

Parameter-free and optimal restart schemes for first-order methods via approximate sharpness

Maksym Neyra-Nesterenko

Department of Mathematics
Simon Fraser University
mneyrane@sfu.ca

Sep 21, 2024

Collaborators



Ben Adcock
(Simon Fraser U.)



Matthew Colbrook
(U. of Cambridge)

Paper: *Restarts subject to approximate sharpness: A parameter-free and optimal scheme for first-order methods.* Foundations of Computational Mathematics (in press, 2024).

Key contribution

We use **approximate sharpness** to design a **meta-algorithm** that **accelerates** the convergence of **any** first-order optimization method.

Remarks:

1. Our approach, based on **restarts**, can be used with essentially any first-order method
2. It applies to **broad classes** of convex problems, e.g. ℓ^1 -minimization
3. We guarantee **fast decay of the objective function error** down to an underlying **error level**

Outline

Motivations for approximate sharpness

Restart schemes for approximately sharp optimization problems

Numerical example

Conclusions

Outline

Motivations for approximate sharpness

Restart schemes for approximately sharp optimization problems

Numerical example

Conclusions

General setup

Problem: Let $f : \mathbb{R}^N \rightarrow \mathbb{R}$ be proper, closed and convex, and $Q \subseteq \mathbb{R}^N$ be a closed, convex set. Consider the problem

$$\min_{\mathbf{x} \in Q} f(\mathbf{x}) \quad (*)$$

and let $\hat{\mathbf{X}}$ be its set of minimizers with function value \hat{f} .

Approximate sharpness: We assume that $(*)$ satisfies

$$\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) \leq \left(\frac{f(\mathbf{x}) - \hat{f} + g_Q(\mathbf{x}) + \eta}{\alpha} \right)^{1/\beta}, \quad \forall \mathbf{x} \in \mathbb{R}^N,$$

where $\alpha > 0$, $\eta \geq 0$ and $\beta \geq 1$, $\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) = \inf_{\mathbf{z} \in \hat{\mathbf{X}}} d(\mathbf{x}, \mathbf{z})$ for some metric d on \mathbb{R}^N and g_Q is a known function satisfying if $\text{dist}(\mathbf{x}_i, Q) \rightarrow 0$, then $g_Q(\mathbf{x}_0) \rightarrow 0$.

General setup

Problem: Let $f : \mathbb{R}^N \rightarrow \mathbb{R}$ be proper, closed and convex, and $Q \subseteq \mathbb{R}^N$ be a closed, convex set. Consider the problem

$$\min_{\mathbf{x} \in Q} f(\mathbf{x}) \quad (*)$$

and let $\hat{\mathbf{X}}$ be its set of minimizers with function value \hat{f} .

Approximate sharpness: We assume that $(*)$ satisfies

$$\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) \leq \left(\frac{f(\mathbf{x}) - \hat{f} + g_Q(\mathbf{x}) + \eta}{\alpha} \right)^{1/\beta}, \quad \forall \mathbf{x} \in \mathbb{R}^N,$$

where $\alpha > 0$, $\eta \geq 0$ and $\beta \geq 1$, $\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) = \inf_{\mathbf{z} \in \hat{\mathbf{X}}} d(\mathbf{x}, \mathbf{z})$ for some metric d on \mathbb{R}^N and g_Q is a known function satisfying if $\text{dist}(\mathbf{x}_i, Q) \rightarrow 0$, then $g_Q(\mathbf{x}_0) \rightarrow 0$.

Related work - sharpness

Approximate sharpness generalizes the well-known condition

$$\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) \leq \left(\frac{f(\mathbf{x}) - \hat{f}}{\alpha} \right)^{1/\beta}, \quad \forall \mathbf{x} \in \mathbb{R}^N,$$

dubbed **sharpness**, **Hölderian growth** or **Lojasiewicz-type inequality**.

Hoffman (1952), Lojasiewicz (1963), Robinson (1975), Mangasarian (1985), Auslender & Crouzeix (1988), Burke & Ferris (1993), Burke & Deng (2002), Bolte, A. Daniilidis & Lewis (2007), ...

Various works have used these conditions to quantify/accelerate convergence:

Nemirovskii & Nesterov (1985), Attouch, Bolte, Redont & Soubeyran (2010), Bolte, Nguyen, Peypouquet & Suter (2017), Bolte, Sabah & Teboulle (2014), Frankel, Garrigos & Peypouquet (2015), Karimi, Nutini & Schmidt (2016), Kerdreux, d'Aspremont & Pokutta (2019), ...

Recent works employing restart schemes specifically:

Roulet & d'Aspremont (2020), Roulet, Boumal & d'Aspremont (2020), Renegar & Grimmer (2021)

Approximate sharpness and constants

$$\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) \leq \left(\frac{f(\mathbf{x}) - \hat{f} + g_Q(\mathbf{x}) + \eta}{\alpha} \right)^{1/\beta}, \quad \forall \mathbf{x} \in \mathbb{R}^N.$$

Generalizations:

- We allow $\eta > 0$.
- We incorporate a **feasibility gap function** g_Q , which means the optimization method need not produce feasible iterates.
- We **do not assume** the constants α , β and η are known to apply our restart scheme.

Motivation from compressed sensing

Compressed sensing concerns the recovery of (approximately) **sparse vectors** from incomplete sets of noisy, linear measurements.

Typical setup:

- The vector $\mathbf{x} \in \mathbb{C}^N$ to recover
- Measurement matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ (often with $m \ll N$)
- Linear measurements $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e} \in \mathbb{C}^m$, where $\mathbf{e} \in \mathbb{C}^m$ is noise
- **Goal:** Recover the vector \mathbf{x} from the measurements \mathbf{y}

Sparsity and ℓ^1 -minimization

Sparsity: \mathbf{x} is s -sparse if it has at most s nonzero entries.

Approximate sparsity: “ $\sigma_s(\mathbf{x})_1 := \min\{\|\mathbf{x} - \mathbf{z}\|_1 : \mathbf{z} \text{ is } s\text{-sparse}\}$ is small”.

Standard approaches to recover (approximately) sparse \mathbf{x} in compressed sensing involve ℓ^1 -minimization, e.g. *Quadratically Constrained Basis Pursuit (QCBP)*

$$\min_{\mathbf{z} \in \mathbb{C}^N} \|\mathbf{z}\|_1 \text{ subject to } \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 \leq \varsigma.$$

Equivalent to:

$$\min_{\mathbf{z} \in Q} f(\mathbf{z}), \quad f(\mathbf{z}) = \|\mathbf{z}\|_1, \quad Q = \{\mathbf{z} : \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 \leq \varsigma\}.$$

Sparsity and ℓ^1 -minimization

Sparsity: \mathbf{x} is s -sparse if it has at most s nonzero entries.

Approximate sparsity: “ $\sigma_s(\mathbf{x})_1 := \min\{\|\mathbf{x} - \mathbf{z}\|_1 : \mathbf{z} \text{ is } s\text{-sparse}\}$ is small”.

Standard approaches to recover (approximately) sparse \mathbf{x} in compressed sensing involve ℓ^1 -minimization, e.g. *Quadratically Constrained Basis Pursuit (QCBP)*

$$\min_{\mathbf{z} \in \mathbb{C}^N} \|\mathbf{z}\|_1 \text{ subject to } \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 \leq \varsigma.$$

Equivalent to:

$$\min_{\mathbf{z} \in Q} f(\mathbf{z}), \quad f(\mathbf{z}) = \|\mathbf{z}\|_1, \quad Q = \{\mathbf{z} : \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 \leq \varsigma\}.$$

Sparsity and ℓ^1 -minimization

Sparsity: \mathbf{x} is s -sparse if it has at most s nonzero entries.

Approximate sparsity: “ $\sigma_s(\mathbf{x})_1 := \min\{\|\mathbf{x} - \mathbf{z}\|_1 : \mathbf{z} \text{ is } s\text{-sparse}\}$ is small”.

Standard approaches to recover (approximately) sparse \mathbf{x} in compressed sensing involve ℓ^1 -minimization, e.g. *Quadratically Constrained Basis Pursuit (QCBP)*

$$\min_{\mathbf{z} \in \mathbb{C}^N} \|\mathbf{z}\|_1 \text{ subject to } \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 \leq \varsigma.$$

Equivalent to:

$$\min_{\mathbf{z} \in Q} f(\mathbf{z}), \quad f(\mathbf{z}) = \|\mathbf{z}\|_1, \quad Q = \{\mathbf{z} : \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 \leq \varsigma\}.$$

Compressed sensing theory

Definition (Restricted Isometry Property)

Let $1 \leq s \leq N$. The s th *Restricted Isometry Constant (RIC)* δ_s of a matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ is the smallest $\delta \geq 0$ such that

$$(1 - \delta)\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta)\|\mathbf{x}\|_2^2,$$

for all s -sparse vectors \mathbf{x} . If $0 \leq \delta_s \leq 1$, then \mathbf{A} is said to have the *Restricted Isometry Property (RIP)* of order s .

Intuition: \mathbf{A} approximately preserves the norm of any s -sparse vector.

Adcock & Hansen (2021), Foucart & Rauhut (2013)

Approximate sharpness in compressed sensing

Lemma

Suppose that $\mathbf{A} \in \mathbb{C}^{m \times N}$ has the RIP of order $2s$ with constant $\delta = \delta_{2s} < \sqrt{2} - 1$. Then the QCBP problem satisfies

$$\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) \leq \left(\frac{f(\mathbf{x}) - \hat{f} + g_Q(\mathbf{x}) + \eta}{\alpha} \right)^{1/\beta}, \quad \forall \mathbf{x} \in \mathbb{C}^N,$$

where $g_Q(\mathbf{z}) = \sqrt{s} \max\{\|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 - \varsigma, 0\}$, $\alpha = C_1\sqrt{s}$, $\beta = 1$, $\eta = C_2\sigma_s(\mathbf{x})_1 + C_3\sqrt{s}\varsigma$, and the constants C_1, C_2, C_3 depend on δ only.

Approximate sharpness: the distance to $\hat{\mathbf{X}}$ is bounded by:

- the **error** in the objective function $f(\mathbf{x}) - \hat{f}$
- the **feasibility gap** $g_Q(\mathbf{x})$
- the underlying **compressed sensing error** η

Approximate sharpness in compressed sensing

$$\text{dist}(\mathbf{x}, \hat{\mathbf{X}}) \leq \left(\frac{f(\mathbf{x}) - \hat{f} + g_Q(\mathbf{x}) + \eta}{\alpha} \right)^{1/\beta}, \quad \forall \mathbf{x} \in \mathbb{R}^N$$

In the compressed sensing example with

$$\alpha = C_1 \sqrt{s}, \quad \beta = 1, \quad \eta = C_2 \sigma_s(\mathbf{x})_1 + C_3 \sqrt{s} \varsigma.$$

- the order s of the RIP may be **unknown**
- $\sigma_s(\mathbf{x})_1$ is typically **unknown**
- C_1, C_2, C_3 depend on the RIC δ
- moreover, given \mathbf{A} and s , finding the RIC δ is **NP-hard**

Outline

Motivations for approximate sharpness

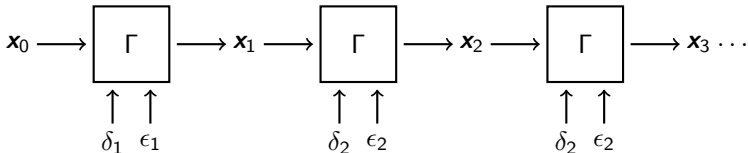
Restart schemes for approximately sharp optimization problems

Numerical example

Conclusions

Restart scheme

Let Γ be a first-order method that takes input $(\mathbf{x}, \delta, \epsilon) \in \mathbb{C}^N \times \mathbb{R}_+ \times \mathbb{R}_+$.



- Run multiple instances of Γ , where the **output** \mathbf{x}_k of the k th instance is used as the **input** of the $(k + 1)$ th instance
- Update the **parameters** $(\delta, \epsilon) = (\delta_{k+1}, \epsilon_{k+1})$ using the approximate sharpness condition
- Restarts can be extended to perform a **grid search** over α or β if their values are **unknown**, while **preserving the order of convergence**

Restart scheme for unknown α, β, η

Algorithm 2: Restart scheme for unknown α, β and η in (1.2) via grid search.

Input : Optimization algorithm Γ for (1.1), bijection ϕ as in Definition 3.1, initial vector $x^{(0)} \in D$, upper bound ϵ_0 such that $f(x^{(0)}) - \hat{f} + g_Q(x^{(0)}) \leq \epsilon_0$, constants $a, b > 1, r \in (0, 1), \alpha_0 > 0, \beta_0 \geq 1$ and total number of inner iterations $t \in \mathbb{N}$.

Output: Final iterate $x^{(t)}$ approximating a solution to (1.1).

```

1 Initialize  $x^{(0)} = x_0, U_{i,j} = 0, V_{i,j} = 0, \epsilon_{i,j,0} = \epsilon_0$  for all  $i \in \mathbb{Z}, j \in \mathbb{N}_0$ ;
2 for  $m = 0, 1, \dots, t - 1$  do
3    $(i, j, k) \leftarrow \phi(m + 1)$ ;
4    $\alpha_i \leftarrow a^i \alpha_0, \beta_j \leftarrow b^j \beta_0, U \leftarrow U_{i,j}, V \leftarrow V_{i,j}$ ;
5    $\epsilon_{i,j,U+1} \leftarrow r \epsilon_{i,j,U}$ ;
6   if  $2\epsilon_{i,j,U} > \alpha_i$  then
7      $\delta_{i,j,U+1} \leftarrow \left( \frac{2\epsilon_{i,j,U}}{\alpha_i} \right)^{\min\{b/\beta_j, 1/\beta_0\}}$ ;
8   else
9      $\delta_{i,j,U+1} \leftarrow \left( \frac{2\epsilon_{i,j,U}}{\alpha_i} \right)^{1/\beta_j}$ ;
10  end
11  if  $V + C_\Gamma(\delta_{i,j,U+1}, \epsilon_{i,j,U+1}) \leq k$  then
12     $z^{(m)} \leftarrow \Gamma(\delta_{i,j,U+1}, \epsilon_{i,j,U+1}, x^{(m)})$ ;
13     $x^{(m+1)} \leftarrow \operatorname{argmin}\{f(x) + g_Q(x) : x = z^{(m)} \text{ or } x = x^{(m)}\}$ ;
14     $V_{i,j} \leftarrow V + C_\Gamma(\delta_{i,j,U+1}, \epsilon_{i,j,U+1})$ ;
15     $U_{i,j} \leftarrow U + 1$ ;
16  else
17     $x^{(m+1)} = x^{(m)}$ ;
18  end
19 end

```

Theorem (Unknown α, β, η)

Suppose the number of iterations computed by Γ is at most $C\delta^{d_1}/\epsilon^{d_2} + 1$, for all $\delta, \epsilon > 0$. Then there is a restart scheme such that after at most

$$\hat{C} \cdot \epsilon^{d_1/\beta_* - d_2} \cdot \begin{cases} \lceil \log(1/\epsilon) \rceil, & \text{if } d_2 \leq d_1/\beta_*, \\ 1, & \text{if } d_2 > d_1/\beta_*, \end{cases}$$

iterations of Γ , where β_* is the scheme's closest grid point to β , the restart scheme produces an output \mathbf{x}^* satisfying

$$f(\mathbf{x}^*) - \hat{f} + g_Q(\mathbf{x}^*) \leq \max\{\epsilon, \eta\}.$$

Here \hat{C} depends on C, α, β_*, d_1 and d_2 .

Remark: There are restart schemes for the special cases when α or β are known, and analogous results can be stated.

Adcock, Colbrook & Neyra-Nesterenko (2024)

Rates for different problem classes

Objective function class/structure	Asymptotic bound for $K(\varepsilon)$	Example method
L -smooth See Definition 4.2 (NB: must have $\beta \geq 2$)	$\beta = 2$: $\sqrt{L/\alpha} \cdot \log(1/\varepsilon)$ <hr style="border-top: 1px dashed #ccc;"/> $\beta > 2$: $\frac{\sqrt{L}}{\alpha^{1/\beta_*}} \cdot \frac{1}{\varepsilon^{1/2-1/\beta_*}}$	Nesterov's method $d_1 = 1, d_2 = 1/2$ See Section 4.1
(u, v) -smoothable See Definition 4.5	$\beta = 1$: $\frac{\sqrt{ab}}{\alpha} \cdot \log(1/\varepsilon)$ <hr style="border-top: 1px dashed #ccc;"/> $\beta > 1$: $\frac{\sqrt{ab}}{\alpha^{1/\beta_*}} \cdot \frac{1}{\varepsilon^{1-1/\beta_*}}$	Nesterov's method with smoothing $d_1 = 1, d_2 = 1$ See Section 4.2
Hölder smooth, parameter $\nu \in [0, 1]$ See Definition 4.8 (NB: must have $\beta \geq 1 + \nu$)	$\beta = 1 + \nu$: $\frac{M\nu^{\frac{2}{1+3\nu}}}{\alpha^{(1+3\nu)}} \cdot \log(1/\varepsilon)$ <hr style="border-top: 1px dashed #ccc;"/> $\beta > 1 + \nu$: $\frac{M\nu^{\frac{2}{1+3\nu}}}{\alpha^{\beta_*(1+3\nu)}} \cdot \frac{1}{\varepsilon^{\frac{2(\beta_*-1-\nu)}{\beta_*(1+3\nu)}}}$	Universal fast gradient method $d_1 = (2 + 2\nu)/(1 + 3\nu)$ $d_2 = 2/(1 + 3\nu)$ See Section 4.3
$f(x) = q(x) + g(x) + h(Bx)$, q is L_q -smooth, $\sup_{x \in \text{dom}(h)} \inf_{y \in \partial h(x)} \ y\ \leq L_h$, $\ B\ \leq L_B$	$\beta = 1$: $\frac{L_B L_h + L_q}{\alpha} \cdot \log(1/\varepsilon)$ <hr style="border-top: 1px dashed #ccc;"/> $\beta > 1$: $\frac{L_B L_h + L_q}{\alpha^{1/\beta_*}} \cdot \frac{1}{\varepsilon^{1-1/\beta_*}}$	Primal-dual algorithm $d_1 = 1, d_2 = 1$ See Section 4.4
$f(x) = q(x) + g(x) + h(Bx)$, q is L_q -smooth, $\sup_{x \in \text{dom}(h)} \inf_{y \in \partial h(x)} \ y\ \leq L_h$, $\ A\ \leq L_A, \ B\ \leq L_B$, $Q = \{x : Ax \in C\}, g_Q(x) = \kappa \inf_{z \in C} \ Ax - z\ $	$\beta = 1$: $\frac{\kappa L_A + L_B L_h + L_q}{\alpha} \cdot \log(1/\varepsilon)$ <hr style="border-top: 1px dashed #ccc;"/> $\beta > 1$: $\frac{\kappa L_A + L_B L_h + L_q}{\alpha^{1/\beta_*}} \cdot \frac{1}{\varepsilon^{1-1/\beta_*}}$	Primal-dual algorithm with constraints $d_1 = 1, d_2 = 1$ See Section 4.5

Table 1: Asymptotic cost bounds (as $\varepsilon \downarrow 0$ for $\eta \lesssim \varepsilon$) and suitable first-order methods for Algorithm 2 when applied to different classes of objective functions. Note also that whenever the bound is a polynomial in $\log(1/\varepsilon)$, we have $\beta_* = \beta$.

Remark: The first three lines constitute **optimal rates**.

Outline

Motivations for approximate sharpness

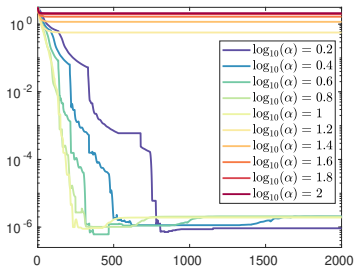
Restart schemes for approximately sharp optimization problems

Numerical example

Conclusions

Compressed sensing example

If $d_1 = d_2/\beta$, then the cost bound reduces to $\hat{C} \cdot \log(1/\varepsilon)$, yielding **linear decay** to η .



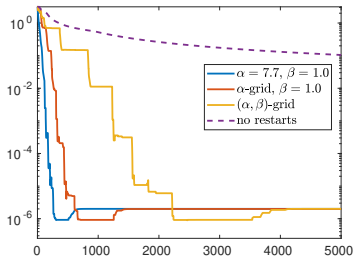
Recovery error vs. # iterations

Example: Γ is primal-dual iteration applied to QCBP, where $\beta = d_1 = d_2 = 1$.

This is applied to our compressed sensing problem, where \mathbf{A} is a Gaussian random matrix and $\varsigma = 10^{-6}$. The ground truth \mathbf{x} is exactly sparse, hence $\eta \approx \varsigma$.

Issue: Restarts are **brittle** with respect to fixed α .

Compressed sensing example



Recovery error vs. # iterations

A direct comparison of restart schemes with tuned constants and the nonrestarted optimization method Γ (primal-dual iterations).

Grid searching maintains linear decay and still outperforms the non-restarted optimization method.

Outline

Motivations for approximate sharpness

Restart schemes for approximately sharp optimization problems

Numerical example

Conclusions

Conclusions

Convex optimization problems arising in applications, such as image and signal reconstruction, matrix completion, feature selection) satisfy an **approximate sharpness** condition with **unknown** constants.

In this setting, our goal is to obtain **fast convergence** down to the **(unknown) approximate sharpness constant η** .

We introduced an algorithm for accelerating any convex optimization method, based on **restarts** and **grid searching**.

This leads to **optimal** rates for various convex optimization problems and competitive practical performance.

Paper: *Restarts subject to approximate sharpness: A parameter-free and optimal scheme for first-order methods*. Foundations of Computational Mathematics (in press, 2024).