Analysis of the Frank-Wolfe Method for Convex Composite Optimization involving a Logarithmically-Homogeneous Barrier

Renbo Zhao

MIT Operations Research Center

Joint work with Robert M. Freund (MIT Sloan School of Management)

SIAM Conference on Optimization July, 2021

Consider the following convex composite optimization problem:

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right] \tag{P}$$

2/20

Consider the following convex composite optimization problem:

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

2/20

 $ightarrow f: \mathbb{R}^m \to \mathbb{R} \cup \{+\infty\}$ is a θ -logarithmically-homogeneous self-concordant barrier (θ -LHSCB) for some regular cone $\mathcal{K} \subseteq \mathbb{R}^m$,

Consider the following convex composite optimization problem:

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right] \tag{P}$$

2/20

 $ightarrow f: \mathbb{R}^m \to \mathbb{R} \cup \{+\infty\}$ is a θ -logarithmically-homogeneous self-concordant barrier (θ -LHSCB) for some regular cone $\mathcal{K} \subseteq \mathbb{R}^m$,

 $\triangleright \mathsf{A}: \mathbb{R}^n \to \mathbb{R}^m$ is a linear operator (not necessarily invertible),

Consider the following convex composite optimization problem:

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

2 / 20

- $ightarrow f: \mathbb{R}^m \to \mathbb{R} \cup \{+\infty\}$ is a θ -logarithmically-homogeneous self-concordant barrier (θ -LHSCB) for some regular cone $\mathcal{K} \subseteq \mathbb{R}^m$,
- $\triangleright \mathsf{A}: \mathbb{R}^n \to \mathbb{R}^m$ is a linear operator (not necessarily invertible),
- $ightarrow h: \mathbb{R}^n \to \mathbb{R} \cup \{+\infty\}$ is a proper, closed and convex (but possibly non-smooth) function, and dom h is nonempty convex and compact.

Consider the following convex composite optimization problem:

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

- $ightarrow f: \mathbb{R}^m \to \mathbb{R} \cup \{+\infty\}$ is a θ -logarithmically-homogeneous self-concordant barrier (θ -LHSCB) for some regular cone $\mathcal{K} \subseteq \mathbb{R}^m$,
- $\triangleright \mathsf{A}: \mathbb{R}^n \to \mathbb{R}^m$ is a linear operator (not necessarily invertible),
- $ightarrow h: \mathbb{R}^n \to \mathbb{R} \cup \{+\infty\}$ is a proper, closed and convex (but possibly non-smooth) function, and dom h is nonempty convex and compact.
- \triangleright We recover the traditional problem setting for Frank-Wolfe when h is the indicator function $h := \iota_{\mathcal{X}}$ of a compact convex set \mathcal{X} .

Consider the following convex composite optimization problem:

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right] \tag{P}$$

- $ightarrow f: \mathbb{R}^m \to \mathbb{R} \cup \{+\infty\}$ is a θ -logarithmically-homogeneous self-concordant barrier (θ -LHSCB) for some regular cone $\mathcal{K} \subseteq \mathbb{R}^m$,
- $\triangleright \mathsf{A}: \mathbb{R}^n \to \mathbb{R}^m$ is a linear operator (not necessarily invertible),
- $ightarrow h: \mathbb{R}^n \to \mathbb{R} \cup \{+\infty\}$ is a proper, closed and convex (but possibly non-smooth) function, and dom h is nonempty convex and compact.
- \triangleright We recover the traditional problem setting for Frank-Wolfe when h is the indicator function $h := \iota_{\mathcal{X}}$ of a compact convex set \mathcal{X} .
- ▷ Assume dom $F \neq \emptyset$, so at least one minimizer $x^* \in \text{dom } F$ exists, and define $F^* := F(x^*)$.

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

 $\triangleright \mathcal{X}$ is a nonempty convex and compact set.

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

3 / 20

 $\triangleright \mathcal{X}$ is a nonempty convex and compact set.

 $\triangleright f$ is L-smooth w.r.t. $\|\cdot\|$ on \mathcal{X} , which then implies

$$f(x') \leq f(x) + \langle \nabla f(x), x' - x \rangle + (L/2) \|x' - x\|^2, \quad \forall x', x \in \mathcal{X}.$$
 (LSm)

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

3 / 20

 $\triangleright \mathcal{X}$ is a nonempty convex and compact set.

 $\succ f \text{ is } L\text{-smooth w.r.t. } \|\cdot\| \text{ on } \mathcal{X} \text{, which then implies}$ $f(x') \leq f(x) + \langle \nabla f(x), x' - x \rangle + (L/2) \|x' - x\|^2, \quad \forall x', x \in \mathcal{X}.$ (LSm)

 \triangleright At iteration k of FW, $x^k \in \mathcal{X}$ and the method does the following:

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

 $\triangleright \mathcal{X}$ is a nonempty convex and compact set.

 $\triangleright f$ is L-smooth w.r.t. $\|\cdot\|$ on \mathcal{X} , which then implies

$$f(x') \leq f(x) + \langle \nabla f(x), x' - x \rangle + (L/2) \|x' - x\|^2, \quad \forall x', x \in \mathcal{X}. \quad (\texttt{LSm})$$

 \triangleright At iteration k of FW, $x^k \in \mathcal{X}$ and the method does the following:

• Compute

$$v^k \in \operatorname{arg\,min}_{x \in \mathcal{X}} \langle \nabla f(x^k), x \rangle$$

by solving a linear-optimization sub-problem.

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

 $\triangleright \mathcal{X}$ is a nonempty convex and compact set.

 $\triangleright f$ is L-smooth w.r.t. $\|\cdot\|$ on \mathcal{X} , which then implies

$$f(x') \leq f(x) + \langle \nabla f(x), x' - x \rangle + (L/2) \|x' - x\|^2, \quad \forall x', x \in \mathcal{X}. \quad (\texttt{LSm})$$

 \triangleright At iteration k of FW, $x^k \in \mathcal{X}$ and the method does the following:

• Compute

$$v^k \in \operatorname{arg\,min}_{x \in \mathcal{X}} \langle \nabla f(x^k), x \rangle$$

by solving a linear-optimization sub-problem.

• Determine step-length
$$\alpha^k \in [0, 1]$$
.

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

 $\triangleright \mathcal{X}$ is a nonempty convex and compact set.

 $\triangleright f$ is L-smooth w.r.t. $\|\cdot\|$ on \mathcal{X} , which then implies

$$f(x') \leq f(x) + \langle \nabla f(x), x' - x \rangle + (L/2) \|x' - x\|^2, \quad \forall x', x \in \mathcal{X}. \quad (\texttt{LSm})$$

 \triangleright At iteration k of FW, $x^k \in \mathcal{X}$ and the method does the following:

• Compute

$$v^k \in \operatorname{arg\,min}_{x \in \mathcal{X}} \langle \nabla f(x^k), x \rangle$$

by solving a linear-optimization sub-problem.

- Determine step-length $\alpha^k \in [0, 1]$.
- Update $x^{k+1} = (1 \alpha_k)x^k + \alpha^k v^k$.

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

4 / 20

 \triangleright The step-size α_k is typically chosen in one of two ways:

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

4 / 20

 \triangleright The step-size α_k is typically chosen in one of two ways:

• Fixed step-size, such as the standard step-size $\alpha_k = 2/(k+2)$, or

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

4 / 20

 \triangleright The step-size α_k is typically chosen in one of two ways:

- Fixed step-size, such as the standard step-size $\alpha_k = 2/(k+2)$, or
- Adaptive step-size, such as $\alpha_k = \min\{G_k/C_k, 1\}$, where

$$G_k := \langle \nabla f(x^k), x^k - v^k \rangle$$
 and $C_k := L \|v^k - x^k\|^2$.

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

- \triangleright The step-size α_k is typically chosen in one of two ways:
 - Fixed step-size, such as the standard step-size $\alpha_k = 2/(k+2)$, or
 - Adaptive step-size, such as $\alpha_k = \min\{G_k/C_k, 1\}$, where

$$G_k := \langle \nabla f(x^k), x^k - v^k \rangle$$
 and $C_k := L \| v^k - x^k \|^2$.

 \triangleright FW is very useful in "sparse" or otherwise "structured" optimization where \mathcal{X} has special structure, e.g., probability simplex or spectrahedron.

$$\min_{x \in \mathcal{X}} f(x) \tag{tP}$$

 \triangleright The step-size α_k is typically chosen in one of two ways:

- Fixed step-size, such as the standard step-size $\alpha_k = 2/(k+2)$, or
- Adaptive step-size, such as $\alpha_k = \min\{G_k/C_k, 1\}$, where

$$G_k := \langle \nabla f(x^k), x^k - v^k \rangle$$
 and $C_k := L \| v^k - x^k \|^2$.

- \triangleright FW is very useful in "sparse" or otherwise "structured" optimization where \mathcal{X} has special structure, e.g., probability simplex or spectrahedron.

However, note that all of these works assume that f is L-smooth.

Renbo Zhao (MIT ORC)

5 / 20

Khachiyan, L.G.: Rounding of polytopes in the real number model of computation. *Mathematics of Operations Research* 21(2), 307–320 (1996) (Elegant analysis of the FW method with exact line-search for D-optimal design)

- Khachiyan, L.G.: Rounding of polytopes in the real number model of computation. *Mathematics of Operations Research* 21(2), 307–320 (1996) (Elegant analysis of the FW method with exact line-search for D-optimal design)
- \rhd Dvurechensky, P., Ostroukhov, P., Safin, K., Shtern, S., Staudigl, M.: Self-concordant analysis of Frank-Wolfe algorithms. *Proc. ICML* , pp. 2814–2824 (2020)

- Khachiyan, L.G.: Rounding of polytopes in the real number model of computation. *Mathematics of Operations Research* 21(2), 307–320 (1996) (Elegant analysis of the FW method with exact line-search for D-optimal design)
- \rhd Dvurechensky, P., Ostroukhov, P., Safin, K., Shtern, S., Staudigl, M.: Self-concordant analysis of Frank-Wolfe algorithms. *Proc. ICML* , pp. 2814–2824 (2020)
- ▷ Dvurechensky et al. (2020) proposed and analyzed a FW method for the *whole* class of self-concordant functions. However, when specialized to D-optimal design, their complexity bound is very different from Khachiyan's result, and lacks the affine-invariance property.

- Khachiyan, L.G.: Rounding of polytopes in the real number model of computation. *Mathematics of Operations Research* 21(2), 307–320 (1996) (Elegant analysis of the FW method with exact line-search for D-optimal design)
- \rhd Dvurechensky, P., Ostroukhov, P., Safin, K., Shtern, S., Staudigl, M.: Self-concordant analysis of Frank-Wolfe algorithms. *Proc. ICML* , pp. 2814–2824 (2020)
- ▷ Dvurechensky et al. (2020) proposed and analyzed a FW method for the *whole* class of self-concordant functions. However, when specialized to D-optimal design, their complexity bound is very different from Khachiyan's result, and lacks the affine-invariance property.
- \triangleright We identified the *logarithmic-homogeneity* as the key element in Khachiyan's analysis, and proposed a (generalized) FW method with adaptive step-size for the much broader problem class (P).

- Khachiyan, L.G.: Rounding of polytopes in the real number model of computation. *Mathematics of Operations Research* 21(2), 307–320 (1996) (Elegant analysis of the FW method with exact line-search for D-optimal design)
- \rhd Dvurechensky, P., Ostroukhov, P., Safin, K., Shtern, S., Staudigl, M.: Self-concordant analysis of Frank-Wolfe algorithms. *Proc. ICML* , pp. 2814–2824 (2020)
- ▷ Dvurechensky et al. (2020) proposed and analyzed a FW method for the *whole* class of self-concordant functions. However, when specialized to D-optimal design, their complexity bound is very different from Khachiyan's result, and lacks the affine-invariance property.
- \triangleright We identified the *logarithmic-homogeneity* as the key element in Khachiyan's analysis, and proposed a (generalized) FW method with adaptive step-size for the much broader problem class (P).
- ▷ Our complexity bound essentially recovers Khachiyan's result, and is affine-invariant (along with other desirable properties).

▷ Let $\mathcal{K} \subsetneq \mathbb{R}^m$ be a regular cone, i.e., \mathcal{K} is closed, convex, pointed and has nonempty interior.

6 / 20

- ▷ Let $\mathcal{K} \subsetneq \mathbb{R}^m$ be a regular cone, i.e., \mathcal{K} is closed, convex, pointed and has nonempty interior.
- $\triangleright f$ is a θ -LHSCB on \mathcal{K} with *complexity parameter* $\theta \ge 1$ if f is three-times differentiable and strictly convex on int \mathcal{K} , and satisfies

6 / 20

- ▷ Let $\mathcal{K} \subsetneq \mathbb{R}^m$ be a regular cone, i.e., \mathcal{K} is closed, convex, pointed and has nonempty interior.
- $\triangleright f$ is a θ -LHSCB on \mathcal{K} with *complexity parameter* $\theta \ge 1$ if f is three-times differentiable and strictly convex on int \mathcal{K} , and satisfies

$$\left| D^3 f(u)[w,w,w] \right| \le 2(\langle H(u)w,w\rangle)^{3/2} \quad \forall \, u \in \operatorname{int} \mathcal{K}, \, \forall \, w \in \mathbb{R}^m,$$

2 $f(u_k) \to \infty$ for any $\{u_k\}_{k \ge 1} \subseteq \operatorname{int} \mathcal{K}$ such that $u_k \to u \in \operatorname{bd} \mathcal{K}$,

$$3 \ f(tu) = f(u) - \theta \ln(t) \ \forall u \in \operatorname{int} \mathcal{K}, \, \forall t > 0 ,$$

where H(u) denotes the Hessian of f at $u \in int \mathcal{K}$.

- ▷ Let $\mathcal{K} \subsetneq \mathbb{R}^m$ be a regular cone, i.e., \mathcal{K} is closed, convex, pointed and has nonempty interior.
- $\triangleright f$ is a θ -LHSCB on \mathcal{K} with *complexity parameter* $\theta \ge 1$ if f is three-times differentiable and strictly convex on int \mathcal{K} , and satisfies

$$\left| D^3 f(u)[w,w,w] \right| \le 2(\langle H(u)w,w\rangle)^{3/2} \quad \forall \, u \in \operatorname{int} \mathcal{K}, \, \forall \, w \in \mathbb{R}^m,$$

2 $f(u_k) \to \infty$ for any $\{u_k\}_{k \ge 1} \subseteq \operatorname{int} \mathcal{K}$ such that $u_k \to u \in \operatorname{bd} \mathcal{K}$,

3
$$f(tu) = f(u) - \theta \ln(t) \quad \forall u \in \operatorname{int} \mathcal{K}, \forall t > 0$$
,

where H(u) denotes the Hessian of f at $u \in int \mathcal{K}$.

 $\,\triangleright\,$ Two prototypical examples:

- ▷ Let $\mathcal{K} \subsetneq \mathbb{R}^m$ be a regular cone, i.e., \mathcal{K} is closed, convex, pointed and has nonempty interior.
- $\triangleright f$ is a θ -LHSCB on \mathcal{K} with *complexity parameter* $\theta \ge 1$ if f is three-times differentiable and strictly convex on int \mathcal{K} , and satisfies

$$\left| D^3 f(u)[w,w,w] \right| \le 2(\langle H(u)w,w\rangle)^{3/2} \quad \forall \, u \in \operatorname{int} \mathcal{K}, \, \forall \, w \in \mathbb{R}^m,$$

2 $f(u_k) \to \infty$ for any $\{u_k\}_{k \ge 1} \subseteq \operatorname{int} \mathcal{K}$ such that $u_k \to u \in \operatorname{bd} \mathcal{K}$,

3
$$f(tu) = f(u) - \theta \ln(t) \quad \forall u \in \operatorname{int} \mathcal{K}, \, \forall t > 0$$
,

where H(u) denotes the Hessian of f at $u \in int \mathcal{K}$.

 \triangleright Two prototypical examples:

•
$$f(U) = -\ln \det(U)$$
 for $U \in \mathcal{K} := \mathbb{S}^k_+$ and $\theta = k_+$

- ▷ Let $\mathcal{K} \subsetneq \mathbb{R}^m$ be a regular cone, i.e., \mathcal{K} is closed, convex, pointed and has nonempty interior.
- $\triangleright f$ is a θ -LHSCB on \mathcal{K} with *complexity parameter* $\theta \geq 1$ if f is three-times differentiable and strictly convex on int \mathcal{K} , and satisfies

$$\left| D^3 f(u)[w,w,w] \right| \le 2(\langle H(u)w,w\rangle)^{3/2} \quad \forall \, u \in \operatorname{int} \mathcal{K}, \, \forall \, w \in \mathbb{R}^m,$$

2 $f(u_k) \to \infty$ for any $\{u_k\}_{k \ge 1} \subseteq \operatorname{int} \mathcal{K}$ such that $u_k \to u \in \operatorname{bd} \mathcal{K}$,

3
$$f(tu) = f(u) - \theta \ln(t) \quad \forall u \in \operatorname{int} \mathcal{K}, \, \forall t > 0$$
,

where H(u) denotes the Hessian of f at $u \in int \mathcal{K}$.

 \triangleright Two prototypical examples:

•
$$f(U) = -\ln \det(U)$$
 for $U \in \mathcal{K} := \mathbb{S}^k_+$ and $\theta = k$,

• $f(u) = -\sum_{j=1}^{m} w_j \ln(u_j)$ for $u \in \mathcal{K} := \mathbb{R}^m_+$ and $\theta = \sum_{j=1}^{m} w_j$ where $w_1, \ldots, w_n \ge 1$.

A Motivating Example: *D*-optimal Design

$$\begin{aligned} \max_{p} h(p) &\triangleq \ln \det \left(\sum_{i=1}^{m} p_{i} a_{i} a_{i}^{\top} \right) \\ \text{s.t.} \quad \sum_{i=1}^{m} p_{i} = 1, \ p_{i} \geq 0, \ \forall i \in [m]. \end{aligned} \tag{D-OPT}$$

A Motivating Example: *D*-optimal Design

$$\begin{aligned} \max_{p} \quad h(p) &\triangleq \ln \det \left(\sum_{i=1}^{m} p_{i} a_{i} a_{i}^{\top} \right) \\ \text{s.t.} \quad \sum_{i=1}^{m} p_{i} = 1, \ p_{i} \geq 0, \ \forall i \in [m]. \end{aligned} \tag{D-OPT}$$

 \triangleright Problem data: $\{a_i\}_{i=1}^m \subseteq \mathbb{R}^n$.

A Motivating Example: *D*-optimal Design

$$\begin{aligned} \max_{p} h(p) &\triangleq \ln \det \left(\sum_{i=1}^{m} p_{i} a_{i} a_{i}^{\top} \right) \\ \text{s.t.} \quad \sum_{i=1}^{m} p_{i} = 1, \ p_{i} \geq 0, \ \forall i \in [m]. \end{aligned} \tag{D-OPT}$$

7 / 20

- \triangleright Problem data: $\{a_i\}_{i=1}^m \subseteq \mathbb{R}^n$.
- ▷ Arises in many places, including optimal experimental design, and as the dual problem of the minimum volume enclosing ellipsoid (MVEE) problem.
A Motivating Example: *D*-optimal Design

$$\begin{aligned} \max_{p} h(p) &\triangleq \ln \det \left(\sum_{i=1}^{m} p_{i} a_{i} a_{i}^{\top} \right) \\ \text{s.t.} \quad \sum_{i=1}^{m} p_{i} = 1, \ p_{i} \geq 0, \ \forall i \in [m]. \end{aligned} \tag{D-OPT}$$

- \triangleright Problem data: $\{a_i\}_{i=1}^m \subseteq \mathbb{R}^n$.
- ▷ Arises in many places, including optimal experimental design, and as the dual problem of the minimum volume enclosing ellipsoid (MVEE) problem.
- ▷ Khachiyan (1996) proposed a "barycentric coordinate ascent" method with exact line-search, which is actually FW with exact line-search. Method works remarkably well both in theory and practice: it computes an ε -optimal solution of (D-OPT) in (essentially) $O(n^2/\varepsilon)$ iterations.

A Motivating Example: *D*-optimal Design

$$\begin{aligned} \max_{p} h(p) &\triangleq \ln \det \left(\sum_{i=1}^{m} p_{i} a_{i} a_{i}^{\top} \right) \\ \text{s.t.} \quad \sum_{i=1}^{m} p_{i} = 1, \ p_{i} \geq 0, \ \forall i \in [m]. \end{aligned} \tag{D-OPT}$$

- \triangleright Problem data: $\{a_i\}_{i=1}^m \subseteq \mathbb{R}^n$.
- ▷ Arises in many places, including optimal experimental design, and as the dual problem of the minimum volume enclosing ellipsoid (MVEE) problem.
- ▷ Khachiyan (1996) proposed a "barycentric coordinate ascent" method with exact line-search, which is actually FW with exact line-search. Method works remarkably well both in theory and practice: it computes an ε -optimal solution of (D-OPT) in (essentially) $O(n^2/\varepsilon)$ iterations.
- ▷ The theoretical success of this method has been a mysterious outlier for more than 20 years, since (D-OPT) does not satisfy the usual *L*-smooth curvature condition in (LSm). What problem structure actually drives the complexity bound? And might such structure exist anywhere else?

A Motivating Example: *D*-optimal Design

$$\begin{aligned} \max_{p} h(p) &\triangleq \ln \det \left(\sum_{i=1}^{m} p_{i} a_{i} a_{i}^{\top} \right) \\ \text{s.t.} \quad \sum_{i=1}^{m} p_{i} = 1, \ p_{i} \geq 0, \ \forall i \in [m]. \end{aligned} \tag{D-OPT}$$

- \triangleright Problem data: $\{a_i\}_{i=1}^m \subseteq \mathbb{R}^n$.
- ▷ Arises in many places, including optimal experimental design, and as the dual problem of the minimum volume enclosing ellipsoid (MVEE) problem.
- ▷ Khachiyan (1996) proposed a "barycentric coordinate ascent" method with exact line-search, which is actually FW with exact line-search. Method works remarkably well both in theory and practice: it computes an ε -optimal solution of (D-OPT) in (essentially) $O(n^2/\varepsilon)$ iterations.
- ▷ The theoretical success of this method has been a mysterious outlier for more than 20 years, since (D-OPT) does not satisfy the usual *L*-smooth curvature condition in (LSm). What problem structure actually drives the complexity bound? And might such structure exist anywhere else?
- \triangleright We resolve this mystery and generalize his method to the much broader class of problems in (P), even while relaxing the exact line-search requirement.

▷ Let an $m \times n$ matrix X denote the true representation of an image, such that $0 \le X_{ij} \le M$ denotes the pixel level at location (i, j).

- ▷ Let an $m \times n$ matrix X denote the true representation of an image, such that $0 \le X_{ij} \le M$ denotes the pixel level at location (i, j).
- \triangleright Let $A : \mathbb{R}^{m \times n} \to \mathbb{R}^{m \times n}$ denote the 2D discrete convolutional (linear) operator, which is assumed to be known.

- ▷ Let an $m \times n$ matrix X denote the true representation of an image, such that $0 \le X_{ij} \le M$ denotes the pixel level at location (i, j).
- \triangleright Let $A : \mathbb{R}^{m \times n} \to \mathbb{R}^{m \times n}$ denote the 2D discrete convolutional (linear) operator, which is assumed to be known.
- \vartriangleright The observed image Y is obtained by first passing X through A, and then is assumed to be subject to additive independent (entry-wise) Poisson noise.

- ▷ Let an $m \times n$ matrix X denote the true representation of an image, such that $0 \le X_{ij} \le M$ denotes the pixel level at location (i, j).
- \triangleright Let $A : \mathbb{R}^{m \times n} \to \mathbb{R}^{m \times n}$ denote the 2D discrete convolutional (linear) operator, which is assumed to be known.
- \vartriangleright The observed image Y is obtained by first passing X through A, and then is assumed to be subject to additive independent (entry-wise) Poisson noise.
- ▷ For convenience, we also represent A in its matrix form $A \in \mathbb{R}^{N \times N}$, where N := mn, and vectorize Y and X into $y \in \mathbb{R}^N$ and $x \in \mathbb{R}^N$, respectively. Notation: we write $x = \operatorname{vec}(X)$ and $X = \operatorname{mat}(x)$, etc.

Poisson Image Deblurring with TV Regularization, continued

Poisson Image Deblurring with TV Regularization, continued

 \triangleright We seek to recover X from Y (equivalently x from y) using maximum-likelihood estimation on the TV-regularized problem:

$$\begin{split} \min_{x \in \mathbb{R}^N} \ \bar{F}(x) &:= -\sum_{l=1}^N y_l \ln(a_l^\top x) + (\sum_{l=1}^N a_l)^\top x + \lambda \mathrm{TV}(x) \\ \text{s.t.} \ 0 &\leq x \leq Me \ , \end{split} \tag{Deblur}$$

Poisson Image Deblurring with TV Regularization, continued

 \triangleright We seek to recover X from Y (equivalently x from y) using maximum-likelihood estimation on the TV-regularized problem:

$$\begin{split} \min_{x \in \mathbb{R}^N} \ \bar{F}(x) &:= -\sum_{l=1}^N y_l \ln(a_l^\top x) + (\sum_{l=1}^N a_l)^\top x + \lambda \mathrm{TV}(x) \\ \text{s.t.} \ 0 &\leq x \leq Me \ , \end{split} \tag{Deblur}$$

▷ (Deblur) has a (standard) total-variation (TV) regularization term to recover a smooth image with sharp edges. The TV term is given by

$$\begin{split} \mathrm{TV}(x) &:= \sum_{i=1}^{m} \sum_{j=1}^{n-1} |[\mathsf{mat}(x)]_{i,j} - [\mathsf{mat}(x)]_{i,j+1}| \\ &+ \sum_{i=1}^{m-1} \sum_{j=1}^{n} |[\mathsf{mat}(x)]_{i,j} - [\mathsf{mat}(x)]_{i+1,j}| \end{split}$$

Some Other Applications

 \triangleright Positron emission tomography (PET)

 \triangleright Positron emission tomography (PET)

 \triangleright Optimal expected log investment (Cover (1984))

 \triangleright Positron emission tomography (PET)

 \triangleright Optimal expected log investment (Cover (1984))

 $\,\vartriangleright\,$ Computation of the analytic center of a polytope

$$F^* := \min_{x \in \mathbb{R}^n} [F(x) := f(Ax) + h(x)]$$
 (P)

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right] \tag{P}$$

11 / 20

▶ Initialize: $x^0 \in \text{dom } F, k := 0$

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

▶ Initialize: $x^0 \in \text{dom } F, k := 0$

▶ **Repeat** (until some convergence criterion is met)

 $v^k \in \arg\min_{x \in \mathbb{R}^n} \langle \nabla f(\mathsf{A}x^k), \mathsf{A}x \rangle + h(x)$ (Solve Lin. subproblem)

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

▶ Initialize: $x^0 \in \text{dom } F, k := 0$

$$\begin{split} v^{k} &\in \arg\min_{x \in \mathbb{R}^{n}} \langle \nabla f(\mathbf{A}x^{k}), \mathbf{A}x \rangle + h(x) \qquad \text{(Solve Lin. subproblem)} \\ G_{k} &:= \langle \nabla f(\mathbf{A}x^{k}), \mathbf{A}(x^{k} - v^{k}) \rangle + h(x^{k}) - h(v^{k}) \qquad \text{(FW Gap)} \end{split}$$

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

11 / 20

▶ Initialize: $x^0 \in \text{dom } F, k := 0$

$$v^{k} \in \arg\min_{x \in \mathbb{R}^{n}} \langle \nabla f(\mathsf{A}x^{k}), \mathsf{A}x \rangle + h(x) \qquad \text{(Solve Lin. subproblem)}$$
$$G_{k} := \langle \nabla f(\mathsf{A}x^{k}), \mathsf{A}(x^{k} - v^{k}) \rangle + h(x^{k}) - h(v^{k}) \qquad \text{(FW Gap)}$$
$$D_{k} := D_{k} := \|\mathsf{A}(v^{k} - x^{k})\|_{\mathsf{A}x^{k}} \qquad \text{(Local Distance)}$$

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

▶ Initialize: $x^0 \in \text{dom } F, k := 0$

$$v^{k} \in \arg\min_{x \in \mathbb{R}^{n}} \langle \nabla f(\mathsf{A}x^{k}), \mathsf{A}x \rangle + h(x) \qquad \text{(Solve Lin. subproblem)}$$

$$G_{k} := \langle \nabla f(\mathsf{A}x^{k}), \mathsf{A}(x^{k} - v^{k}) \rangle + h(x^{k}) - h(v^{k}) \qquad \text{(FW Gap)}$$

$$D_{k} := D_{k} := \|\mathsf{A}(v^{k} - x^{k})\|_{\mathsf{A}x^{k}} \qquad \text{(Local Distance)}$$

$$\alpha_{k} := \min\left\{\frac{G_{k}}{D_{k}(G_{k} + D_{k})}, 1\right\} \qquad \text{(Stepsize)}$$

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

▶ Initialize: $x^0 \in \text{dom } F, k := 0$

$$\begin{split} v^{k} &\in \arg\min_{x \in \mathbb{R}^{n}} \langle \nabla f(\mathbf{A}x^{k}), \mathbf{A}x \rangle + h(x) \qquad \text{(Solve Lin. subproblem)} \\ G_{k} &:= \langle \nabla f(\mathbf{A}x^{k}), \mathbf{A}(x^{k} - v^{k}) \rangle + h(x^{k}) - h(v^{k}) \qquad \text{(FW Gap)} \\ D_{k} &:= D_{k} &:= \|\mathbf{A}(v^{k} - x^{k})\|_{\mathbf{A}x^{k}} \qquad \text{(Local Distance)} \\ \alpha_{k} &:= \min\left\{ \frac{G_{k}}{D_{k}(G_{k} + D_{k})}, 1 \right\} \qquad \text{(Stepsize)} \\ x^{k+1} &:= x^{k} + \alpha_{k}(v^{k} - x^{k}) \qquad \text{(Update)} \end{split}$$

$$F^* := \min_{x \in \mathbb{R}^n} \left[F(x) := f(\mathsf{A}x) + h(x) \right]$$
(P)

▶ Initialize: $x^0 \in \text{dom } F, k := 0$

$$\begin{aligned} v^{k} &\in \arg\min_{x \in \mathbb{R}^{n}} \langle \nabla f(\mathsf{A}x^{k}), \mathsf{A}x \rangle + h(x) & \text{(Solve Lin. subproblem)} \\ G_{k} &:= \langle \nabla f(\mathsf{A}x^{k}), \mathsf{A}(x^{k} - v^{k}) \rangle + h(x^{k}) - h(v^{k}) & \text{(FW Gap)} \\ D_{k} &:= D_{k} &:= \|\mathsf{A}(v^{k} - x^{k})\|_{\mathsf{A}x^{k}} & \text{(Local Distance)} \\ \alpha_{k} &:= \min\left\{ \frac{G_{k}}{D_{k}(G_{k} + D_{k})}, 1 \right\} & \text{(Stepsize)} \\ x^{k+1} &:= x^{k} + \alpha_{k}(v^{k} - x^{k}) & \text{(Update)} \\ k &:= k + 1 \end{aligned}$$

▷ When h is the indicator function $h = \iota_{\mathcal{X}}$, then gFW-LHSCB specializes exactly to the algorithm of Dvurechensky et al. (2020).

- \triangleright When h is the indicator function $h = \iota_{\mathcal{X}}$, then gFW-LHSCB specializes exactly to the algorithm of Dvurechensky et al. (2020).
- \triangleright For most applications (including all of the applications mentioned previously), D_k in (Local Distance) can be computed in O(n) time.

- ▷ When h is the indicator function $h = \iota_{\mathcal{X}}$, then gFW-LHSCB specializes exactly to the algorithm of Dvurechensky et al. (2020).
- \triangleright For most applications (including all of the applications mentioned previously), D_k in (Local Distance) can be computed in O(n) time.
- ▷ The step-size rule in (Stepsize) is derived from the "curvature property" of a (standard) self-concordant function:

$$f(x^k + \alpha(v^k - x^k)) \le f(x^k) - \alpha G_k + \omega(\alpha D_k),$$
 (Curvature)

where $\omega(t) := -t - \ln(1-t)$ for t < 1.

- \triangleright When h is the indicator function $h = \iota_{\mathcal{X}}$, then gFW-LHSCB specializes exactly to the algorithm of Dvurechensky et al. (2020).
- \triangleright For most applications (including all of the applications mentioned previously), D_k in (Local Distance) can be computed in O(n) time.
- ▷ The step-size rule in (Stepsize) is derived from the "curvature property" of a (standard) self-concordant function:

$$f(x^k + \alpha(v^k - x^k)) \le f(x^k) - \alpha G_k + \omega(\alpha D_k),$$
 (Curvature)

where $\omega(t) := -t - \ln(1-t)$ for t < 1.

 \triangleright Neither the algorithm nor (Curvature) use the special properties of the barrier or the logarithmic homogeneity of f. However, these properties drive our complexity analysis.

Computational Guarantees

Define $\delta_k := F(x^k) - F^*$ for $k \ge 0$ (hence δ_0 is the initial optimality gap) Define $R_h := \max_{x,y \in \text{dom } h} |h(x) - h(y)|$ (the variation of h on its domain)

Theorem:

Computational Guarantees

Define $\delta_k := F(x^k) - F^*$ for $k \ge 0$ (hence δ_0 is the initial optimality gap)

Define $R_h := \max_{x,y \in \text{dom } h} |h(x) - h(y)|$ (the variation of h on its domain)

Theorem:

 $[\text{Iteration complexity for } \varepsilon \text{-optimality gap}) \text{ Let } K_{\varepsilon} \text{ denote the number of iterations required by gFW-LHSCB to obtain } \delta_k \leq \varepsilon. \text{ Then:} \\ K_{\varepsilon} \leq \left\lceil 5.3(\delta_0 + \theta + R_h) \ln(10.6\delta_0) \right\rceil + \left\lceil 12(\theta + R_h)^2 \max\left\{\frac{1}{\varepsilon} - \frac{1}{\delta_0}, 0\right\} \right\rceil .$

Computational Guarantees

Define $\delta_k := F(x^k) - F^*$ for $k \ge 0$ (hence δ_0 is the initial optimality gap)

Define $R_h := \max_{x,y \in \text{dom } h} |h(x) - h(y)|$ (the variation of h on its domain)

Theorem:

▷ (Iteration complexity for ε -optimality gap) Let K_{ε} denote the number of iterations required by gFW-LHSCB to obtain $\delta_k \leq \varepsilon$. Then:

$$K_{\varepsilon} \leq \left\lceil 5.3(\delta_0 + \theta + R_h) \ln(10.6\delta_0) \right\rceil + \left| 12(\theta + R_h)^2 \max\left\{ \frac{1}{\varepsilon} - \frac{1}{\delta_0}, 0 \right\} \right| .$$

 \triangleright (Iteration complexity for ε -FW gap) Let FWGAP $_{\varepsilon}$ denote the number of iterations required by gFW-LHSCB to obtain $G_k \leq \varepsilon$. Then:

$$FWGAP_{\varepsilon} \leq \left\lceil 5.3(\delta_0 + \theta + R_h) \ln(10.6\delta_0) \right\rceil + \left\lceil \frac{24(\theta + R_h)^2}{\varepsilon} \right\rceil$$

 \triangleright Our computational guarantees only depend on three (natural) quantities:

 \triangleright Our computational guarantees only depend on three (natural) quantities:

• the initial optimality gap δ_0 ,

- \triangleright Our computational guarantees only depend on three (natural) quantities:
 - the initial optimality gap δ_0 ,
 - the complexity parameter θ of the barrier f,

▷ Our computational guarantees only depend on three (natural) quantities:

- the initial optimality gap δ_0 ,
- the complexity parameter θ of the barrier f,
- the variation of h on its domain dom $h (= 0 \text{ if } h = \iota_{\mathcal{X}})$.
- ▷ Our computational guarantees only depend on three (natural) quantities:
 - the initial optimality gap δ_0 ,
 - the complexity parameter θ of the barrier f,
 - the variation of h on its domain dom $h (= 0 \text{ if } h = \iota_{\mathcal{X}})$.
- \triangleright Comparison with Khachiyan's results for (D-OPT):

- ▷ Our computational guarantees only depend on three (natural) quantities:
 - the initial optimality gap δ_0 ,
 - the complexity parameter θ of the barrier f,
 - the variation of h on its domain dom $h (= 0 \text{ if } h = \iota_{\mathcal{X}})$.
- \triangleright Comparison with Khachiyan's results for (D-OPT):
 - In (D-OPT), we have $\theta = n$, $R_h = 0$, and if $x^0 = (1/m)e$, then $\delta_0 \leq n \ln(m/n)$.

- \triangleright Our computational guarantees only depend on three (natural) quantities:
 - the initial optimality gap δ_0 ,
 - the complexity parameter θ of the barrier f,
 - the variation of h on its domain dom $h (= 0 \text{ if } h = \iota_{\mathcal{X}})$.
- \triangleright Comparison with Khachiyan's results for (D-OPT):
 - In (D-OPT), we have $\theta = n$, $R_h = 0$, and if $x^0 = (1/m)e$, then $\delta_0 \leq n \ln(m/n)$.
 - Using the adaptive step-size, our complexity bound specializes to $O\left(n\ln(m/n)(\ln n + \ln\ln(m/n)) + n^2/\varepsilon\right). \tag{Ours}$

- \triangleright Our computational guarantees only depend on three (natural) quantities:
 - the initial optimality gap δ_0 ,
 - the complexity parameter θ of the barrier f,
 - the variation of h on its domain dom $h (= 0 \text{ if } h = \iota_{\mathcal{X}})$.
- \triangleright Comparison with Khachiyan's results for (D-OPT):
 - In (D-OPT), we have $\theta = n$, $R_h = 0$, and if $x^0 = (1/m)e$, then $\delta_0 \leq n \ln(m/n)$.
 - Using the adaptive step-size, our complexity bound specializes to $O\left(n\ln(m/n)(\ln n + \ln\ln(m/n)) + n^2/\varepsilon\right).$
 - Using exact line-search, Khachiyan's bound is

$$O\left(n(\ln n + \ln \ln(m/n)) + n^2/\varepsilon\right)$$
 . (Kha)

(Ours)

- \triangleright Our computational guarantees only depend on three (natural) quantities:
 - the initial optimality gap δ_0 ,
 - the complexity parameter θ of the barrier f,
 - the variation of h on its domain dom $h (= 0 \text{ if } h = \iota_{\mathcal{X}})$.
- \triangleright Comparison with Khachiyan's results for (D-OPT):
 - In (D-OPT), we have $\theta = n$, $R_h = 0$, and if $x^0 = (1/m)e$, then $\delta_0 \leq n \ln(m/n)$.
 - Using the adaptive step-size, our complexity bound specializes to $O\left(n\ln(m/n)(\ln n + \ln\ln(m/n)) + n^2/\varepsilon\right).$
 - Using exact line-search, Khachiyan's bound is

$$Oig(n(\ln n + \ln\ln(m/n)) + n^2/arepsilonig)$$
 . (Kha)

• Observe that (Ours) has the exact same dependence on ε as (Kha), namely $O(n^2/\varepsilon)$, but the "fixed" term is slightly inferior to (Kha) by the factor $O(\ln(m/n))$.

(Ours)

$$\begin{split} \min_{x \in \mathbb{R}^N} \ \bar{F}(x) &:= \underbrace{-\sum_{l=1}^N y_l \ln(a_l^\top x)}_{=f(\mathsf{A}x)} + \underbrace{\langle \sum_{l=1}^N a_l, x \rangle + \lambda \mathrm{TV}(x)}_{=h(x)} \\ \text{s.t.} \ 0 \le x \le Me \ , \end{split}$$
(Deblur)

15 / 20

$$\begin{split} \min_{x \in \mathbb{R}^N} \ \bar{F}(x) &:= \underbrace{-\sum_{l=1}^N y_l \ln(a_l^\top x)}_{=f(\mathsf{A}x)} + \underbrace{\langle \sum_{l=1}^N a_l, x \rangle + \lambda \mathrm{TV}(x)}_{=h(x)} \\ \text{s.t.} \ 0 \le x \le Me \ , \end{split} \tag{Deblur}$$

Very few principled first-order methods have been proposed to solve (Deblur), because:

$$\begin{split} \min_{x \in \mathbb{R}^N} \ \bar{F}(x) &:= \underbrace{-\sum_{l=1}^N y_l \ln(a_l^\top x)}_{=f(\mathsf{A}x)} + \underbrace{\langle \sum_{l=1}^N a_l, x \rangle + \lambda \mathrm{TV}(x)}_{=h(x)} \\ \text{s.t.} \ 0 \le x \le Me \ , \end{split} \tag{Deblur}$$

- Very few principled first-order methods have been proposed to solve (Deblur), because:
 - $f: u \mapsto -\sum_{l=1}^{N} y_l \ln(u_l)$ is neither Lipschitz nor *L*-smooth on the set $\{u \in \mathbb{R}^N : u = \mathsf{A}x, \ 0 \le x \le Me\}$, and

$$\begin{split} \min_{x \in \mathbb{R}^N} \ \bar{F}(x) &:= \underbrace{-\sum_{l=1}^N y_l \ln(a_l^\top x)}_{=f(\mathsf{A}x)} + \underbrace{\langle \sum_{l=1}^N a_l, x \rangle + \lambda \mathrm{TV}(x)}_{=h(x)} \\ \text{s.t.} \ 0 \le x \le Me \ , \end{split} \tag{Deblur}$$

- Very few principled first-order methods have been proposed to solve (Deblur), because:
 - $f: u \mapsto -\sum_{l=1}^{N} y_l \ln(u_l)$ is neither Lipschitz nor *L*-smooth on the set $\{u \in \mathbb{R}^N : u = \mathsf{A}x, \ 0 \le x \le Me\}$, and
 - $TV(\cdot)$ does not have an efficiently computable proximal operator.

15 / 20

$$\begin{split} \min_{x \in \mathbb{R}^N} \ \bar{F}(x) &:= \underbrace{-\sum_{l=1}^N y_l \ln(a_l^\top x)}_{=f(\mathsf{A}x)} + \underbrace{\langle \sum_{l=1}^N a_l, x \rangle + \lambda \mathrm{TV}(x)}_{=h(x)} \\ \text{s.t.} \ 0 \le x \le Me \ , \end{split} \tag{Deblur}$$

- Very few principled first-order methods have been proposed to solve (Deblur), because:
 - $f: u \mapsto -\sum_{l=1}^{N} y_l \ln(u_l)$ is neither Lipschitz nor *L*-smooth on the set $\{u \in \mathbb{R}^N : u = \mathsf{A}x, \ 0 \le x \le Me\}$, and
 - $TV(\cdot)$ does not have an efficiently computable proximal operator.
- \vartriangleright However, $\mathrm{TV}(\cdot)$ is a polyhedral function, and the linear-optimization sub-problem

 $v^k \in \arg\min_{0 \le x \le Me} \langle \nabla f(\mathsf{A}x^k), \mathsf{A}x \rangle + \langle \sum_{l=1}^N a_l, x \rangle + \lambda \mathrm{TV}(x)$

can be formulated as a relatively simple LP and solved easily using a standard LP solver such as Gurobi.

Implementation Details/Issues

16 / 20

Implementation Details/Issues

 \triangleright We evaluate the numerical performance of our FW method gFW-LHSCB (with adaptive stepsize) which we call FW-Adapt.

- \rhd We evaluate the numerical performance of our FW method gFW-LHSCB (with adaptive stepsize) which we call FW-Adapt.
- \rhd It turns out that an exact line-search step-size for gFW-LHSCB can be computed for this particular problem, which we call FW-Exact.

- \rhd We evaluate the numerical performance of our FW method gFW-LHSCB (with adaptive stepsize) which we call FW-Adapt.
- \rhd It turns out that an exact line-search step-size for gFW-LHSCB can be computed for this particular problem, which we call FW-Exact.
- ▷ We tested FW-Adapt and FW-Exact on the Shepp-Logan phantom image of size 100×100 (hence N = 10,000).

- \rhd We evaluate the numerical performance of our FW method gFW-LHSCB (with adaptive stepsize) which we call FW-Adapt.
- \triangleright It turns out that an exact line-search step-size for gFW-LHSCB can be computed for this particular problem, which we call FW-Exact.
- \triangleright We tested FW-Adapt and FW-Exact on the Shepp-Logan phantom image of size 100×100 (hence N = 10,000).
- \triangleright We chose the starting point $x^0 = \text{vec}(Y)$, and we set $\lambda = 0.01$.

- \rhd We evaluate the numerical performance of our FW method gFW-LHSCB (with adaptive stepsize) which we call FW-Adapt.
- \rhd It turns out that an exact line-search step-size for gFW-LHSCB can be computed for this particular problem, which we call FW-Exact.
- \triangleright We tested FW-Adapt and FW-Exact on the Shepp-Logan phantom image of size 100×100 (hence N = 10,000).
- \triangleright We chose the starting point $x^0 = \text{vec}(Y)$, and we set $\lambda = 0.01$.
- \triangleright We used CVXPY to (approximately) compute the optimal objective value \bar{F}^* of (**Deblur**) in order to compute optimality gaps.

Results: Recovered Images



Figure 1:

True, noisy and recovered Shepp-Logan phantom image.

Results: Optimality Gaps versus Time and Iterations



(a) Optimality gap versus time (in seconds)

(b) Optimality gap versus iterations

Figure 2:

Comparison of empirical optimality gaps of FW-Adapt (FW-A) and FW-Exact (FW-E) for image recovery of the Shepp-Logan phantom image.

Thank you!

 $\,\triangleright\,$ To find an $\varepsilon\text{-optimal solution, the complexity bound in Dvurechensky et al. (2020) reads:$

$$O\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\ln\left(\delta_0/\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\right)\right) + L(x^0)D_{\mathcal{X},\|\cdot\|_2}^2/\varepsilon\right), \quad (\mathsf{Dvu})$$

where $\mathcal{S}(x^0):=\{x\in \mathsf{dom}\, F\cap \mathcal{X}\,:\, F(x)\leq F(x^0)\}$ denotes the initial level-set and

$$L(x^{0}) := \max_{x \in \mathcal{S}(x^{0})} \|\nabla^{2} \bar{F}(x)\|_{2} < +\infty \text{ , and } D_{\mathcal{X}, \|\cdot\|_{2}} := \max_{x, y \in \mathcal{X}} \|x - y\|_{2}$$

 $\,\vartriangleright\,$ To find an $\varepsilon\text{-optimal solution, the complexity bound in Dvurechensky et al. (2020) reads:$

$$O\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\ln\left(\delta_0/\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\right)\right) + L(x^0)D_{\mathcal{X},\|\cdot\|_2}^2/\varepsilon\right), \quad (\mathsf{Dvu})$$

where $\mathcal{S}(x^0) := \{x \in \mathsf{dom} \ F \cap \mathcal{X} \ : \ F(x) \le F(x^0)\}$ denotes the initial level-set and

$$L(x^{0}) := \max_{x \in \mathcal{S}(x^{0})} \|\nabla^{2} \bar{F}(x)\|_{2} < +\infty \text{ , and } D_{\mathcal{X}, \|\cdot\|_{2}} := \max_{x, y \in \mathcal{X}} \|x - y\|_{2}$$

▷ Specialized to the traditional setting, our complexity bound reads:

$$O((\delta_0 + \theta) \ln(\delta_0) + (\theta)^2 / \varepsilon).$$
 (Ours)

 $\,\triangleright\,$ To find an $\varepsilon\text{-optimal solution, the complexity bound in Dvurechensky et al. (2020) reads:$

$$O\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\ln\left(\delta_0/\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\right)\right) + L(x^0)D_{\mathcal{X},\|\cdot\|_2}^2/\varepsilon\right), \quad (\mathsf{Dvu})$$

where $S(x^0) := \{x \in \mathsf{dom} \ F \cap \mathcal{X} : F(x) \le F(x^0)\}$ denotes the initial level-set and

$$L(x^{0}) := \max_{x \in \mathcal{S}(x^{0})} \|\nabla^{2} \bar{F}(x)\|_{2} < +\infty \text{ , and } D_{\mathcal{X}, \|\cdot\|_{2}} := \max_{x, y \in \mathcal{X}} \|x - y\|_{2}$$

▷ Specialized to the traditional setting, our complexity bound reads:

$$O((\delta_0 + \theta) \ln(\delta_0) + (\theta)^2 / \varepsilon).$$
 (Ours)

 $\,\vartriangleright\,$ To find an $\varepsilon\text{-optimal solution, the complexity bound in Dvurechensky et al. (2020) reads:$

$$O\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\ln\left(\delta_0/\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\right)\right) + L(x^0)D_{\mathcal{X},\|\cdot\|_2}^2/\varepsilon\right), \quad (\mathsf{Dvu})$$

where $S(x^0) := \{x \in \mathsf{dom} \ F \cap \mathcal{X} : F(x) \le F(x^0)\}$ denotes the initial level-set and

$$L(x^{0}) := \max_{x \in \mathcal{S}(x^{0})} \|\nabla^{2} \bar{F}(x)\|_{2} < +\infty \text{ , and } D_{\mathcal{X}, \|\cdot\|_{2}} := \max_{x, y \in \mathcal{X}} \|x - y\|_{2}$$

▷ Specialized to the traditional setting, our complexity bound reads:

$$O((\delta_0 + \theta) \ln(\delta_0) + (\theta)^2 / \varepsilon).$$
 (Ours)

 \triangleright Our bound (**Ours**) has the following merits:

• Affine-invariance

 $\,\vartriangleright\,$ To find an $\varepsilon\text{-optimal solution, the complexity bound in Dvurechensky et al. (2020) reads:$

$$O\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\ln\left(\delta_0/\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\right)\right) + L(x^0)D_{\mathcal{X},\|\cdot\|_2}^2/\varepsilon\right), \quad (\mathsf{Dvu})$$

where $S(x^0) := \{x \in \mathsf{dom} \ F \cap \mathcal{X} : F(x) \le F(x^0)\}$ denotes the initial level-set and

$$L(x^{0}) := \max_{x \in \mathcal{S}(x^{0})} \|\nabla^{2} \bar{F}(x)\|_{2} < +\infty \text{ , and } D_{\mathcal{X}, \|\cdot\|_{2}} := \max_{x, y \in \mathcal{X}} \|x - y\|_{2}$$

▷ Specialized to the traditional setting, our complexity bound reads:

$$O((\delta_0 + \theta) \ln(\delta_0) + (\theta)^2 / \varepsilon).$$
 (Ours)

- Affine-invariance
- Norm-invariance

 $\,\vartriangleright\,$ To find an $\varepsilon\text{-optimal solution, the complexity bound in Dvurechensky et al. (2020) reads:$

$$O\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\ln\left(\delta_0/\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\right)\right) + L(x^0)D_{\mathcal{X},\|\cdot\|_2}^2/\varepsilon\right), \quad (\mathsf{Dvu})$$

where $S(x^0) := \{x \in \mathsf{dom} \ F \cap \mathcal{X} : F(x) \le F(x^0)\}$ denotes the initial level-set and

$$L(x^{0}) := \max_{x \in \mathcal{S}(x^{0})} \|\nabla^{2} \bar{F}(x)\|_{2} < +\infty \text{ , and } D_{\mathcal{X}, \|\cdot\|_{2}} := \max_{x, y \in \mathcal{X}} \|x - y\|_{2}$$

▷ Specialized to the traditional setting, our complexity bound reads:

$$O((\delta_0 + \theta) \ln(\delta_0) + (\theta)^2 / \varepsilon).$$
 (Ours)

- Affine-invariance
- Norm-invariance
- Interpretability

 $\,\vartriangleright\,$ To find an $\varepsilon\text{-optimal solution, the complexity bound in Dvurechensky et al. (2020) reads:$

$$O\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\ln\left(\delta_0/\left(\sqrt{L(x^0)}D_{\mathcal{X},\|\cdot\|_2}\right)\right) + L(x^0)D_{\mathcal{X},\|\cdot\|_2}^2/\varepsilon\right), \qquad (\mathsf{Dvu})$$

where $S(x^0) := \{x \in \mathsf{dom} \ F \cap \mathcal{X} : F(x) \le F(x^0)\}$ denotes the initial level-set and

$$L(x^{0}) := \max_{x \in \mathcal{S}(x^{0})} \|\nabla^{2} \bar{F}(x)\|_{2} < +\infty \text{ , and } D_{\mathcal{X}, \|\cdot\|_{2}} := \max_{x, y \in \mathcal{X}} \|x - y\|_{2}$$

▷ Specialized to the traditional setting, our complexity bound reads:

$$O((\delta_0 + \theta) \ln(\delta_0) + (\theta)^2 / \varepsilon).$$
 (Ours)

- Affine-invariance
- Norm-invariance
- Interpretability
- Ease of parameter estimation