Avoiding saddle points in nonsmooth optimization

Damek Davis School of Operations Research and Information Engineering Cornell University

Joint with L. Jiang (Cornell) and D. Drusvyatskiy (U. Washington)

One World Optimization Seminar Nov 2021

Saddle point avoidance

Recent Realization:

Simple algorithms for minimizing C^2 functions avoid all strict saddle points, when randomly initialized.¹

- Simple algorithms: Gradient descent (GD), coordinate descent....
- Strict saddle points: Critical points that have negative curvature.



Saddle point avoidance

Recent Realization:

Simple algorithms for minimizing C^2 functions avoid all strict saddle points, when randomly initialized.²

- Simple algorithms: Gradient descent (GD), coordinate descent....
- *Strict* saddle points: Critical points that have negative curvature.

Motivation:

For a wealth of estimation and learning problems, all spurious critical points are strict saddles and therefore avoidable!

(Sun-Qu-Wright '15-'18, Ge-Lee-Ma '16, Bhojanapalli-Neyshabur-Srebro '16, Ge-Jin-Zheng '17...)

²Lee-Simchowitz-Jordan-Recht '16

Saddle point avoidance

Recent Realization:

Simple algorithms for minimizing C^2 functions avoid all strict saddle points, when randomly initialized.²

- Simple algorithms: Gradient descent (GD), coordinate descent....
- *Strict* saddle points: Critical points that have negative curvature.

Motivation:

For a wealth of estimation and learning problems, all spurious critical points are strict saddles and therefore avoidable!

(Sun-Qu-Wright '15-'18, Ge-Lee-Ma '16, Bhojanapalli-Neyshabur-Srebro '16, Ge-Jin-Zheng '17...)

This talk:

Do first-order methods avoid "strict saddles" of nonsmooth functions?

²Lee-Simchowitz-Jordan-Recht '16

Weak convexity: an amenable problem class

 $\underset{x \in \mathbb{R}^d}{\text{minimize}} \ F(x)$

Running assumption: weak convexity

$$F(\cdot) + \frac{\rho}{2} \| \cdot \|^2 \qquad \text{is convex}.$$

Weak convexity: an amenable problem class

 $\underset{x \in \mathbb{R}^d}{\text{minimize}} \ F(x)$

Running assumption: weak convexity

$$F(\cdot) + \frac{\rho}{2} \|\cdot\|^2$$
 is convex.

Main example:

 $(convex) \circ (smooth)$

h(c(x))

h is convex and L-Lipschitz; c is smooth with ℓ -Lipschitz Jacobian ($\rho = L\ell$) (Fletcher '80, Powell '83, Burke '85, Wright '90, Lewis-Wright '08, Cartis-Gould-Toint '11,...)

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

 $\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$

Goal: Recover M_{\sharp} from b_i

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

 $\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

 $\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

Natural Nonconvex Penalty Formulation:³

 $\min_{M \in \mathbb{R}^{d \times d}} \ \||\mathcal{A}(M) - b\|| \qquad \text{subject to: } M \text{ is rank} \leq r$

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

$$\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

Natural Nonconvex Penalty Formulation:³

$$M = XX^T \qquad X \in \mathbb{R}^{d \times r}$$

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

$$\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

Natural Nonconvex Penalty Formulation:³

$$\min_{X \in \mathbb{R}^{d \times r}} h(c(X)) := \||\mathcal{A}(XX^{\top}) - b\||$$

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

 $\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

Natural Nonconvex Penalty Formulation:³

$$\min_{X \in \mathbb{R}^{d \times r}} h(c(X)) := \| |\mathcal{A}(XX^{\top}) - b \| |$$

Question: Is there a natural norm $\| \cdot \|$ that enables recovery?

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

$$\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

Natural Nonconvex Penalty Formulation:³

$$\min_{X \in \mathbb{R}^{d \times r}} h(c(X)) := \||\mathcal{A}(XX^{\top}) - b\||$$

Question: Is there a natural norm $\|\cdot\|$ that enables recovery?

Typical norms⁴: $\||\cdot||| = \frac{1}{\sqrt{m}} \|\cdot\|_2$ and $\||\cdot||| = \frac{1}{m} \|\cdot\|_1$

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

$$\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

Natural Nonconvex Penalty Formulation:³

$$\min_{X \in \mathbb{R}^{d \times r}} h(c(X)) := \||\mathcal{A}(XX^{\top}) - b\||$$

Question: Is there a natural norm $\|\cdot\|$ that enables recovery?

Typical norms⁴: $\||\cdot|| = \frac{1}{\sqrt{m}} \|\cdot\|_2$ and $\||\cdot|| = \frac{1}{m} \|\cdot\|_1$

• ℓ_2 : Gaussian A_i /Gaussian noise, leads to smooth problems.

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Set-up: Fix rank r matrix $M_{\sharp} \succeq 0$ and observe measurements

$$\langle A_i, M_{\sharp} \rangle \approx b_i \qquad \forall i = 1, \dots, m.$$

Goal: Recover M_{\sharp} from b_i

Examples: Matrix completion, robust PCA, phase retrieval...

Natural Nonconvex Penalty Formulation:³

$$\min_{X \in \mathbb{R}^{d \times r}} h(c(X)) := \||\mathcal{A}(XX^{\top}) - b\||$$

Question: Is there a natural norm $\|\cdot\|$ that enables recovery?

Typical norms⁴: $\||\cdot|| = \frac{1}{\sqrt{m}} \|\cdot\|_2$ and $\||\cdot|| = \frac{1}{m} \|\cdot\|_1$

- ℓ_2 : Gaussian A_i /Gaussian noise, leads to smooth problems.
- ℓ_1 : structured A_i /sparse corruption, leads to nonsmooth problems.

³Burer-Monteiro '01

⁴Candes-Tao '05, Chen-Chi-Goldsmith '13

Common iterative methods take form

 $x_{t+1} = \operatorname*{arg\,min}_{y} F_{x_t}(y)$

where F_{x_t} = nonsmooth strongly convex model of F.

Common iterative methods take form

 $x_{t+1} = \operatorname*{arg\,min}_{y} F_{x_t}(y)$

where F_{x_t} = nonsmooth strongly convex model of F.

Example: Proximal point



Common iterative methods take form

 $x_{t+1} = \operatorname*{arg\,min}_{y} F_{x_t}(y)$

where F_{x_t} = nonsmooth strongly convex model of F.

Example: Proximal linear (for $F = h \circ c$)



 $F_{x_t}(y) = h(c(x_t) + \nabla c(x_t)(y - x_t)) + \frac{1}{2\eta} \|y - x_t\|^2$

Common iterative methods take form

 $x_{t+1} = \operatorname*{arg\,min}_{y} F_{x_t}(y)$

where F_{x_t} = nonsmooth strongly convex model of F.

Example:

Algorithm	Objective F	Update function $F_x(y)$
Prox-point	F(x)	$F(y) + rac{1}{2\eta} \ y - x\ ^2$
Prox-linear	h(c(x)) + r(x)	$h(c(x) + \nabla c(x)(y-x)) + r(y) + \frac{1}{2\eta} y-x ^2$
Prox-gradient	f(x) + r(x)	$\left \begin{array}{c} f(x) + \langle abla f(x), y - x angle + rac{r(y)}{2\eta} \left\ y - x ight\ ^2 \end{array} ight ^2$

Table: h is convex and Lipschitz, r is weakly convex, and f and c are C^2 -smooth.

 $\ensuremath{\mathbf{Q}}\xspace$: What is an avoidable saddle point in nonsmooth optimization?⁵

⁵(D-Drusvyatskiy '19)

Recall C^2 case: A strict saddle is critical point with negative curvature:

$$\nabla F(x) = 0$$
 and $\lambda_{\min}(\nabla^2 F(x)) < 0$

⁵(D-Drusvyatskiy '19)

Recall C^2 case: A strict saddle is critical point with negative curvature:

$$\nabla F(x) = 0$$
 and $\lambda_{\min}(\nabla^2 F(x)) < 0$

Generalization Attempt: A strict saddle is critical point such that

⁵(D-Drusvyatskiy '19)

Recall C^2 case: A strict saddle is critical point with negative curvature:

$$\nabla F(x) = 0$$
 and $\lambda_{\min}(\nabla^2 F(x)) < 0$

Generalization Attempt: A strict saddle is critical point such that

• There exists direction v s.t.

$$g(t) := F(x + tv) \text{ is } C^2.$$

⁵(D-Drusvyatskiy '19)

Recall C^2 case: A strict saddle is critical point with negative curvature:

$$\nabla F(x) = 0$$
 and $\lambda_{\min}(\nabla^2 F(x)) < 0$

Generalization Attempt: A strict saddle is critical point such that

• There exists direction v s.t.

$$g(t) := F(x + tv) \text{ is } C^2.$$

• Function g has negative curvature:

g''(0) < 0.

⁵(D-Drusvyatskiy '19)

Recall C^2 case: A strict saddle is critical point with negative curvature:

$$\nabla F(x) = 0$$
 and $\lambda_{\min}(\nabla^2 F(x)) < 0$

Generalization Attempt: A strict saddle is critical point such that

• There exists direction v s.t.

$$g(t) := F(x + tv) \text{ is } C^2.$$

• Function g has negative curvature:

g''(0) < 0.

Equivalent when F is C^2 .

⁵(D-Drusvyatskiy '19)

Negative curvature is not enough even for C^1 functions



Negative curvature: $F(0,y) = -\alpha y^2$

Negative curvature is not enough even for C^1 functions



Negative curvature: $F(0,y) = -\alpha y^2$

Problem: do not reach *y* axis fast enough to benefit from curvature!

An extra ingredient: sharpness

Idea: Require F to grow **sharply** away from axis:

 $\inf\{\|\nabla F(x,y)\|\colon \text{ for } (x,y) \text{ off of } y \text{ axis}\}>0$

Benefit: Ensures grad. flow aims towards axis with (at least) constant speed.

An extra ingredient: sharpness

Idea: Require F to grow **sharply** away from axis:

 $\inf\{\|\nabla F(x,y)\|: \text{ for } (x,y) \text{ off of } y \text{ axis}\} > 0$

Benefit: Ensures grad. flow aims towards axis with (at least) constant speed.



Negative curvature: $F(0, y) = -\alpha y^2$

An extra ingredient: sharpness

Idea: Require F to grow **sharply** away from axis:

 $\inf\{\|\nabla F(x,y)\|\colon \text{ for } (x,y) \text{ off of } y \text{ axis}\} > 0$

Benefit: Ensures grad. flow aims towards axis with (at least) constant speed.



Negative curvature: $F(0, y) = -\alpha y^2$

Question: How to generalize?

Idea: Replace axis with "active manifold" of smoothness.

Idea: Replace axis with "active manifold" of smoothness.

Defn: Critical point lies on C^2 -smooth "active manifold \mathcal{M} ":

1. F varies C^2 -smoothly along \mathcal{M} .

Idea: Replace axis with "active manifold" of smoothness.

Defn: Critical point lies on C^2 -smooth "active manifold \mathcal{M} ":

- 1. F varies C^2 -smoothly along \mathcal{M} .
- 2. F grows sharply normal to \mathcal{M} :

 $\inf\{\|v\|: v \in \partial F(z): z \in U \setminus \mathcal{M}\} > 0.$

Idea: Replace axis with "active manifold" of smoothness.

Defn: Critical point lies on C^2 -smooth "active manifold \mathcal{M} ":

- 1. F varies C^2 -smoothly along \mathcal{M} .
- 2. F grows sharply normal to \mathcal{M} :

$$\inf\{\|v\|: v \in \partial F(z): z \in U \setminus \mathcal{M}\} > 0.$$



Idea: Replace axis with "active manifold" of smoothness.

Defn: Critical point lies on C^2 -smooth "active manifold \mathcal{M} ":

- 1. F varies C^2 -smoothly along \mathcal{M} .
- 2. F grows sharply normal to \mathcal{M} :

$$\inf\{\|v\|: v \in \partial F(z): z \in U \setminus \mathcal{M}\} > 0.$$



Question: What about curvature? (Wright '93, Lemaréchal-Oustry-Sagastizábal '96, Bonnans-Shapiro '00, Lewis '03, Drusvyatskiy-Lewis '14...)

Putting it all together: the active strict saddle property



(a) A nonsmooth loss F


Defn: (D-Drusvyatskiy '19) a critical point \bar{x} of F is an active strict saddle if



Defn: (D-Drusvyatskiy '19) a critical point \bar{x} of F is an active strict saddle if

1. F admits active manifold \mathcal{M} containing \bar{x} .



Defn: (D-Drusvyatskiy '19) a critical point \bar{x} of F is an active strict saddle if

- 1. F admits active manifold \mathcal{M} containing \bar{x} .
- 2. The smooth extension $F \circ P_{\mathcal{M}}$ has a strict saddle point at \bar{x} :

$$\lambda_{\min}(\nabla^2(\boldsymbol{F} \circ \boldsymbol{P}_{\mathcal{M}})(\bar{x})) < 0.$$



Although it may seem stringent, this property is generic:

Theorem (Drusvyatskiy-loffe-Lewis '16, D-Drusvyatskiy '19)

If F is semi-algebraic and weakly convex, then for full Lebesgue measure set of perturbations $v \in \mathbb{R}^d$ every critical point of

$$F_v(x) = F(x) - \langle v, x \rangle$$

is either an active strict saddle or a local minimizer.

Although it may seem stringent, this property is generic:

Theorem (Drusvyatskiy-loffe-Lewis '16, D-Drusvyatskiy '19) If F is semi-algebraic and weakly convex, then for full Lebesgue measure set of perturbations $v \in \mathbb{R}^d$ every critical point of

$$F_v(x) = F(x) - \langle v, x \rangle$$

is either an active strict saddle or a local minimizer.



Example is Highly Unstable: small linear tilts do not exhibit this behavior!

Question: Do the three proximal methods avoid active strict saddles?

⁶For the algorithms considered thus far, critical points are fixed points of the iteration.

Question: Do the three proximal methods avoid active strict saddles?

Strategy: Borrow "stable manifold theorem" argument from smooth setting!

⁶For the algorithms considered thus far, critical points are fixed points of the iteration.

Question: Do the three proximal methods avoid active strict saddles?

Strategy: Borrow "stable manifold theorem" argument from smooth setting!

Key: view algorithms

 $x_{t+1} = \operatorname*{arg\,min}_{y} F_{x_t}(y),$

as fixed-point iteration of well-behaved operator $T.^6$

⁶For the algorithms considered thus far, critical points are fixed points of the iteration.

Fixed point iteration

 $x_{t+1} = T(x_t)$

Fixed point iteration

 $x_{t+1} = T(x_t)$ [Grad descent is $T = I - \eta \nabla F$]

Fixed point iteration

 $x_{t+1} = T(x_t)$ [Grad descent is $T = I - \eta \nabla F$]

Recipe:

• Strict saddles \bar{x}

$$abla F(\bar{x}) = 0$$
 and $\lambda_{\min}(\nabla^2 F(\bar{x})) < 0$

are unstable fixed points:

 $\nabla T(\bar{x})$ has EigVal of magnitude $\,>1$

Fixed point iteration

 $x_{t+1} = T(x_t)$ [Grad descent is $T = I - \eta \nabla F$]

Recipe:

- Strict saddles \bar{x}

 $abla F(\bar{x}) = 0$ and $\lambda_{\min}(\nabla^2 F(\bar{x})) < 0$

are unstable fixed points:

 $abla T(ar{x})$ has EigVal of magnitude > 1

Classical center-stable manifold theorem implies

 $W:=\left\{x\colon \lim_{k\to\infty}T^k(x) \text{ is unstable }\right\} \quad \text{ has Lebesgue measure zero.}$

Fixed point iteration

 $x_{t+1} = T(x_t)$ [Grad descent is $T = I - \eta \nabla F$]

Recipe:

• Strict saddles \bar{x}

 $abla F(ar{x}) = 0$ and $\lambda_{\min}(
abla^2 F(ar{x})) < 0$

are unstable fixed points:

 $abla T(ar{x})$ has EigVal of magnitude > 1

Classical center-stable manifold theorem implies

 $W:=\left\{x\colon \lim_{k\to\infty}T^k(x) \text{ is unstable }\right\} \quad \text{ has Lebesgue measure zero.}$

- Since random init will not land in W, algorithm avoids strict saddles

Fixed point iteration

 $x_{t+1} = T(x_t)$ [Grad descent is $T = I - \eta \nabla F$]

Recipe:

• Strict saddles \bar{x}

 $\nabla F(\bar{x})=0 \qquad \text{and} \qquad \lambda_{\min}(\nabla^2 F(\bar{x}))<0$ are unstable fixed points:

 $\nabla T(\bar{x})$ has EigVal of magnitude > 1

Classical center-stable manifold theorem implies

 $W:=\left\{x\colon \lim_{k\to\infty}T^k(x) \text{ is unstable }\right\} \quad \text{ has Lebesgue measure zero.}$

• Since random init will not land in W, algorithm avoids strict saddles

Important: Argument requires that *T* is local diffeomorphism.

Beyond gradient descent

To apply argument, need

1. Local Smoothness: The update mapping

$$S(x) = \operatorname*{arg\,min}_{y} F_x(y),$$

is a local ${\cal C}^1$ diffeomorphism near active strict saddle points.

Beyond gradient descent

To apply argument, need

1. Local Smoothness: The update mapping

$$S(x) = \operatorname*{arg\,min}_{y} F_x(y),$$

is a local C^1 diffeomorphism near active strict saddle points.

2. Unstable: Active strict saddle points \bar{x} are unstable:

 $\nabla S(\bar{x})$ has EigVal of magnitude > 1.

Beyond gradient descent

To apply argument, need

1. Local Smoothness: The update mapping

$$S(x) = \operatorname*{arg\,min}_{y} F_x(y),$$

is a local ${\cal C}^1$ diffeomorphism near active strict saddle points.

2. Unstable: Active strict saddle points \bar{x} are unstable:

 $\nabla S(\bar{x})$ has EigVal of magnitude > 1.

Focus on Local Smoothness, since other calculation complex.

Surprising: Function F is nonsmooth, yet S is C^1 around strict saddles. Why?

⁷Lemaréchal-Sagastizábal '97

Surprising: Function F is nonsmooth, yet S is C^1 around strict saddles. Why?

Sharpness \implies Identification $S(x) \in \mathcal{M}$ near $\bar{x}!$

Example: Prox-point



⁷Lemaréchal-Sagastizábal '97

Surprising: Function F is nonsmooth, yet S is C^1 around strict saddles. Why?

Sharpness \implies Identification $S(x) \in \mathcal{M}$ near $\bar{x}!$

Important: Do not need to know \mathcal{M} !

⁷Lemaréchal-Sagastizábal '97

Surprising: Function F is nonsmooth, yet S is C^1 around strict saddles. Why?

Sharpness \implies Identification $S(x) \in \mathcal{M}$ near $\bar{x}!$

Important: Do not need to know \mathcal{M} !

Consequence (Prox-point Method):

$$S(x) = \underset{y}{\arg\min} F(y) + \frac{1}{2\eta} \|y - x\|^{2} = \underset{y \in \mathcal{M}}{\arg\min} F(y) + \frac{1}{2\eta} \|y - x\|^{2}.$$

⁷Lemaréchal-Sagastizábal '97

Surprising: Function F is nonsmooth, yet S is C^1 around strict saddles. Why?

Sharpness \implies Identification $S(x) \in \mathcal{M}$ near $\bar{x}!$

Important: Do not need to know \mathcal{M} !

Consequence (Prox-point Method):

$$S(x) = \underset{y}{\arg\min} F(y) + \frac{1}{2\eta} \|y - x\|^2 = \underset{y \in \mathcal{M}}{\arg\min} F(y) + \frac{1}{2\eta} \|y - x\|^2.$$

 \implies minimizing smooth function over smooth manifold!

⁷Lemaréchal-Sagastizábal '97

Surprising: Function F is nonsmooth, yet S is C^1 around strict saddles. Why?

Sharpness \implies Identification $S(x) \in \mathcal{M}$ near $\bar{x}!$

Important: Do not need to know \mathcal{M} !

Consequence (Prox-point Method):

$$S(x) = \underset{y}{\arg\min} F(y) + \frac{1}{2\eta} \|y - x\|^2 = \underset{y \in \mathcal{M}}{\arg\min} F(y) + \frac{1}{2\eta} \|y - x\|^2.$$

 \implies minimizing smooth function over smooth manifold!

Then Weak convexity + classical perturbation theory $\implies S$ is C^1 near \bar{x} .⁷

⁷Lemaréchal-Sagastizábal '97

Proof extends to the three methods:

Algorithm	Objective F	Update function $F_x(y)$
Prox-point	F(x)	$F(y)+rac{1}{2\eta}\ y-x\ ^2$
Prox-linear	h(c(x)) + r(x)	$h(c(x) + \nabla c(x)(y - x)) + r(y) + \frac{1}{2\eta} y - x ^2$
Prox-gradient	f(x) + r(x)	$ f(x)+\langle abla f(x),y-x angle+r(y)+rac{1}{2\eta}\ y-x\ ^2$

Table: h is convex and Lipschitz, r is weakly convex, and f and c are C^2 -smooth.

Theorem: (Local smoothness, D-Drusvyatskiy '19)

Around each active strict saddle \bar{x} of F, the iteration mapping

 $S(x) = \operatorname*{arg\,min}_{y} F_x(y),$

is C^1 and the Jacobian $\nabla S(\bar{x})$ has a real EigVal strictly greater than 1

Proof more interesting/surprising for prox-gradient and prox-linear.

Problem: S may not be Local diffeomorphism

Problem: S may not be Local diffeomorphism

Easy solution: Add damping

 $T = (1 - \lambda)I + \lambda S.$

Corollary: (Random initialization, D-Drusvyatskiy '19)

Randomly initialized three methods with small damping

$$x_{t+1} = (1 - \lambda)x_t + \lambda S(x_t),$$

locally escape active strict saddles.

Globalization:

- Results hold globally when S is Lipschitz (prox-point, prox-gradient)
- Open Problem: Is prox-linear update globally Lipschitz?

Limitation of result: Only applies to three "proximal methods."

Algorithm	Objective F	Update function $F_x(y)$
Prox-point	F(x)	$F(y) + rac{1}{2\eta} \ y - x\ ^2$
Prox-linear	h(c(x)) + r(x)	$h(c(x) + \nabla c(x)(y - x)) + r(y) + \frac{1}{2\eta} y - x ^2$
Prox-gradient	f(x) + r(x)	$ f(x)+\langle abla f(x),y-x angle+r(y)+rac{1}{2\eta}\ y-x\ ^2$

Table: h is convex and Lipschitz, r is weakly convex, and f and c are C^2 -smooth.

Limitation of result: Only applies to three "proximal methods."

Algorithm	Objective F	Update function $F_x(y)$
Prox-point	F(x)	$F(y)+rac{1}{2\eta}\ y-x\ ^2$
Prox-linear	h(c(x)) + r(x)	$h(c(x) + \nabla c(x)(y - x)) + r(y) + \frac{1}{2\eta} y - x ^2$
Prox-gradient	f(x) + r(x)	$f(x)+\langle abla f(x),y-x angle+oldsymbol{r}(y)+rac{1}{2\eta}ig\ y-x\ ^2$

Table: h is convex and Lipschitz, r is weakly convex, and f and c are C^2 -smooth.

Drawbacks:

1. Numerical Difficulties: need exact solutions to subproblems.

Limitation of result: Only applies to three "proximal methods."

Algorithm	Objective F	Update function $F_x(y)$
Prox-point	F(x)	$F(y)+rac{1}{2\eta}\ y-x\ ^2$
Prox-linear	h(c(x)) + r(x)	$h(c(x) + \nabla c(x)(y - x)) + r(y) + \frac{1}{2\eta} y - x ^2$
Prox-gradient	f(x) + r(x)	$f(x)+\langle abla f(x),y-x angle+oldsymbol{r}(y)+rac{1}{2\eta}ig\ y-x\ ^2$

Table: h is convex and Lipschitz, r is weakly convex, and f and c are C^2 -smooth.

Drawbacks:

- 1. Numerical Difficulties: need exact solutions to subproblems.
- 2. Decomposable structure not always available.

Limitation of result: Only applies to three "proximal methods."

Algorithm	Objective F	Update function $F_x(y)$
Prox-point	F(x)	$F(y)+rac{1}{2\eta}\ y-x\ ^2$
Prox-linear	h(c(x)) + r(x)	$h(c(x) + \nabla c(x)(y - x)) + r(y) + \frac{1}{2\eta} y - x ^2$
Prox-gradient	f(x) + r(x)	$f(x)+\langle abla f(x),y-x angle+oldsymbol{r}(y)+rac{1}{2\eta}ig\ y-x\ ^2$

Table: h is convex and Lipschitz, r is weakly convex, and f and c are C^2 -smooth.

Drawbacks:

- 1. Numerical Difficulties: need exact solutions to subproblems.
- 2. Decomposable structure not always available.

Alternative: subgradient method

The subdifferential of a weakly convex function

Fact: For any $F \colon \mathbb{R}^d \to \mathbb{R}$, have equivalence:

- F is ρ -weakly convex
- Subgradient inequality: $\forall x \exists v_x$ satisfying

 $F(y) \ge F(x) + \langle v_x, y - x \rangle - \frac{\rho}{2} \|y - x\|^2$



The subdifferential of a weakly convex function

Fact: For any $F \colon \mathbb{R}^d \to \mathbb{R}$, have equivalence:

- F is ρ -weakly convex
- Subgradient inequality: $\forall x \exists v_x$ satisfying

$$F(y) \ge F(x) + \langle v_x, y - x \rangle - \frac{\rho}{2} \|y - x\|^2$$

Subdifferential: $\partial F(x) := \{v_x\}$



The subdifferential of a weakly convex function

Fact: For any $F \colon \mathbb{R}^d \to \mathbb{R}$, have equivalence:

- F is ρ -weakly convex
- Subgradient inequality: $\forall x \exists v_x$ satisfying

$$F(y) \ge F(x) + \langle v_x, y - x \rangle - \frac{\rho}{2} \|y - x\|^2$$



Subdifferential:

$$\partial F(x) := \{v_x\}$$

Calculus:

$$\partial(h \circ c)(x) := \nabla c(x)^T \partial h(c(x))$$
The subdifferential of a weakly convex function

Fact: For any $F \colon \mathbb{R}^d \to \mathbb{R}$, have equivalence:

- F is ρ -weakly convex
- Subgradient inequality: $\forall x \exists v_x$ satisfying

$$F(y) \ge F(x) + \langle v_x, y - x \rangle - \frac{\rho}{2} \|y - x\|^2$$



$$\partial F(x) := \{v_x\}$$

Calculus:

$$\partial(h \circ c)(x) := \nabla c(x)^T \partial h(c(x))$$

Fermat's rule: If \bar{x} is a local minimizer of F then

 $0 \in \partial F(\bar{x}).$



Idea: At time t

⁸Bolte-Pauwels '19-'20

Idea: At time t

1. "Linearize F:" choose $v_t \in \partial F(x_t)$ and form

$$F_{x_t,\alpha_t}(y) = F(x_t) + \langle \mathbf{v}_t, y - x_t \rangle + \frac{1}{2\alpha_t} \|y - x_t\|^2.$$

⁸Bolte-Pauwels '19-'20

Idea: At time t

1. "Linearize F:" choose $v_t \in \partial F(x_t)$ and form

$$F_{x_t,\alpha_t}(y) = F(x_t) + \langle \mathbf{v}_t, y - x_t \rangle + \frac{1}{2\alpha_t} \|y - x_t\|^2.$$

2. Next iterate minimizes:

$$\begin{aligned} x_{t+1} &= \operatorname*{arg\,min}_{y} F_{x_t,\alpha_t}(y) \\ &= x_t - \alpha_t v_t. \end{aligned}$$

⁸Bolte-Pauwels '19-'20

Idea: At time t

1. "Linearize F:" choose $v_t \in \partial F(x_t)$ and form

$$F_{x_t,\alpha_t}(y) = F(x_t) + \langle \mathbf{v}_t, y - x_t \rangle + \frac{1}{2\alpha_t} \|y - x_t\|^2.$$

2. Next iterate minimizes:



Idea: At time t

1. "Linearize F:" choose $v_t \in \partial F(x_t)$ and form

$$F_{x_t,\alpha_t}(y) = F(x_t) + \langle v_t, y - x_t \rangle + \frac{1}{2\alpha_t} \|y - x_t\|^2.$$

2. Next iterate minimizes:

$$\begin{aligned} x_{t+1} &= \operatorname*{arg\,min}_{y} F_{x_t,\alpha_t}(y) \\ &= x_t - \alpha_t v_t. \end{aligned}$$

Benefits:

1. Computable with extensive calculus: $\partial(h \circ c)(x) := \nabla c(x)^T \partial h(c(x))$

⁸Bolte-Pauwels '19-'20

Idea: At time t

1. "Linearize F:" choose $v_t \in \partial F(x_t)$ and form

$$F_{x_t,\alpha_t}(y) = F(x_t) + \langle v_t, y - x_t \rangle + \frac{1}{2\alpha_t} \|y - x_t\|^2.$$

2. Next iterate minimizes:

$$\begin{aligned} x_{t+1} &= \operatorname*{arg\,min}_{y} F_{x_t,\alpha_t}(y) \\ &= x_t - \alpha_t v_t. \end{aligned}$$

Benefits:

- 1. Computable with extensive calculus: $\partial(h \circ c)(x) := \nabla c(x)^T \partial h(c(x))$
- 2. Can often replace v_t with result of auto-differentiation procedure.⁸

⁸Bolte-Pauwels '19-'20

Extension: Subgradient method

Question: Does subgradient method avoid active strict saddle points?

 $x_{t+1} \in x_t - \alpha_t \partial F(x_t)$

⁹D-Drusvyatskiy-Jiang '21

Extension: Subgradient method

Question: Does subgradient method avoid active strict saddle points?

```
x_{t+1} \in x_t - \alpha_t \partial F(x_t)
```

Difficulties:

- Identification fails: $x_t \notin \mathcal{M}$.
- Unclear how to leverage smoothness on the manifold.

Our recent work 9 overcomes these difficulties.

⁹D-Drusvyatskiy-Jiang '21

Extension: Subgradient method

Question: Does subgradient method avoid active strict saddle points?

```
x_{t+1} \in x_t - \alpha_t \partial F(x_t)
```

Difficulties:

- Identification fails: $x_t \notin \mathcal{M}$.
- Unclear how to leverage smoothness on the manifold.

Our recent work 9 overcomes these difficulties.

Key: "orthogonal decomposition" of trajectory.

⁹D-Drusvyatskiy-Jiang '21



¹⁰Mifflin-Sagastizábal '05



Decompose trajectory:

¹⁰Mifflin-Sagastizábal '05



Decompose trajectory:

1. Tangent directions:

$$P_{\mathcal{M}}(x_{t+1}) \approx P_{\mathcal{M}}(x_t) - \alpha_t \nabla F_{\mathcal{U}}(x_t)$$

¹⁰Mifflin-Sagastizábal '05



Decompose trajectory:

1. Tangent directions:

$$P_{\mathcal{M}}(x_{t+1}) \approx P_{\mathcal{M}}(x_t) - \alpha_t \nabla F_{\mathcal{U}}(x_t)$$

2. Normal directions:

$$x_{t+1} - P_{\mathcal{M}}(x_{t+1}) \approx x_t - P_{\mathcal{M}}(x_t) - \alpha_t \widetilde{\nabla} F_{\mathcal{V}}(x_t)$$

¹⁰Mifflin-Sagastizábal '05

The two regularity assumptions

1. Aiming: Negative subgradients aim towards manifold:

Sharpness
$$\Longrightarrow \langle \widetilde{\nabla} F_{\mathcal{V}}(x_t), x_t - P_{\mathcal{M}}(x_t) \rangle \ge \mu \operatorname{dist}(x_t, \mathcal{M})$$



The two regularity assumptions

1. Aiming: Negative subgradients aim towards manifold:

Sharpness
$$\Longrightarrow \langle \widetilde{\nabla} F_{\mathcal{V}}(x_t), x_t - P_{\mathcal{M}}(x_t) \rangle \geq \mu \operatorname{dist}(x_t, \mathcal{M})$$

2. Smooth in tangent directions:

$$\|P_{\mathcal{T}_{\mathcal{M}}(y)}\tilde{\nabla}F_{\mathcal{V}}(x_t)\| \leq C\|x_t - y\|$$
 for $y \in \mathcal{M}$.



The two regularity assumptions

1. Aiming: Negative subgradients aim towards manifold:

Sharpness
$$\implies \langle \widetilde{\nabla} F_{\mathcal{V}}(x_t), x_t - P_{\mathcal{M}}(x_t) \rangle \ge \mu \operatorname{dist}(x_t, \mathcal{M})$$

2. Smooth in tangent directions:

$$\|P_{\mathcal{T}_{\mathcal{M}}(y)}\tilde{\nabla}F_{\mathcal{V}}(x_t)\| \leq C\|x_t - y\|$$
 for $y \in \mathcal{M}$.



Prevalent: true generically for weakly convex semialgebraic problems.

The two pillars

The two pillars: For a wide class of problems

• Subgradient method quickly approaches the active manifold:

 $\operatorname{dist}(x_t, \mathcal{M}) = O(\alpha_t).$



(a) Quickly approach manifold

The two pillars

The two pillars: For a wide class of problems

• Subgradient method quickly approaches the active manifold:

 $\operatorname{dist}(x_t, \mathcal{M}) = O(\alpha_t).$

• The shadow $y_t = P_{\mathcal{M}}(x_t)$ forms inexact Riemannian gradient sequence:

$$\mathbf{y_{t+1}} = \mathbf{y_t} - \alpha_t \nabla_{\mathcal{M}} F(\mathbf{y_t}) + O(\alpha_t \operatorname{dist}(x_t, \mathcal{M}) + \alpha_t^2).$$



(b) "Smooth in tangent directions"

(a) Quickly approach manifold

The two pillars

The two pillars: For a wide class of problems

• Subgradient method quickly approaches the active manifold:

 $\operatorname{dist}(x_t, \mathcal{M}) = O(\alpha_t).$

• The shadow $y_t = P_{\mathcal{M}}(x_t)$ forms inexact Riemannian gradient sequence:

$$y_{t+1} = y_t - \alpha_t \nabla_{\mathcal{M}} F(y_t) + O(\alpha_t \operatorname{dist}(x_t, \mathcal{M}) + \alpha_t^2).$$



Conclusion: Get to the manifold quick enough to leverage smoothness of F!

Due to inexactness, must analyze "perturbed" subgradient method¹¹:

 $x_{t+1} \in x_t - \alpha_t (\partial F(x_t) + \nu_t)$ where $\nu_t \sim \text{Unif}(B)$.

¹¹D-Drusvyatskiy-Jiang '21

¹²Concurrent work: Bianchi-Hachem-Schechtman'21.

Due to inexactness, must analyze "perturbed" subgradient method¹¹:

 $x_{t+1} \in x_t - \alpha_t (\partial F(x_t) + \nu_t)$ where $\nu_t \sim \text{Unif}(B)$.

 $\implies y_{t+1} = y_t - \alpha_t (\nabla_{\mathcal{M}} F(y_t) + \nu_t) + O(\alpha_t \operatorname{dist}(x_t, \mathcal{M}) + \alpha_t^2).$

¹¹D-Drusvyatskiy-Jiang '21

¹²Concurrent work: Bianchi-Hachem-Schechtman'21.

Due to inexactness, must analyze "perturbed" subgradient method¹¹:

 $x_{t+1} \in x_t - \alpha_t (\partial F(x_t) + \nu_t)$ where $\nu_t \sim \text{Unif}(B)$.

Under mild conditions, we show

Theorem: (D-Drusvyatskiy-Jiang '19)¹²

Almost surely, x_t does not converge to an active strict saddle point.

¹¹D-Drusvyatskiy-Jiang '21

¹²Concurrent work: Bianchi-Hachem-Schechtman'21.

Due to inexactness, must analyze "perturbed" subgradient method¹¹:

 $x_{t+1} \in x_t - \alpha_t (\partial F(x_t) + \nu_t)$ where $\nu_t \sim \text{Unif}(B)$.

Under mild conditions, we show

Theorem: (D-Drusvyatskiy-Jiang '19)¹²

Almost surely, x_t does not converge to an active strict saddle point.

Corollary: (D-Drusvyatskiy-Jiang '19)

Perturbed subgradient method converges only to local minimizers of generic semialgebraic weakly convex functions.

¹¹D-Drusvyatskiy-Jiang '21

¹²Concurrent work: Bianchi-Hachem-Schechtman'21.

Due to inexactness, must analyze "perturbed" subgradient method¹¹:

 $x_{t+1} \in x_t - \alpha_t (\partial F(x_t) + \nu_t)$ where $\nu_t \sim \text{Unif}(B)$.

Under mild conditions, we show

Theorem: (D-Drusvyatskiy-Jiang '19)¹²

Almost surely, x_t does not converge to an active strict saddle point.

Corollary: (D-Drusvyatskiy-Jiang '19)

Perturbed subgradient method converges only to local minimizers of generic semialgebraic weakly convex functions.

Extensions.

1. Algorithms: Proximal/projected subgradient methods.

¹¹D-Drusvyatskiy-Jiang '21

¹²Concurrent work: Bianchi-Hachem-Schechtman'21.

Due to inexactness, must analyze "perturbed" subgradient method¹¹:

 $x_{t+1} \in x_t - \alpha_t (\partial F(x_t) + \nu_t)$ where $\nu_t \sim \text{Unif}(B)$.

Under mild conditions, we show

Theorem: (D-Drusvyatskiy-Jiang '19)¹²

Almost surely, x_t does not converge to an active strict saddle point.

Corollary: (D-Drusvyatskiy-Jiang '19)

Perturbed subgradient method converges only to local minimizers of generic semialgebraic weakly convex functions.

Extensions.

- 1. Algorithms: Proximal/projected subgradient methods.
- 2. Beyond weak convexity: Clarke regularity.

¹¹D-Drusvyatskiy-Jiang '21

¹²Concurrent work: Bianchi-Hachem-Schechtman'21.

Thank you!

References

- Proximal methods avoid active strict saddles of weakly convex functions
 D, Drusvyatskiy. Found. Comput. Math. arxiv.org/abs/1912.07146.
- Subgradient methods near active manifolds: saddle point avoidance, local convergence, and asymptotic normality

D, Drusvyatskiy, Jiang. https://arxiv.org/abs/2108.11832.