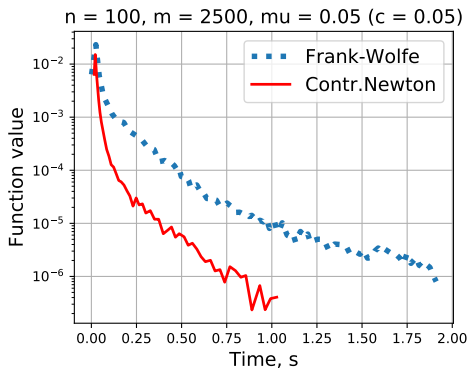
- ▶ At each step we can solve the subproblem inexactly by the first-order Conditional Gradient Method.
- ▶ We have the full control over the required accuracy.

Theorem. To reach $F(x_K) - F^* \leq \varepsilon$ we need $K = \mathcal{O}\left(\frac{1}{\sqrt{\varepsilon}}\right)$ outer iterations (oracle calls for f).

The total number of linear minimization oracle calls for ψ is $\mathcal{O}\left(\frac{1}{\varepsilon}\right)$.

Experiment: Log-sum-exp over the Simplex

$$\min_{x \in \mathbb{S}_n} f(x), \quad \mathbb{S}_n \stackrel{\text{def}}{=} \left\{ x \in \mathbb{R}_+^n : \sum_{i=1}^n x^{(i)} = 1 \right\}.$$



► two times faster.

Finite-sum minimization: $f(x) = \frac{1}{M} \sum_{i=1}^M f_i(x)$.

- ▶ M can be very big in modern applications (several millions).
- ▶ Machine Learning: M is the size of the dataset.

It is expensive to compute the full gradient and Hessian:

$$\nabla f(x) = \frac{1}{M} \sum_{i=1}^M \nabla f_i(x), \quad \nabla^2 f(x) = \frac{1}{M} \sum_{i=1}^M \nabla^2 f_i(x).$$

Random estimators:

$$\begin{aligned} \nabla f(x_k) &\approx g_k := \frac{1}{m_k^g} \sum_{i \in S_k^g} \nabla f_i(x_k), \\ \nabla^2 f(x_k) &\approx H_k := \frac{1}{m_k^H} \sum_{i \in S_k^H} \nabla^2 f_i(x_k). \end{aligned}$$

$S_k^g, S_k^H \subseteq \{1, \dots, M\}$ are random subsets (sampled uniformly) for a fixed **batchsize** $m_k^g = |S_k^g|$, and $m_k^H = |S_k^H|$.

Stochastic Contracting Newton

$$\begin{aligned}g_k &:= \frac{1}{m_k^g} \sum_{i \in S_k^g} \nabla f_i(x_k), \\H_k &:= \frac{1}{m_k^H} \sum_{i \in S_k^H} \nabla^2 f_i(x_k).\end{aligned}$$

Stochastic Contracting Newton:

$$x_{k+1} = \operatorname{argmin}_y \left\{ \langle g_k, y - x_k \rangle + \frac{1}{2} \langle H_k(y - x_k), y - x_k \rangle + \gamma_k \psi(x_k + \frac{1}{\gamma_k}(y - x_k)) \right\}.$$

Theorem. At iteration k , set $m_k^g = \mathcal{O}(k^4)$, $m_k^H = \mathcal{O}(k^2)$. Then,

$$\mathbb{E}[F(x_k) - F^*] \leq \mathcal{O}(1/k^2).$$

Variance Reduction

- ▶ **Idea:** at some iterations, recompute the full gradient [Schmidt-Roux-Bach, 2017].

$$\hat{g}_k := \frac{1}{m_k^g} \sum_{i \in S_k^g} (\nabla f_i(x_k) - \nabla f_i(z_k) + \nabla f(z_k)),$$
$$H_k := \frac{1}{m_k^H} \sum_{i \in S_k^H} \nabla^2 f_i(x_k),$$

where z_k is being updated not often.

$$z_k := x_{\pi(k)}, \quad \pi(k) \stackrel{\text{def}}{=} \begin{cases} 2^{\lfloor \log_2 k \rfloor}, & k > 0 \\ 0, & k = 0. \end{cases}$$

- ▶ During N iterations, we recompute the full gradient only $\log_2 N$ times.

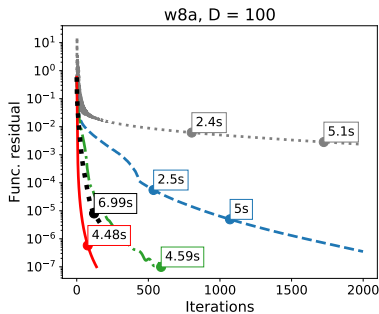
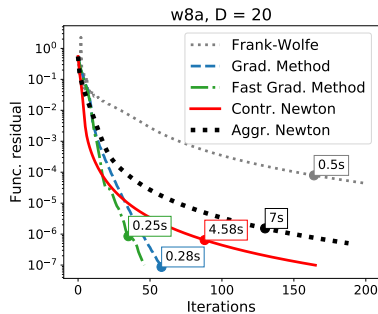
Theorem. It is enough to set $m_k^g = m_k^H = \mathcal{O}(k^2)$. Then we have

$$\mathbb{E}[F(x_k) - F^*] \leq \mathcal{O}(1/k^2).$$

Experiments: Logistic Regression

$$\min_{\|x\|_2 \leq \frac{D}{2}} \sum_{i=1}^M f_i(x), \quad f_i(x) = \log(1 + \exp(\langle a_i, x \rangle)).$$

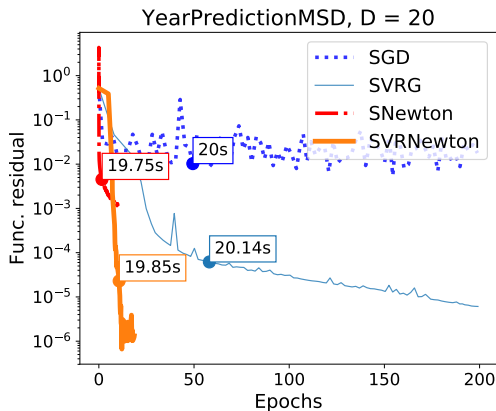
D plays the role of **regularization parameter**.



For bigger D the problem becomes more *ill-conditioned*.

Stochastic Methods for Logistic Regression

Approximate $\nabla f(x)$, $\nabla^2 f(x)$ by stochastic estimates.



The problem with big dataset size ($M = 463715$) and small dimension ($n = 90$).

1. Introduction
2. Contracting-Point Methods
3. Inexact and Stochastic Algorithms
4. Conclusions

Takeaway Points

Using the contraction of the objective

$$g_k(x) := f(\gamma_k x + (1 - \gamma_k)x_k),$$

we are able to construct new algorithms for Convex Optimization, endowed with the global complexity bounds.

1. First-order Taylor's approximation \Rightarrow Frank-Wolfe algorithm.
2. Second-order approximation \Rightarrow **Contracting Newton Method**.
3. Third-order \Rightarrow ... ?

- ▶ The methods are affine-invariant (do not depend on a norm).
- ▶ There is a complementary *Proximal-Point approach*:

$$g_k(x) := f(x) + \frac{\alpha_k}{2} \|x - x_k\|^2.$$

- ▶ Lower complexity bounds?

Note: Frank-Wolfe algorithm is near-optimal for $\|\cdot\|_\infty$ -balls [Guzmán-Nemirovski, 2015].

- ▶ Implementation for $p = 3$ (the subproblem is not convex)?

Third-order Proximal-type Tensor Methods admits effective implementation [Grapiglia-Nesterov, 2019].

- ▶ Variance reduction for the Hessian.

Nikita Doikov and Yurii Nesterov. “Affine-invariant contracting-point methods for Convex Optimization”. In: *arXiv:2009.08894* (2020)

Nikita Doikov and Yurii Nesterov. “Convex optimization based on global lower second-order models”. In: *arXiv:2006.08518* (2020)

Thank you for your attention!