

A superlinearly convergent subgradient method for sharp semismooth problems

Vasilis Charisopoulos

Joint work with Damek Davis

International Conference on Continuous Optimization, 2022

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

under the following assumptions:

- f nonsmooth, locally Lipschitz.
- f satisfies classical **sharp growth** condition:

$$f(x) \geq \mu \operatorname{dist}(x, \mathcal{X}_*), \quad \mathcal{X}_* = \operatorname{argmin} f.$$

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

under the following assumptions:

- f nonsmooth, locally Lipschitz.
- f satisfies classical **sharp growth** condition:

$$f(x) \geq \mu \operatorname{dist}(x, \mathcal{X}_*), \quad \mathcal{X}_* = \operatorname{argmin} f.$$

Classical example: Hoffman bound for LPs / linear inequalities.¹

$$\operatorname{dist}(x, \{x \mid Ax \leq b\}) \leq H_A \|(Ax - b)_+\|$$

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

under the following assumptions:

- f nonsmooth, locally Lipschitz.
- f satisfies classical **sharp growth** condition:

$$f(x) \geq \mu \operatorname{dist}(x, \mathcal{X}_*), \quad \mathcal{X}_* = \operatorname{argmin} f.$$

Root-finding problems:

$$\operatorname{find} x \text{ s.t. } F(x) = 0 \Leftrightarrow \operatorname{argmin}_x \|F(x)\|.$$

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

under the following assumptions:

- f nonsmooth, locally Lipschitz.
- f satisfies classical **sharp growth** condition:

$$f(x) \geq \mu \operatorname{dist}(x, \mathcal{X}_*), \quad \mathcal{X}_* = \operatorname{argmin} f.$$

Root-finding problems:

$$\text{find } x \text{ s.t. } F(x) = 0 \Leftrightarrow \operatorname{argmin}_x \|F(x)\|.$$

Growth condition known as *metric subregularity*:²

$$\|F(x)\| \geq \mu \operatorname{dist}(x, \mathcal{X}_*).$$

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

under the following assumptions:

- f nonsmooth, locally Lipschitz.
- f satisfies classical **sharp growth** condition:

$$f(x) \geq \mu \operatorname{dist}(x, \mathcal{X}_*), \quad \mathcal{X}_* = \operatorname{argmin} f.$$

Set intersection problems: given $\mathcal{X}_1, \mathcal{X}_2$ closed,

$$\text{find } \bar{x} \in \mathcal{X}_1 \cap \mathcal{X}_2 \Leftrightarrow \operatorname{argmin}_x \{\operatorname{dist}(x, \mathcal{X}_1) + \operatorname{dist}(x, \mathcal{X}_2)\}.$$

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

under the following assumptions:

- f nonsmooth, locally Lipschitz.
- f satisfies classical **sharp growth** condition:

$$f(x) \geq \mu \operatorname{dist}(x, \mathcal{X}_*), \quad \mathcal{X}_* = \operatorname{argmin} f.$$

Set intersection problems: given $\mathcal{X}_1, \mathcal{X}_2$ closed,

$$\text{find } \bar{x} \in \mathcal{X}_1 \cap \mathcal{X}_2 \Leftrightarrow \operatorname{argmin}_x \{ \operatorname{dist}(x, \mathcal{X}_1) + \operatorname{dist}(x, \mathcal{X}_2) \}.$$

Growth condition known as *linear regularity*:

$$\operatorname{dist}(x, \mathcal{X}_1) + \operatorname{dist}(x, \mathcal{X}_2) \geq \mu \operatorname{dist}(x, \mathcal{X}_1 \cap \mathcal{X}_2).$$

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Setting

Goal: fast first-order algorithms for

$$\operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \quad \text{where} \quad \min f = 0,$$

under the following assumptions:

- f nonsmooth, locally Lipschitz.
- f satisfies classical **sharp growth** condition:

$$f(x) \geq \mu \operatorname{dist}(x, \mathcal{X}_*), \quad \mathcal{X}_* = \operatorname{argmin} f.$$

Fast algorithms?

¹Hoffman '52.

²J.S. Pang '93; Ioffe '80s.

Sharp growth and rapid convergence

Subgradient method:

$$x_{k+1} := x_k - \alpha_k \frac{v_k}{\|v_k\|}, \quad v_k \in \underbrace{\partial f(x_k)}_{\text{Clarke subdifferential}}.$$

¹Polyak '69.

²Davis et al. '18.

Sharp growth and rapid convergence

Subgradient method:

$$x_{k+1} := x_k - \alpha_k \frac{v_k}{\|v_k\|}, \quad v_k \in \underbrace{\partial f(x_k)}_{\text{Clarke subdifferential}}.$$

Theorem: $\{x_k\}$ converge *linearly* with *Polyak step size*:

$$\alpha_k = \frac{f(x_k)}{\|v_k\|}.$$

¹Polyak '69.

²Davis et al. '18.

Sharp growth and rapid convergence

Subgradient method:

$$x_{k+1} := x_k - \alpha_k \frac{v_k}{\|v_k\|}, \quad v_k \in \underbrace{\partial f(x_k)}_{\text{Clarke subdifferential}}.$$

Theorem: $\{x_k\}$ converge *linearly* with *Polyak step size*:

$$\alpha_k = \frac{f(x_k)}{\|v_k\|}.$$

- Classically known for convex functions¹.
- Recently generalized for *weakly convex* functions².

Key example : $f = (\text{convex}) \circ (\text{smooth})$.

¹Polyak '69.

²Davis et al. '18.

Sharp growth and rapid convergence

Subgradient method:

$$x_{k+1} := x_k - \alpha_k \frac{v_k}{\|v_k\|}, \quad v_k \in \underbrace{\partial f(x_k)}_{\text{Clarke subdifferential}}.$$

Theorem: $\{x_k\}$ converge *linearly* with *Polyak step size*:

$$\alpha_k = \frac{f(x_k)}{\|v_k\|}.$$

- Classically known for convex functions¹.
- Recently generalized for *weakly convex* functions².

Key example : $f = (\text{convex}) \circ (\text{smooth})$.

Question: *Faster than linear convergence using only subgradients?*

¹Polyak '69.

²Davis et al. '18.

An answer in pictures

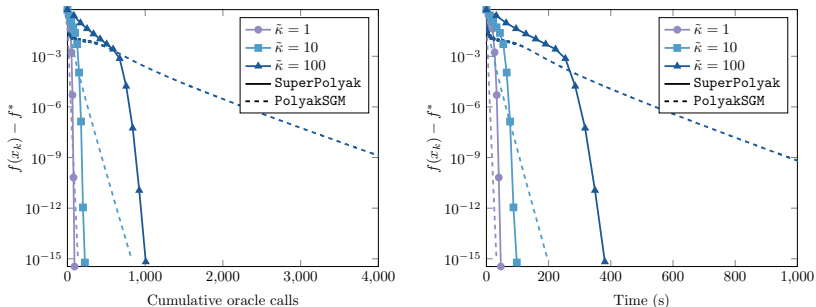


Figure: Our algorithm (SuperPolyak) applied to a matrix sensing problem with dimensions $(d, r) = (2^{15}, 2)$ and $m = 2^{19}$ measurements. Here, $\tilde{\kappa}$ is the condition number of the unknown matrix.

Algorithm converges superlinearly, with fewer subgradient oracle calls. **How?**

Motivation: Polyak bundle

Simplifying assumption: problem has *unique solution* \bar{x} .

Motivation: Polyak bundle

Simplifying assumption: problem has *unique solution* \bar{x} .

Polyak step equivalent to:

$$x_{k+1} := \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(x_k) + \langle v_k, x - x_k \rangle \leq 0 \}, \quad v_k \in \partial f(x_k).$$

Motivation: Polyak bundle

Simplifying assumption: problem has *unique solution* \bar{x} .

Polyak step equivalent to:

$$x_{k+1} := \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(x_k) + \langle v_k, x - x_k \rangle \leq 0 \}, \quad v_k \in \partial f(x_k).$$

To improve convergence, can try to use a **bundle** ($y_i, v_i \in \partial f(y_i)$):

$$x_{k+1} := \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle \leq 0, \text{ for all } i \}.$$

- **Note:** \bar{x} *always feasible* when f is convex.
- Bundle points y_i chosen among the past k iterates.
- Each step requires solving a QP.

Motivation: Polyak bundle

Simplifying assumption: problem has *unique solution* \bar{x} .

Polyak step equivalent to:

$$x_{k+1} := \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(x_k) + \langle v_k, x - x_k \rangle \leq 0 \}, \quad v_k \in \partial f(x_k).$$

To improve convergence, can try to use a **bundle** ($y_i, v_i \in \partial f(y_i)$):

$$x_{k+1} := \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle \leq 0, \text{ for all } i \}.$$

- **Note:** \bar{x} *always feasible* when f is convex.
- Bundle points y_i chosen among the past k iterates.
- Each step requires solving a QP.

Guarantees?

- Good practical performance, but rate similar to subgradient method.³

³Polyak 87'.

Our approach

Main idea.

- Replace QP by a linear system:

$$x_{k+1} = \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle = 0, \text{ for all } i \},$$

Our approach

Main idea.

- Replace QP by a linear system:

$$x_{k+1} = \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle = 0, \text{ for all } i \},$$

- Choose bundle points y_i “carefully”.

Our approach

Main idea.

- Replace QP by a linear system:

$$x_{k+1} = \underset{x}{\operatorname{argmin}} \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle = 0, \text{ for all } i \},$$

- Choose bundle points y_i “carefully”.

Problem admits **closed-form solution**:

$$x_{k+1} = x_k - A^\dagger \begin{bmatrix} f(y_0) + \langle v_0, x_k - y_0 \rangle \\ \vdots \\ f(y_i) + \langle v_i, x_k - y_i \rangle \end{bmatrix}, \quad \text{where } A = \begin{bmatrix} v_0^\top \\ \vdots \\ v_i^\top \end{bmatrix}.$$

Our approach

Main idea.

- Replace QP by a linear system:

$$x_{k+1} = \underset{x}{\operatorname{argmin}} \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle = 0, \text{ for all } i \},$$

- Choose bundle points y_i “carefully”.

Key additional assumption: **Semismoothness.**

$$f(x) + \langle v, \bar{x} - x \rangle = o(\|x - \bar{x}\|), \quad \text{for } v \in \partial f(x) \text{ and as } x \rightarrow \bar{x}.$$

\Rightarrow Implies \bar{x} is *nearly feasible* for system of equations!

Our approach

Main idea.

- Replace QP by a linear system:

$$x_{k+1} = \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle = 0, \text{ for all } i \},$$

- Choose bundle points y_i “carefully”.

Key additional assumption: **Semismoothness**.

$$f(x) + \langle v, \bar{x} - x \rangle = o(\|x - \bar{x}\|), \quad \text{for } v \in \partial f(x) \text{ and as } x \rightarrow \bar{x}.$$

⇒ Implies \bar{x} is *nearly feasible* for system of equations!

Semismoothness is common. Satisfied by:

- Convex and smooth, and (convex) \circ (smooth) functions.
- Any *semialgebraic* function.¹

¹Bolte, Daniilidis & Lewis '09.

Our approach

Main idea.

- Replace QP by a linear system:

$$x_{k+1} = \operatorname{argmin}_x \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle = 0, \text{ for all } i \},$$

- Choose bundle points y_i “carefully”.

Key additional assumption: **Semismoothness.**

$$f(x) + \langle v, \bar{x} - x \rangle = o(\|x - \bar{x}\|), \quad \text{for } v \in \partial f(x) \text{ and as } x \rightarrow \bar{x}.$$

\Rightarrow Implies \bar{x} is *nearly feasible* for system of equations!

Question: How to choose the bundle points $\{y_i\}$?

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix}$ for $v_i \in \partial f(y_i)$

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix}$ for $v_i \in \partial f(y_i)$

return y_s , where $s = \operatorname{argmin}_{i: \|y_i - y_0\| \leq \tau f(y_0)} f(y_i)$

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix}$ for $v_i \in \partial f(y_i)$

return y_s , where $s = \operatorname{argmin}_{i: \|y_i - y_0\| \leq \tau f(y_0)} f(y_i)$

Note: first bundle point recovers Polyak subgradient step:

$$y_1 = y_0 - (v_0^\top)^\dagger (f(y_0) + \langle v_0, y_0 - y_0 \rangle) = y_0 - \frac{f(y_0)}{\|v_0\|^2} v_0.$$

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix}$ for $v_i \in \partial f(y_i)$

return y_s , where $s = \operatorname{argmin}_{i: \|y_i - y_0\| \leq \tau f(y_0)} f(y_i)$

Strategy. Sharpness and semismoothness lead to “lemma of alternatives”:

1. Suppose that y_0, \dots, y_{j-1} have not improved superlinearly.

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix} \text{ for } v_i \in \partial f(y_i)$$

return y_s , where $s = \operatorname{argmin}_{i: \|y_i - y_0\| \leq \tau f(y_0)} f(y_i)$

Strategy. Sharpness and semismoothness lead to “lemma of alternatives”:

1. Suppose that y_0, \dots, y_{j-1} have not improved superlinearly.
2. Then, *either one* of the following must hold:
 - y_j improves superlinearly upon x ;
 - $\operatorname{rank}(A_j) = j + 1$.

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix} \text{ for } v_i \in \partial f(y_i)$$

return y_s , where $s = \operatorname{argmin}_{i: \|y_i - y_0\| \leq \tau f(y_0)} f(y_i)$

Strategy. Sharpness and semismoothness lead to “lemma of alternatives”:

1. Suppose that y_0, \dots, y_{j-1} have not improved superlinearly.
2. Then, *either one* of the following must hold:
 - y_j improves superlinearly upon x ;
 - $\operatorname{rank}(A_j) = j + 1$.
3. Since $A_j \in \mathbb{R}^{(j+1) \times d}$, superlinear improvement achieved within d steps.

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix} \text{ for } v_i \in \partial f(y_i)$$

return y_s , where $s = \operatorname{argmin}_{i: \|y_i - y_0\| \leq \tau f(y_0)} f(y_i)$

Theorem: C. & Davis, 2022 (informal)

Assume x near \bar{x} and τ is sufficiently large. Then

$$f(\text{PolyakBundle}(x, \tau)) = o(f(x)).$$

- ✓ Provable early termination strategies available.
- ✓ Sequence of linear systems solved incrementally (QR-based algorithm).

Algorithm PolyakBundle(x, τ)

$y_0 := x; v_0 \in \partial f(y_0); A_1 = [v_0^\top]$.

for $i = 1, \dots, d$ **do**

$$y_i := y_0 - A_i^\dagger [f(y_j) + \langle v_j, y_0 - y_j \rangle]_{j=0}^{i-1}; \quad A_{i+1} := \begin{bmatrix} A_i \\ v_i^\top \end{bmatrix} \text{ for } v_i \in \partial f(y_i)$$

return y_s , where $s = \operatorname{argmin}_{i: \|y_i - y_0\| \leq \tau f(y_0)} f(y_i)$

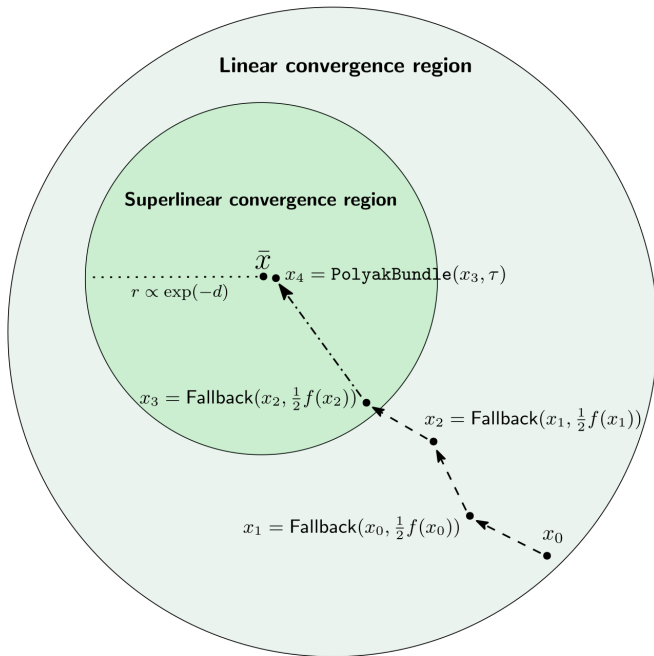
Theorem: C. & Davis, 2022 (informal)

Assume x near \bar{x} and τ is sufficiently large. Then

$$f(\text{PolyakBundle}(x, \tau)) = o(f(x)).$$

- ✓ Provable early termination strategies available.
- ✓ Sequence of linear systems solved incrementally (QR-based algorithm).
- ✗ **Issue:** region of local convergence $\propto \exp(-d)$.
 - **Fix:** couple with linearly convergent method (e.g., subgradient).

High-level overview



Algorithm: SuperPolyak

Input: $x_0 \in \mathbb{R}^d$.

Repeat for $k = 0, 1, \dots$

Algorithm: SuperPolyak

Input: $x_0 \in \mathbb{R}^d$.

Repeat for $k = 0, 1, \dots$

 Attempt a superlinear step:

$$\tilde{x} = \text{PolyakBundle}(x_k, (3/2)^k)$$

Algorithm: SuperPolyak

Input: $x_0 \in \mathbb{R}^d$.

Repeat for $k = 0, 1, \dots$

Attempt a superlinear step:

$$\tilde{x} = \text{PolyakBundle}(x_k, (3/2)^k)$$

If \tilde{x} available **and** sufficient decrease:

$$f(\tilde{x}) \leq \frac{1}{2}f(x_k) \implies x_{k+1} := \tilde{x}.$$

Algorithm: SuperPolyak

Input: $x_0 \in \mathbb{R}^d$.

Repeat for $k = 0, 1, \dots$

Attempt a superlinear step:

$$\tilde{x} = \text{PolyakBundle}(x_k, (3/2)^k)$$

If \tilde{x} available **and** sufficient decrease:

$$f(\tilde{x}) \leq \frac{1}{2}f(x_k) \implies x_{k+1} := \tilde{x}.$$

Else: run Fallback method from x_k until objective halved.

$$x_{k+1} := \text{Fallback}\left(x_k, \frac{1}{2}f(x_k)\right).$$

Algorithm: SuperPolyak

Input: $x_0 \in \mathbb{R}^d$.

Repeat for $k = 0, 1, \dots$

Attempt a superlinear step:

$$\tilde{x} = \text{PolyakBundle}(x_k, (3/2)^k)$$

If \tilde{x} available **and** sufficient decrease:

$$f(\tilde{x}) \leq \frac{1}{2}f(x_k) \implies x_{k+1} := \tilde{x}.$$

Else: run Fallback method from x_k until objective halved.

$$x_{k+1} := \text{Fallback}\left(x_k, \frac{1}{2}f(x_k)\right).$$

Question: which algorithm can we use as fallback method?

Algorithm: SuperPolyak

Input: $x_0 \in \mathbb{R}^d$.

Repeat for $k = 0, 1, \dots$

Attempt a superlinear step:

$$\tilde{x} = \text{PolyakBundle}(x_k, (3/2)^k)$$

If \tilde{x} available **and** sufficient decrease:

$$f(\tilde{x}) \leq \frac{1}{2}f(x_k) \implies x_{k+1} := \tilde{x}.$$

Else: run Fallback method from x_k until objective halved.

$$x_{k+1} := \text{Fallback}\left(x_k, \frac{1}{2}f(x_k)\right).$$

Theorem: C. and Davis, 2022 (informal)

SuperPolyak with the Polyak subgradient method as fallback enters the region of superlinear convergence in $O(d)$ iterations, as long as x_0 is in a dimension-independent region around \bar{x} .

Consequences: root-finding problems

Problem: find x s.t. $F(x) = 0$.

Consequences: root-finding problems

Problem: find x s.t. $F(x) = 0$.

Classical guarantees. Suppose that F satisfies:

1. **Semismoothness:** for all x near \bar{x} and $A \in \partial F(x)$,

$$\|F(x) + A(\bar{x} - x)\| = o(\|x - \bar{x}\|)$$

Consequences: root-finding problems

Problem: find x s.t. $F(x) = 0$.

Classical guarantees. Suppose that F satisfies:

1. **Semismoothness:** for all x near \bar{x} and $A \in \partial F(x)$,

$$\|F(x) + A(\bar{x} - x)\| = o(\|x - \bar{x}\|)$$

2. **Invertibility:** Clarke Jacobian $\partial F(x)$ invertible for x near \bar{x} .

Consequences: root-finding problems

Problem: find x s.t. $F(x) = 0$.

Classical guarantees. Suppose that F satisfies:

1. **Semismoothness:** for all x near \bar{x} and $A \in \partial F(x)$,

$$\|F(x) + A(\bar{x} - x)\| = o(\|x - \bar{x}\|)$$

2. **Invertibility:** Clarke Jacobian $\partial F(x)$ invertible for x near \bar{x} .

Then the *semismooth Newton* method converges locally superlinearly:¹

$$x_{k+1} := x_k - A_k^{-1} F(x_k), \quad A_k \in \partial F(x_k).$$

¹Qi & Sun, '93.

Consequences: root-finding problems

Problem: find x s.t. $F(x) = 0$.

Classical guarantees. Suppose that F satisfies:

1. **Semismoothness:** for all x near \bar{x} and $A \in \partial F(x)$,

$$\|F(x) + A(\bar{x} - x)\| = o(\|x - \bar{x}\|)$$

2. **Invertibility:** Clarke Jacobian $\partial F(x)$ invertible for x near \bar{x} .

Then the *semismooth Newton* method converges locally superlinearly:¹

$$x_{k+1} := x_k - A_k^{-1} F(x_k), \quad A_k \in \partial F(x_k).$$

Question: is superlinear convergence possible *without* invertibility condition?

¹Qi & Sun, '93.

Consequences: root-finding problems

Problem: find x s.t. $F(x) = 0$.

Assumptions. Suppose that F satisfies:

1. **Semismoothness:** for all x near \bar{x} and $A(x) \in \partial F(x)$,

$$\|F(x) + A(\bar{x} - x)\| = o(\|x - \bar{x}\|)$$

2. **Metric subregularity:** for x near \bar{x} ,

$$\|F(x)\| \geq \mu \|x - \bar{x}\|.$$

Consequences: root-finding problems

Problem: find x s.t. $F(x) = 0$.

Assumptions. Suppose that F satisfies:

1. **Semismoothness:** for all x near \bar{x} and $A(x) \in \partial F(x)$,

$$\|F(x) + A(\bar{x} - x)\| = o(\|x - \bar{x}\|)$$

2. **Metric subregularity:** for x near \bar{x} ,

$$\|F(x)\| \geq \mu \|x - \bar{x}\|.$$

Corollary: C. and Davis, 2022 (informal)

Under above assumptions, SuperPolyak converges locally superlinearly.

Natural fallback method: fixed-point iteration

$$z_{i+1} := T(z_i), \quad \text{where } T := I - F.$$

Example: LASSO

Goal: recover s -sparse $x_{\#} \in \mathbb{R}^d$ from $y = Ax_{\#} \in \mathbb{R}^m$ ($m \ll d$).

$$\operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1 \right\}$$

Example: LASSO

Goal: recover s -sparse $x_{\#} \in \mathbb{R}^d$ from $y = Ax_{\#} \in \mathbb{R}^m$ ($m \ll d$).

$$\operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1 \right\}$$

Fallback algorithm: proximal gradient (ISTA).

$$x_{k+1} = T(x_k) := \operatorname{prox}_{\lambda \|\cdot\|_1} (x_k - \tau A^T (Ax_k - y))$$
$$f(x) := \|(I - T)(x)\|.$$

Example: LASSO

Goal: recover s -sparse $x_{\#} \in \mathbb{R}^d$ from $y = Ax_{\#} \in \mathbb{R}^m$ ($m \ll d$).

$$\operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1 \right\}$$

Fallback algorithm: proximal gradient (ISTA).

$$x_{k+1} = T(x_k) := \operatorname{prox}_{\lambda \|\cdot\|_1} (x_k - \tau A^T (Ax_k - y))$$
$$f(x) := \|(I - T)(x)\|.$$

Note: $I - T$ metrically subregular but need not satisfy invertibility condition!

Example: LASSO

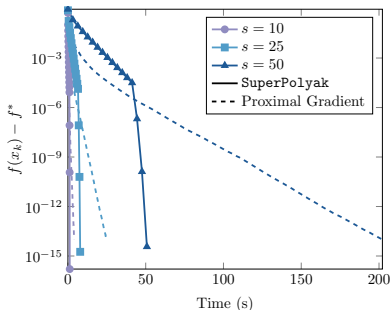
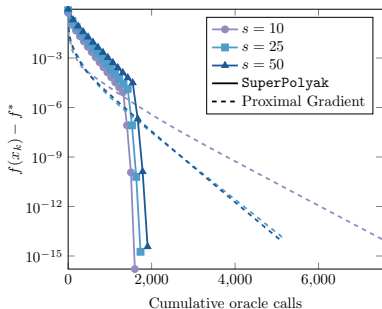
Goal: recover s -sparse $x_{\#} \in \mathbb{R}^d$ from $y = Ax_{\#} \in \mathbb{R}^m$ ($m \ll d$).

$$\operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1 \right\}$$

Fallback algorithm: proximal gradient (ISTA).

$$x_{k+1} = T(x_k) := \operatorname{prox}_{\lambda \|\cdot\|_1} (x_k - \tau A^T(Ax_k - y))$$

$$f(x) := \|(I - T)(x)\|.$$



Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

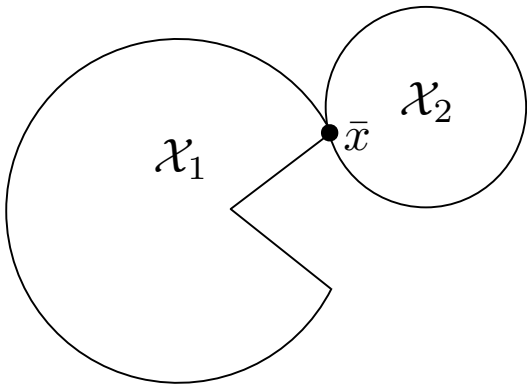


Figure: Semialgebraic sets intersecting at a single point \bar{x} .

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

Setting I: intersections of **semialgebraic sets**.

1. The family $\{\mathcal{X}_1, \mathcal{X}_2\}$ is μ -linearly regular:

$$\text{dist}(x, \mathcal{X}_1) + \text{dist}(x, \mathcal{X}_2) \geq \mu \text{dist}(x, \mathcal{X}_*) \text{ for all } x \text{ near } \bar{x}.$$

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

Setting I: intersections of **semialgebraic sets**.

1. The family $\{\mathcal{X}_1, \mathcal{X}_2\}$ is μ -linearly regular:

$$\text{dist}(x, \mathcal{X}_1) + \text{dist}(x, \mathcal{X}_2) \geq \mu \text{dist}(x, \mathcal{X}_*) \text{ for all } x \text{ near } \bar{x}.$$

2. Every \mathcal{X}_i is semialgebraic and $\mathcal{X}_* = \{\bar{x}\}$.

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

Setting I: intersections of **semialgebraic sets**.

1. The family $\{\mathcal{X}_1, \mathcal{X}_2\}$ is μ -linearly regular:

$$\text{dist}(x, \mathcal{X}_1) + \text{dist}(x, \mathcal{X}_2) \geq \mu \text{dist}(x, \mathcal{X}_*) \text{ for all } x \text{ near } \bar{x}.$$

2. Every \mathcal{X}_i is semialgebraic and $\mathcal{X}_* = \{\bar{x}\}$.

Corollary: C. and Davis, 2022 (informal)

Under above assumptions, SuperPolyak converges locally superlinearly.

Natural fallback method: alternating projections algorithm.¹

$$z_{i+1} := \mathcal{P}_{\mathcal{X}_2}(\mathcal{P}_{\mathcal{X}_1}(z_i)).$$

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

Setting II: intersections of **manifolds**.

1. The family $\{\mathcal{X}_1, \mathcal{X}_2\}$ is μ -linearly regular:

$$\text{dist}(x, \mathcal{X}_1) + \text{dist}(x, \mathcal{X}_2) \geq \mu \text{dist}(x, \mathcal{X}_*) \text{ for all } x \text{ near } \bar{x}.$$

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

Setting II: intersections of **manifolds**.

1. The family $\{\mathcal{X}_1, \mathcal{X}_2\}$ is μ -linearly regular:

$$\text{dist}(x, \mathcal{X}_1) + \text{dist}(x, \mathcal{X}_2) \geq \mu \text{dist}(x, \mathcal{X}_*) \text{ for all } x \text{ near } \bar{x}.$$

2. Every \mathcal{X}_i is a C^2 manifold near \bar{x} .

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

Setting II: intersections of **manifolds**.

1. The family $\{\mathcal{X}_1, \mathcal{X}_2\}$ is μ -linearly regular:

$$\text{dist}(x, \mathcal{X}_1) + \text{dist}(x, \mathcal{X}_2) \geq \mu \text{dist}(x, \mathcal{X}_*) \text{ for all } x \text{ near } \bar{x}.$$

2. Every \mathcal{X}_i is a C^2 manifold near \bar{x} .

Corollary: C. and Davis, 2022 (informal)

Under above assumptions, SuperPolyak converges locally *quadratically*.

Natural fallback method: alternating projections algorithm.²

$$z_{i+1} := \mathcal{P}_{\mathcal{X}_2}(\mathcal{P}_{\mathcal{X}_1}(z_i)).$$

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Consequences: feasibility problems

Problem: find $\bar{x} \in \mathcal{X}_* = \mathcal{X}_1 \cap \mathcal{X}_2$, \mathcal{X}_1 and \mathcal{X}_2 closed.

Setting II: intersections of **manifolds**.

1. The family $\{\mathcal{X}_1, \mathcal{X}_2\}$ is μ -linearly regular:

$$\text{dist}(x, \mathcal{X}_1) + \text{dist}(x, \mathcal{X}_2) \geq \mu \text{dist}(x, \mathcal{X}_*) \text{ for all } x \text{ near } \bar{x}.$$

2. Every \mathcal{X}_i is a C^2 manifold near \bar{x} .

Corollary: C. and Davis, 2022 (informal)

Under above assumptions, SuperPolyak converges locally *quadratically*.

Natural fallback method: alternating projections algorithm.²

$$z_{i+1} := \mathcal{P}_{\mathcal{X}_2}(\mathcal{P}_{\mathcal{X}_1}(z_i)).$$

Related work: QP-based algorithm that converges under similar conditions.³

¹Drusvyatskiy '13.

²Lewis, Luke & Malick '09.

³C.H.J. Pang '15.

Example: complex phase retrieval

Complex phase retrieval: given $y_{\#} \in \mathbb{R}^m$ with $(y_{\#})_i = |\langle a_i, x_{\#} \rangle|$:

find $\hat{y} \in \mathcal{Y}_1 \cap \mathcal{Y}_2$, $\mathcal{Y}_1 := \{u \in \mathbb{C}^m \mid |u| = y\}$, $\mathcal{Y}_2 := \text{Range}(A)$.

Example: complex phase retrieval

Complex phase retrieval: given $y_{\#} \in \mathbb{R}^m$ with $(y_{\#})_i = |\langle a_i, x_{\#} \rangle|$:

find $\hat{y} \in \mathcal{Y}_1 \cap \mathcal{Y}_2$, $\mathcal{Y}_1 := \{u \in \mathbb{C}^m \mid |u| = y\}$, $\mathcal{Y}_2 := \text{Range}(A)$.

Fallback algorithm: alternating projections⁴.

$$y_{k+1} = T(y_k) := AA^{\dagger} (y_{\#} \odot \text{phase}(y_k))$$

Function used by PolyakBundle: $f(y) := \text{dist}(y, \mathcal{Y}_1) + \text{dist}(y, \mathcal{Y}_2)$.

⁴WalDSPurger '18.

Example: complex phase retrieval

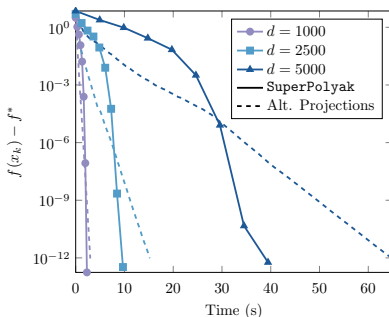
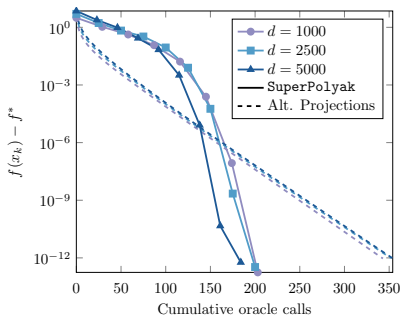
Complex phase retrieval: given $y_{\#} \in \mathbb{R}^m$ with $(y_{\#})_i = |\langle a_i, x_{\#} \rangle|$:

find $\hat{y} \in \mathcal{Y}_1 \cap \mathcal{Y}_2$, $\mathcal{Y}_1 := \{u \in \mathbb{C}^m \mid |u| = y\}$, $\mathcal{Y}_2 := \text{Range}(A)$.

Fallback algorithm: alternating projections⁴.

$$y_{k+1} = T(y_k) := AA^{\dagger} (y_{\#} \odot \text{phase}(y_k))$$

Function used by PolyakBundle: $f(y) := \text{dist}(y, \mathcal{Y}_1) + \text{dist}(y, \mathcal{Y}_2)$.



⁴Waldspurger '18.

Concluding remarks

Not covered in this talk:

- Results for non-isolated \mathcal{X}_* via a uniformization of semismoothness.⁵
- Provable early termination strategies for PolyakBundle loop.
 - In practice, lead to small (even constant-sized) linear systems.

⁵Davis et al. '21

Concluding remarks

Not covered in this talk:

- Results for non-isolated \mathcal{X}_* via a uniformization of semismoothness.⁵
- Provable early termination strategies for PolyakBundle loop.
 - In practice, lead to small (even constant-sized) linear systems.

Open questions:

1. Reduce dim. dependence of local convergence region of PolyakBundle.
2. Remove requirement that $f_* = \min f$ is known.
3. Layer on top of existing large-scale solvers (LPs? QPs?)

⁵Davis et al. '21

Concluding remarks

Not covered in this talk:

- Results for non-isolated \mathcal{X}_* via a uniformization of semismoothness.⁵
- Provable early termination strategies for PolyakBundle loop.
 - In practice, lead to small (even constant-sized) linear systems.

Open questions:

1. Reduce dim. dependence of local convergence region of PolyakBundle.
2. Remove requirement that $f_* = \min f$ is known.
3. Layer on top of existing large-scale solvers (LPs? QPs?)

Thank you!

arXiv:abs/2201.04611

⁵Davis et al. '21

Example: low-rank matrix sensing

Bilinear sensing: recover low-rank factors $U_{\#}, V_{\#}$ from bilinear measurements:

$$y_i = \ell_i^{\top} U_{\#} V_{\#}^{\top} r_i, \quad \ell_i, r_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, I_d).$$

Example: low-rank matrix sensing

Bilinear sensing: recover low-rank factors $U_{\#}, V_{\#}$ from bilinear measurements:

$$y_i = \ell_i^{\top} U_{\#} V_{\#}^{\top} r_i, \quad \ell_i, r_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, I_d).$$

Loss function⁶:

$$f(U, V) = \frac{1}{m} \sum_{i=1}^m |y_i - \ell_i^{\top} UV^{\top} r_i|.$$

Fallback algorithm: Polyak subgradient method.

⁶C. et al., '21

Example: low-rank matrix sensing

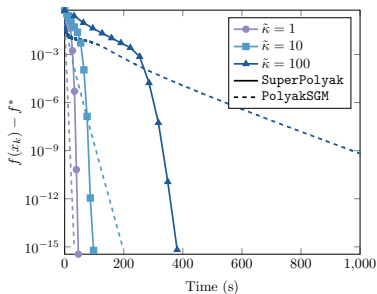
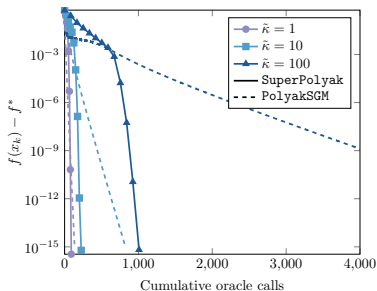
Bilinear sensing: recover low-rank factors $U_{\#}, V_{\#}$ from bilinear measurements:

$$y_i = \ell_i^{\top} U_{\#} V_{\#}^{\top} r_i, \quad \ell_i, r_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, I_d).$$

Loss function⁶:

$$f(U, V) = \frac{1}{m} \sum_{i=1}^m |y_i - \ell_i^{\top} UV^{\top} r_i|.$$

Fallback algorithm: Polyak subgradient method.



⁶C. et al., '21

Example: max-linear regression

Max-linear regression: recover unknown “slopes” $\beta_1^\#, \dots, \beta_r^\#$ from

$$y_i = \max_{j \in [r]} \langle a_i, \beta_j^\# \rangle \quad \text{where} \quad a_i \sim \mathcal{N}(0, I_d).$$

Example: max-linear regression

Max-linear regression: recover unknown “slopes” $\beta_1^\#, \dots, \beta_r^\#$ from

$$y_i = \max_{j \in [r]} \langle a_i, \beta_j^\# \rangle \quad \text{where} \quad a_i \sim \mathcal{N}(0, I_d).$$

Loss function:

$$f(\beta_1, \dots, \beta_r) = \frac{1}{m} \sum_{i=1}^m \left| y_i - \max_{j \in [r]} \langle a_i, \beta_j \rangle \right|.$$

Fallback algorithm: Polyak subgradient method.

Example: max-linear regression

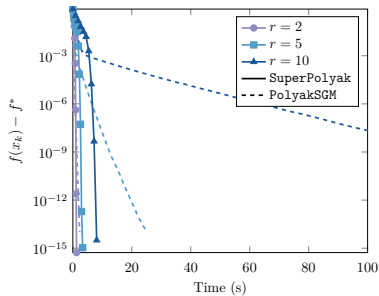
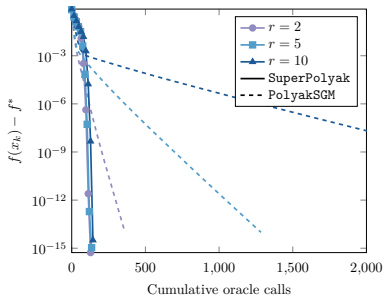
Max-linear regression: recover unknown “slopes” $\beta_1^\#, \dots, \beta_r^\#$ from

$$y_i = \max_{j \in [r]} \langle a_i, \beta_j^\# \rangle \quad \text{where} \quad a_i \sim \mathcal{N}(0, I_d).$$

Loss function:

$$f(\beta_1, \dots, \beta_r) = \frac{1}{m} \sum_{i=1}^m \left| y_i - \max_{j \in [r]} \langle a_i, \beta_j \rangle \right|.$$

Fallback algorithm: Polyak subgradient method.



From projection problem to linear system:

$$\begin{aligned} & \operatorname{argmin} \{ \|x - x_k\|^2 \mid f(y_i) + \langle v_i, x - y_i \rangle = 0 \} \\ &= \operatorname{argmin} \{ \|x - x_k\|^2 \mid \langle v_i, x - x_k \rangle = \langle v_i, y_i - x_k \rangle - f(y_i) \} \\ &= \operatorname{argmin} \{ \|z\|^2 \mid \langle v_i, z \rangle = \langle v_i, y_i - x_k \rangle - f(y_i) \} \\ &= \operatorname{argmin} \{ \|z\|^2 \mid Az + [f(y_i) + \langle v_i, x_k - y_i \rangle]_i = 0 \} \end{aligned}$$

Least-norm solution of $Ax + b = 0$: $x = -A^\dagger b$.