

Projection Efficient Subgradient Method and Optimal Nonsmooth Frank-Wolfe Method

Praneeth Netrapalli

Microsoft Research India



Kiran Thekumparampil

UIUC



Prateek Jain

Google India Research



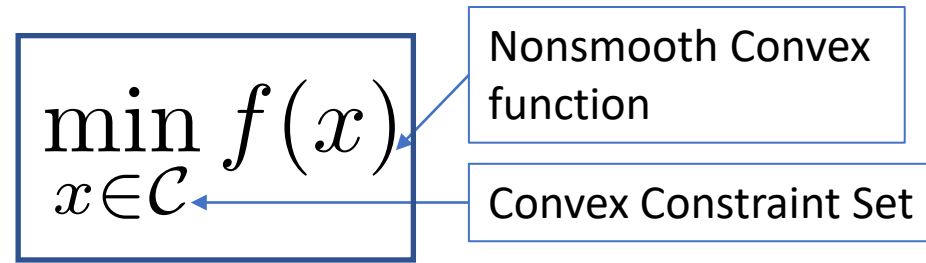
Sewoong Oh

Univ. of Washington

Outline

- Constrained Nonsmooth Convex Problem
- What is currently known
- Our results
- Techniques: Moreau-Yosida smoothing and Gradient sliding

Constrained Nonsmooth Convex Minimization



$$\mathcal{C} \subset \mathbb{R}^d$$

GOAL:

Find an ϵ -suboptimal *feasible* solution \tilde{x}

$$f(\tilde{x}) - \min_{x \in \mathcal{C}} f(x) \leq \epsilon \quad \text{and} \quad \tilde{x} \in \mathcal{C}$$

- High dimensional: $d \gg \text{poly}(\epsilon^{-1})$
- Prefer methods with convergence rate independent of d

Assumptions

$$\min_{x \in \mathcal{C}} f(x)$$

• G-Lipschitz continuous function $|f(x) - f(y)| \leq G\|x - y\|$

• Bounded Constraint Set $D = \max_{x, y \in \mathcal{C}} \|x - y\|$

$$\min_{x \in \mathcal{C}} f(x)$$

Assumptions

- G-Lipschitz continuous function $|f(x) - f(y)| < G\|x - y\|$

- Bounded Constraint

For $f(x) = |x|$

$$\partial f(x) = \begin{cases} -1 & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Black-box Oracle

Oracle Complexity
= Num. of Oracle calls

- First-order Oracle (FO) $\text{FO}(x) \in \partial f(x)$ sub-gradient

- Projection Oracle (PO) $\text{PO}(x) = \mathcal{P}_{\mathcal{C}}(x) = \arg \min_{y \in \mathcal{C}} \|x - y\|^2$

Black-box oracles. Cannot use methods like IPM, Augmented Lagrangian Methods

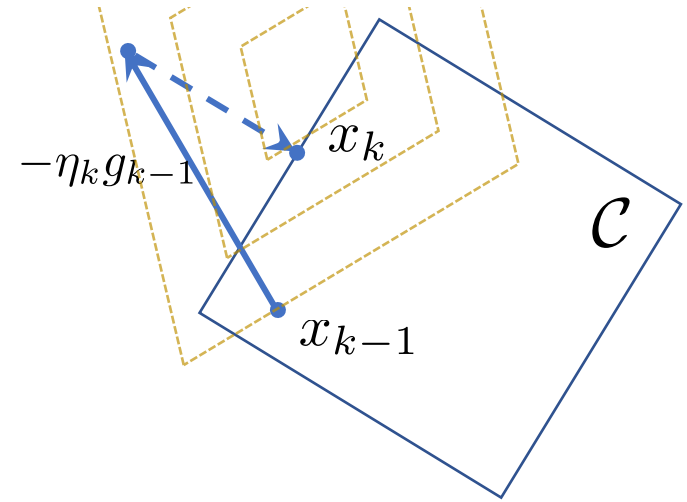
Projected subGradient Descent (PGD)

$$x_0 \in \mathcal{C}$$

For all $k = 1, 2, \dots$

$$g_{k-1} = \text{FO}(x_{k-1}) \in \partial f(x_{k-1})$$

$$x_k = \mathcal{P}_{\mathcal{C}}(x_{k-1} - \eta_k g_{k-1})$$



- Simple and practical. Workhorse of machine learning

- Guarantee: $f\left(\frac{1}{k} \sum_{i=1}^k x_i\right) - \min_{x \in \mathcal{C}} f(x) \leq \epsilon$ if $k \geq \Theta\left(\frac{G^2 D^2}{\epsilon^2}\right)$

- # FO calls = # PO calls = $\Theta\left(\frac{G^2 D^2}{\epsilon^2}\right)$

Is this the best we can do?

- Yes, finding an ε -suboptimal *feasible* solution requires

$$\text{number FO calls} \geq \Omega\left(\frac{G^2 D^2}{\varepsilon^2}\right) \quad [\text{NemYud83}]$$

- No, we don't know how many projections (PO) are required!

$$\text{number PO calls} \geq ?$$

$$\text{PGD: } x_k = \mathcal{P}_C(x_{k-1} - \eta_k g_{k-1})$$

Projection and subgradient descent steps are coupled!

PO calls can be very expensive!

- Example: Nuclear-norm ball constrained low-rank Matrix SVM
- n labeled training samples: sample matrix A_i and label b_i for $i = 1, \dots, n$

$$\min_{\substack{X \in \mathbb{R}^{m \times m} \\ \|X\|_* \leq \lambda}} \frac{1}{n} \sum_{i=1}^n \max(0, 1 - b_i \langle X, A_i \rangle)$$

Hinge loss is
Nonsmooth and Convex

where $\|X\|_* := \sum_{i=1}^m \sigma_i(X)$, where $\sigma_i(X)$ is X 's i -th singular value

- Projection requires full SVD $O(m^3)$

$$\mathcal{C} = \{X \mid \|X\|_* \leq \lambda\}$$

PO can be expensive.
We want PO call efficiency!

Accelerated PGD for L -smooth objectives

L -smooth function

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

- First-order Oracle (FO)

$$\text{FO}(x) = \nabla f(x) \quad \text{gradient}$$

$$f(x_k) - \min_{x \in \mathcal{C}} f(x) \leq \varepsilon \quad \text{if } k \geq \Theta\left(\sqrt{\frac{L}{\varepsilon}}\right)$$

$$\# \text{ FO calls} = \# \text{ PO calls} = O\left(\sqrt{\frac{L}{\varepsilon}}\right)$$

Even $1/\varepsilon$ -smoothness reduces # PO calls to $O\left(\frac{1}{\varepsilon}\right)$

Accelerated Gradient Descent (AGD)

$$z_0 = x_0 \in \mathcal{C}$$

For $k = 1, 2, \dots$

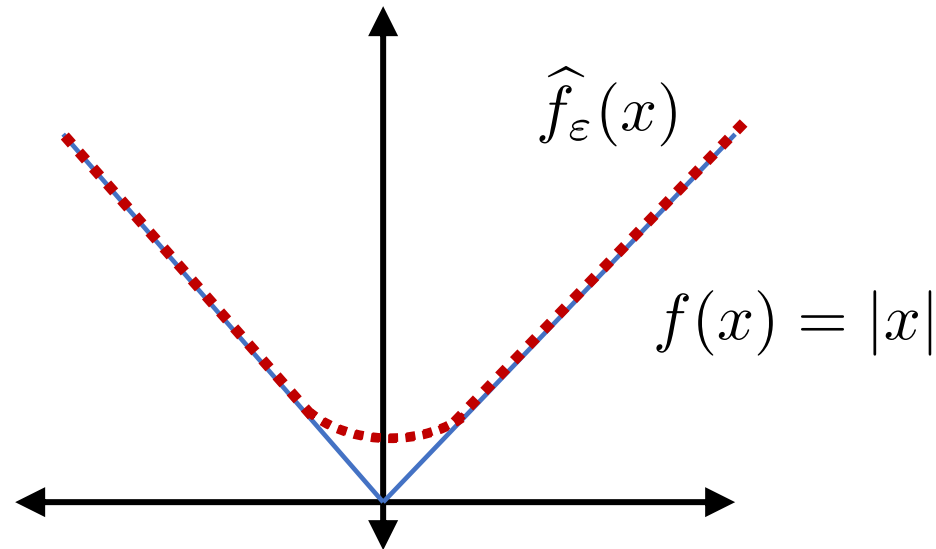
$$y_k = (1 - \gamma_k) x_{k-1} + \gamma_k z_{k-1}$$

$$z_k = \mathcal{P}_{\mathcal{C}}(z_{k-1} - \text{FO}(y_k)/\beta_k)$$

$$x_k = (1 - \gamma_k) x_{k-1} + \gamma_k z_k$$

Randomized Smoothing of Lipschitz functions

$$\hat{f}_\varepsilon(x) := [f * \mathcal{N}(0, \varepsilon^2 d^{-1} \mathbf{I})](x) = \mathbb{E}_{\xi \sim \mathcal{N}(0, \mathbf{I})} [f(x + \varepsilon d^{-\frac{1}{2}} \xi)]$$



Close upper-bound

$$f(x) \leq \hat{f}_\varepsilon(x) \leq f(x) + O(\varepsilon)$$

Convex

$$L = O\left(\frac{\sqrt{d}}{\varepsilon}\right)\text{-smooth}$$

AGD with Randomized Smoothing

$$\widehat{f}_\varepsilon(x) = \mathbb{E}[f(x + \varepsilon d^{-\frac{1}{2}} \xi)]$$

$$\nabla \widehat{f}_\varepsilon(x) \approx \frac{1}{T_\varepsilon} \sum_{i=1}^{T_\varepsilon} \text{FO}(x + \varepsilon d^{-\frac{1}{2}} \xi_i), \text{ where } T_\varepsilon = O\left(\frac{d^{-\frac{1}{4}}}{\varepsilon}\right)$$

AGD + Randomized Smoothing

$$y_k = (1 - \gamma_k) x_{k-1} + \gamma_k z_{k-1}$$

$$\widehat{\nabla}_k = \frac{1}{T_\varepsilon} \sum_{i=1}^{T_\varepsilon} \text{FO}(y_k + \varepsilon d^{-\frac{1}{2}} \xi_i)$$

$$z_k = \mathcal{P}_C \left(z_{k-1} - \widehat{\nabla}_k / \beta_k \right)$$

$$x_k = (1 - \gamma_k) x_{k-1} + \gamma_k z_k$$

$$O\left(\sqrt{\frac{L}{\varepsilon}}\right) = O\left(\frac{d^{1/4}}{\varepsilon}\right) \text{ iterations} \quad L = O\left(\frac{\sqrt{d}}{\varepsilon}\right)$$

$$\times O\left(\frac{d^{-\frac{1}{4}}}{\varepsilon}\right) \text{ FO/grad} = O\left(\frac{1}{\varepsilon^2}\right) \text{ FO calls}$$

$$\times 1 \quad \text{PO/proj} = O\left(\frac{d^{1/4}}{\varepsilon}\right) \text{ PO calls}$$

Dimension dependent due to randomization

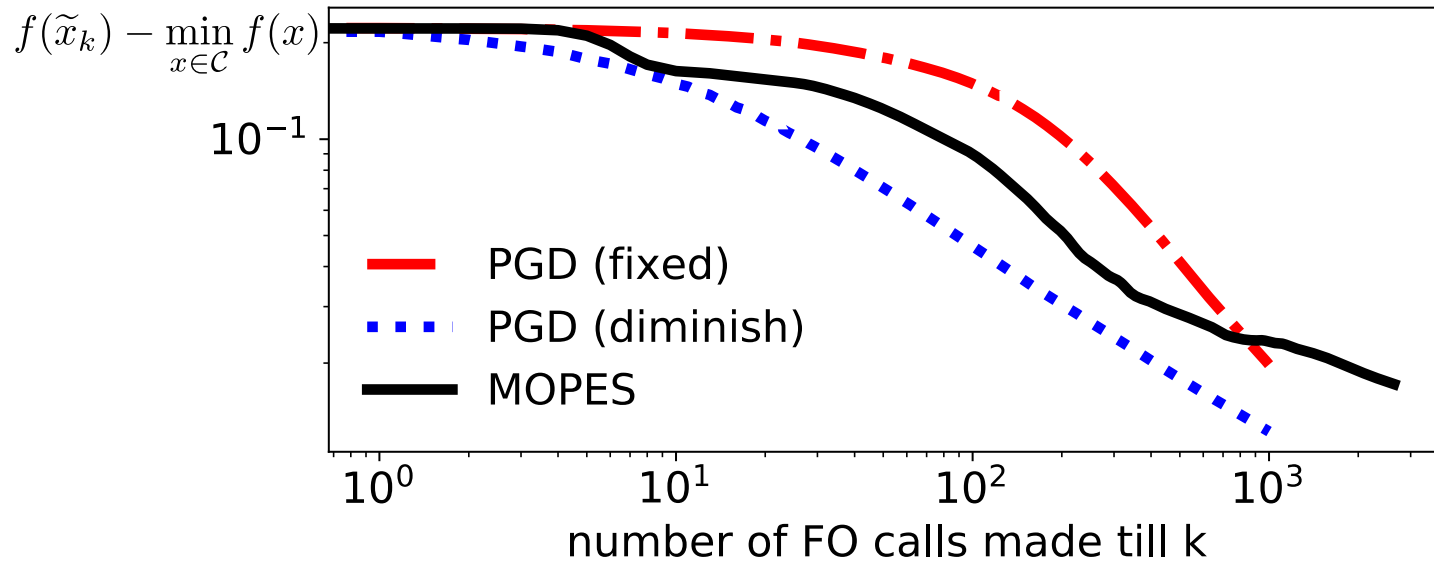
Our result: only $O\left(\frac{1}{\varepsilon}\right)$ projections

Algorithm	Number of FO calls	Number of PO calls
PGD [1]	$O\left(\frac{1}{\varepsilon^2}\right)$	$O\left(\frac{1}{\varepsilon^2}\right)$
Randomized Smoothing [2]	$O\left(\frac{1}{\varepsilon^2}\right)$	$O\left(\frac{d^{1/4}}{\varepsilon}\right)$
MOPES (ours)	$O\left(\frac{1}{\varepsilon^2}\right)$	$O\left(\frac{1}{\varepsilon}\right)$
<i>Lower Bound</i>	$\Omega\left(\frac{1}{\varepsilon^2}\right)$	<i>Open Question!</i>

[1] Nemirovski, and Yudin, "Problem complexity and method efficiency in optimization", Wiley-Interscience 1983

[2] Duchi, Bartlett, and Wainwright, "Randomized Smoothing for stochastic optimization", SIAM Opt 2012

Sub-optimality at k



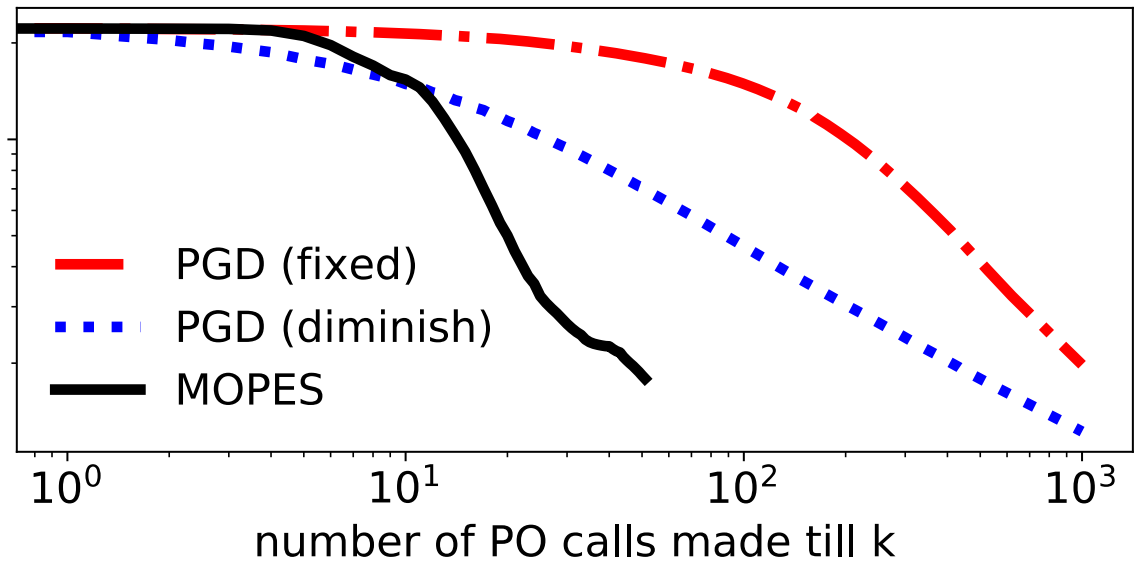
MOPES (our) and PGD use similar number of FO calls

MOPES (our) uses fewer PO calls than PGD!

Sub-optimality at k

$$f(\tilde{x}_k) - \min_{x \in \mathcal{C}} f(x)$$

10^{-1}



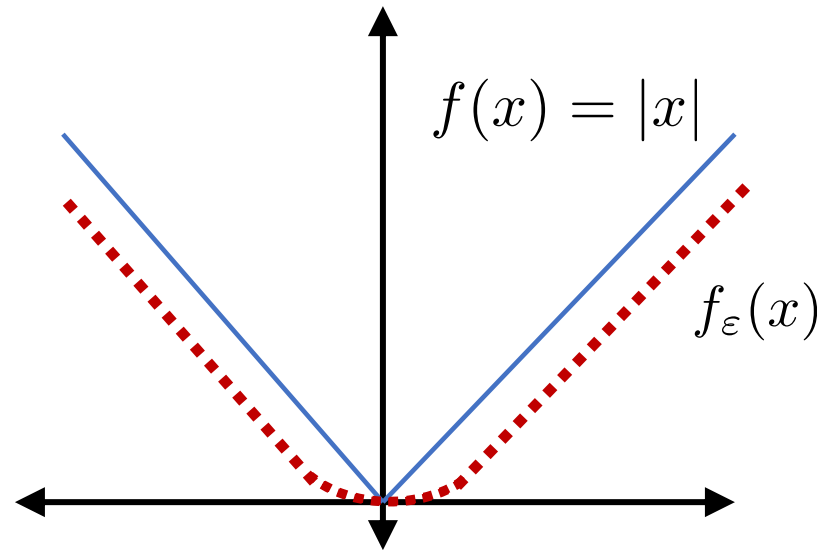
Moreau-Yosida Smoothing of Lipschitz functions

$$f_\varepsilon(x) := \min_{x'} f(x') + \frac{1}{2\varepsilon} \|x' - x\|^2$$

Infimal convolution

Approximate lower-bound

$$f(x) - O(\varepsilon) \leq f_\varepsilon(x) \leq f(x)$$



Convex

$$L = O\left(\frac{1}{\varepsilon}\right)\text{-smooth}$$

AGD + Moreau Smoothing

$$f_\varepsilon(x) := \min_{x'} f(x') + \frac{1}{2\varepsilon} \|x' - x\|^2$$

$$\nabla f_\varepsilon(x) = \frac{x - \text{prox}_f(x)}{\varepsilon}, \text{ where } \text{prox}_f(x) = \arg \min_{x'} f(x') + \frac{1}{2\varepsilon} \|x' - x\|^2$$

subgradient method

AGD + Moreau Smoothing

$$y_k = (1 - \gamma_k) x_{k-1} + \gamma_k z_{k-1}$$

$$y'_k \approx \text{prox}_f(y_k)$$

$$z_k = \mathcal{P}_C \left(z_{k-1} - \frac{y_k - y'_k}{\varepsilon \beta_k} \right)$$

$$x_k = (1 - \gamma_k) x_{k-1} + \gamma_k z_k$$

$$O\left(\sqrt{\frac{L}{\varepsilon}}\right) = O\left(\frac{1}{\varepsilon}\right) \text{ iterations} \quad L = O\left(\frac{1}{\varepsilon}\right)$$

$$\times O\left(\frac{1}{\varepsilon^2}\right) \text{ FO/prox} = O\left(\frac{1}{\varepsilon^3}\right) \text{ FO calls}$$

$$\times 1 \text{ PO/proj} = O\left(\frac{1}{\varepsilon}\right) \text{ PO calls}$$

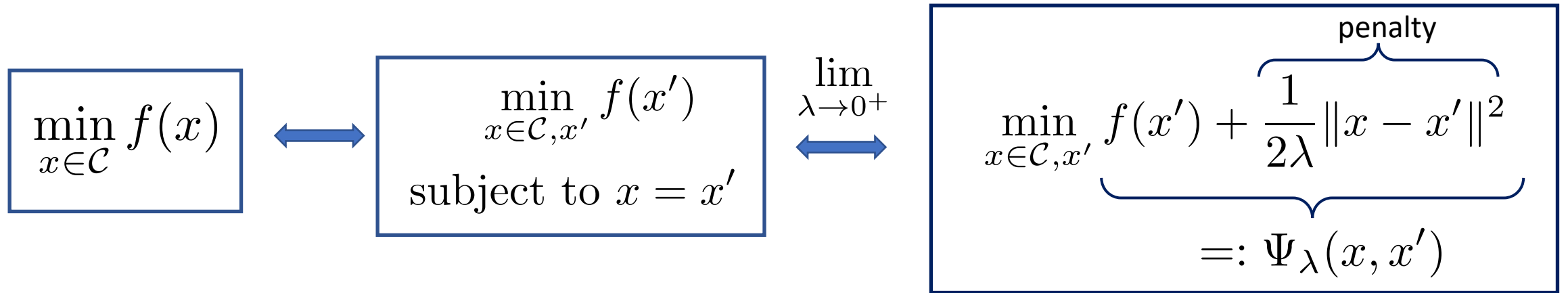
Are we wasting FO calls?

$$\min_{x \in \mathcal{X}} \left[f_\varepsilon(x) := \left[\min_{x'} f(x') + \frac{1}{2\varepsilon} \|x' - x\|^2 \right] \right]$$

Each inner problem requires $O\left(\frac{1}{\varepsilon^2}\right)$ FO calls

$$\min_{x \in \mathcal{X}, x'} \left[\Psi_\varepsilon(x, x') := f(x') + \frac{1}{2\varepsilon} \|x' - x\|^2 \right]$$

Decoupling the objective and the constraint



- Moreau-Yosida regularization theory suggests $\lambda = \varepsilon$

$$f(\tilde{x}) - \min_{x \in \mathcal{C}} f(x) \leq [\Psi_\varepsilon(\tilde{x}, \tilde{x}') - \min_{x \in \mathcal{C}, x'} \Psi_\varepsilon(x)] + O(\varepsilon)$$

- **NEW GOAL:** Find an ε -suboptimal solution (\tilde{x}, \tilde{x}') to

$$\min_{x \in \mathcal{C}, x'} \Psi_\varepsilon(x, x')$$

Reducing Number of Projections

$\min_{x \in \mathcal{C}, x' \in \mathbb{R}^d}$	$\Psi_\varepsilon(x, x')$ <p>Composite Objective</p> <p>Gradient Sliding [Lan16]</p>	=	$f(x')$ <p>Unconstrained Convex</p> <p>subGradient Descent (no projection)</p>	+	$\frac{1}{2\varepsilon} \ x - x'\ ^2$ <p>Constrained $L = 2/\varepsilon$-smooth Convex</p> <p>Accelerated Projected Gradient Descent (f not in objective)</p>
# FO =	$O\left(\frac{1}{\varepsilon^2}\right)$	=	$O\left(\frac{1}{\varepsilon^2}\right)$	+	0
# PO =	$O\left(\frac{1}{\varepsilon}\right)$	=	0	+	$O\left(\sqrt{\frac{L}{\varepsilon}}\right) = O\left(\frac{1}{\varepsilon}\right)$

Projection Efficient Subgradient method (MOPES)

$$\min_{x \in \mathcal{C}, x'} [\Psi_\varepsilon(x, x') = f(x') + \underbrace{\frac{1}{2\varepsilon} \|x - x'\|^2}_{(L = 1/\varepsilon)\text{-Smooth } \phi_\varepsilon(x, x')}]$$

$$\bar{x} := (x, x')$$

Accelerated Proximal Gradient Descent $O\left(\sqrt{\frac{L}{\varepsilon}}\right) = O\left(\frac{1}{\varepsilon}\right)$ iterations

$$\bar{y}_k = (1 - \gamma_k) \bar{x}_{k-1} + \gamma_k \bar{z}_{k-1}$$

$$z_k = \mathcal{P}_{\mathcal{C}}(z_{k-1} - \nabla_x \phi_\varepsilon(\bar{y}_k) / \beta_k)$$

$$\hat{z}'_k = (z'_{k-1} - \nabla_{x'} \phi_\lambda(\bar{y}_k) / \beta_k)$$

$$z'_k \stackrel{\sim}{=} \text{prox}_{f/\beta_k}(\hat{z}'_k) := \arg \min_{z'} f(z') + \frac{\beta_k}{2} \|z' - \hat{z}'_k\|^2$$

$$\bar{x}_k = (1 - \gamma_k) \bar{x}_{k-1} + \gamma_k \bar{z}_k$$

projection

x 1 PO/proj = $O\left(\frac{1}{\varepsilon}\right)$ PO calls

Inexact proximal
[Lan16]

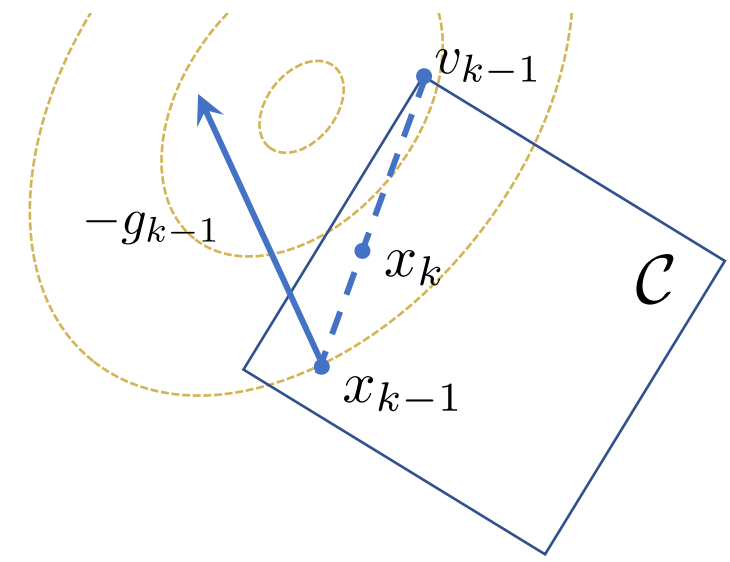
x $O\left(\frac{1}{\varepsilon}\right)$ FO/prox = $O\left(\frac{1}{\varepsilon^2}\right)$ FO calls

Frank-Wolfe (FW) method for Smooth Convex

- Alternative to PO is Linear Minimization Oracle (LMO) $\text{LMO}(g) \in \arg \min_{x \in \mathcal{C}} \langle g, x \rangle$
- LMO is often cheaper than PO, e.g. nuclear-norm ball

$$\begin{aligned} &x_0 \in \mathcal{C} \\ &\text{For all } k = 1, 2, \dots \\ &\quad g_{k-1} = \text{FO}(x_{k-1}) \in \partial f(x_{k-1}) \\ &\quad v_{k-1} = \text{LMO}(g_{k-1}) \in \arg \min_{x \in \mathcal{C}} \langle g_{k-1}, x \rangle \\ &\quad x_k = (1 - \alpha_k) \cdot x_{k-1} + \alpha_k \cdot v_{k-1} \end{aligned}$$

$O\left(\frac{L}{\varepsilon}\right)$ iterations for L -smooth function



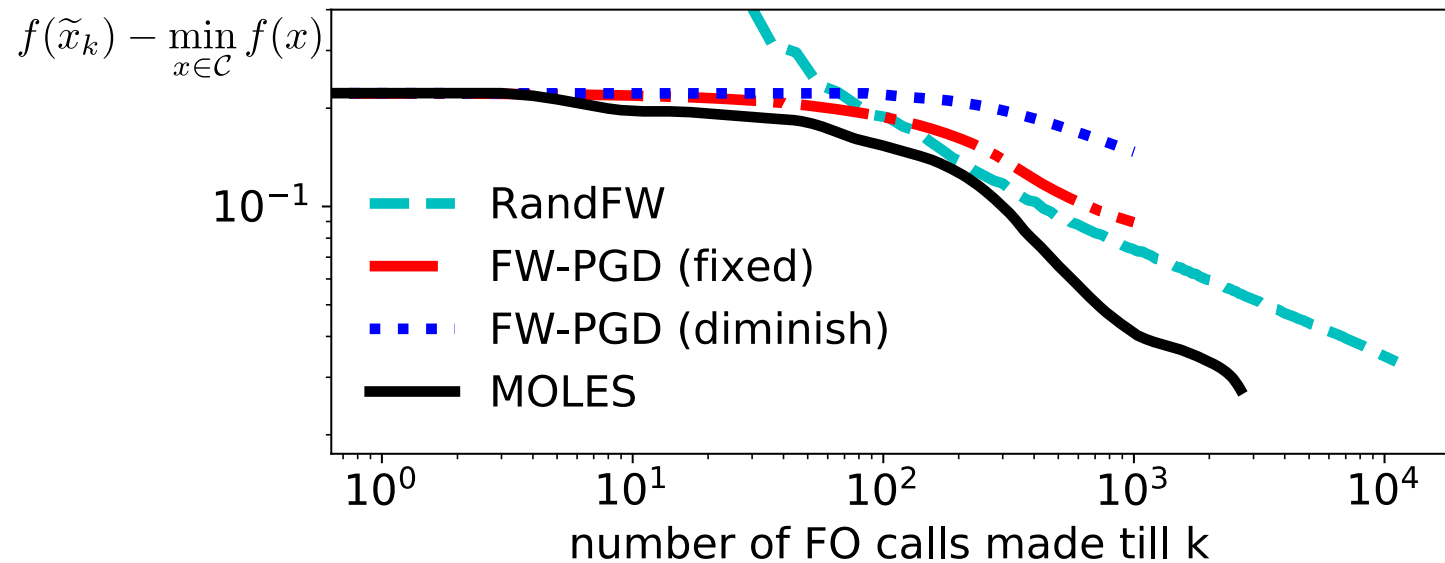
FAILS for Nonsmooth objectives!

Optimal Nonsmooth FW method uses $O\left(\frac{1}{\varepsilon^2}\right)$ LMO and FO calls

Algorithm	Number of FO calls	Number of LMO calls
FW-PGD [folklore]	$O\left(\frac{1}{\varepsilon^2}\right)$	$O\left(\frac{1}{\varepsilon^4}\right)$
Randomized Smoothing [3]	$O\left(\frac{\sqrt{d}}{\varepsilon^4}\right)$	$O\left(\frac{\sqrt{d}}{\varepsilon^2}\right)$
MOLES (ours)	$O\left(\frac{1}{\varepsilon^2}\right)$	$O\left(\frac{1}{\varepsilon^2}\right)$
<i>Lower Bound</i> [3]	$\Omega\left(\frac{1}{\varepsilon^2}\right)$	$\Omega\left(\frac{1}{\varepsilon^2}\right)$

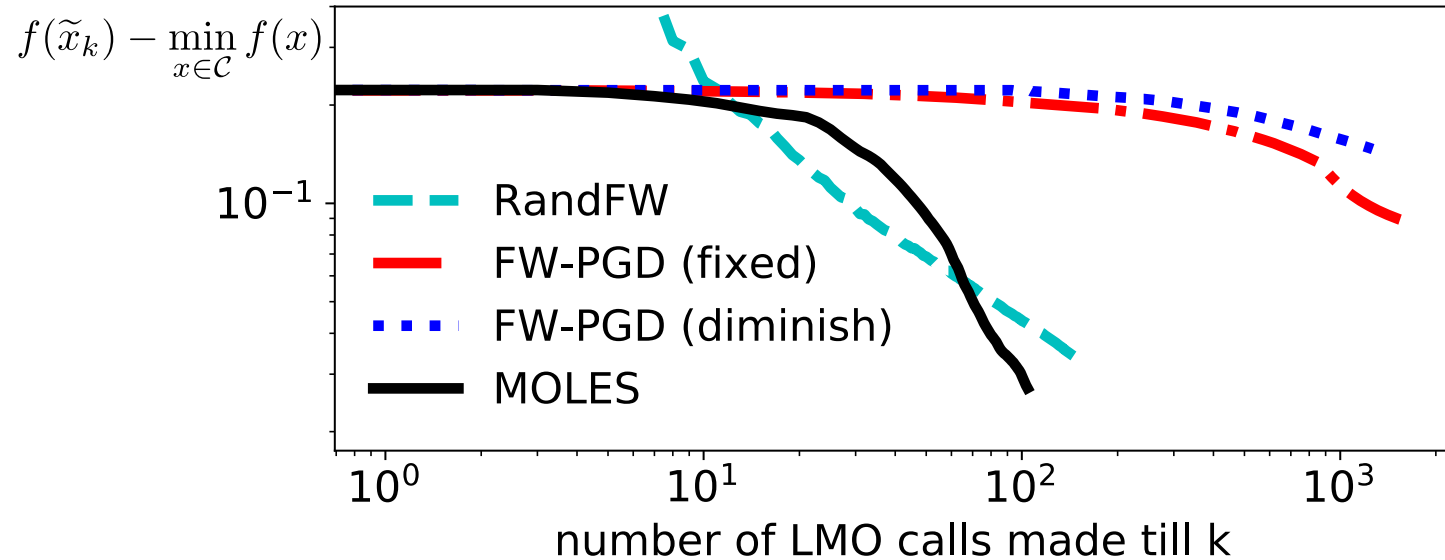
MOLES = MOPES + Inexact Projections using FW method

Sub-optimality at k



MOLES (our) uses the least number of FO and LMO calls.

Sub-optimality at k



Optimal Nonsmooth Frank-Wolfe method (MOLES)

$$\min_{x \in \mathcal{C}, x'} [\Psi_\varepsilon(x, x') = f(x') + \underbrace{\frac{1}{2\varepsilon} \|x - x'\|^2}_{(L = 1/\varepsilon)\text{-Smooth } \phi_\varepsilon(x, x')}]$$

$$\bar{x} := (x, x')$$

Projection is a $\frac{1}{\varepsilon}$ -smooth problem

$$\mathcal{P}_{\mathcal{X}}(x) = \arg \min_{x_{\mathcal{P}} \in \mathcal{X}} \frac{L}{2} \|x_{\mathcal{P}} - x\|^2$$

Inexact MOPES

Inexact
projection
[LZ16]

$$\bar{y}_k = (1 - \gamma_k) \bar{x}_{k-1} + \gamma_k \bar{z}_{k-1}$$

$$z_k \stackrel{\sim}{=} \mathcal{P}_{\mathcal{C}}(z_{k-1} - \nabla_x \phi_\varepsilon(\bar{y}_k) / \beta_k)$$

Inexact
proximal

$$\hat{z}'_k = (z'_{k-1} - \nabla_{x'} \phi_\lambda(\bar{y}_k) / \beta_k)$$

$$z'_k \stackrel{\approx}{=} \text{prox}_{f/\beta_k}(\hat{z}'_k)$$

$$\bar{x}_k = (1 - \gamma_k) \bar{x}_{k-1} + \gamma_k \bar{z}_k$$

$$O\left(\sqrt{\frac{L}{\varepsilon}}\right) = O\left(\frac{1}{\varepsilon}\right) \text{ iterations}$$

$$\times O\left(\frac{1}{\varepsilon}\right) \text{ LMO/proj} = O\left(\frac{1}{\varepsilon^2}\right) \text{ LMO calls}$$

$$\times O\left(\frac{1}{\varepsilon}\right) \text{ FO/prox} = O\left(\frac{1}{\varepsilon^2}\right) \text{ FO calls}$$

Summary

- Constrained convex minimization method using $O(1/\varepsilon)$ projections
- Optimal nonsmooth FW method using $O(1/\varepsilon^2)$ FO and LMO calls
- We obtain similar complexities for Stochastic FO

Open questions

- Lower bound for PO calls complexity
- Extension to non-Euclidean geometry

Thank you!