1   We sincerely thank all the reviewers for their detailed comments and queries, and give clarifications and answers below.

2   **Reviewer 1:** We will revise according to all the comments on typos, clarity, and rigor of the mathematical writing.
3   First, we focus on the correctness of Lemma 1, and then address the others comments.
4   **Correctness.** We agree with the reviewer that Lemma 1 is not precise, as we (mistakenly) did not include the precise
5   range of $x$ in the statements. We provide here the precise version of Lem.1, where the differences are colored in BLUE:
6   **Lemma 1.** *For a closed convex set $\mathbb{X}$, a convex proper l.s.c. function $f : \mathbb{X} \to \mathbb{R} \cup \{\infty\}$ and $\lambda > 0$ define $f_\lambda : \mathbb{R}^d \to \mathbb{R}$*
7   *as $f_\lambda(x) := \min_{x' \in \mathbb{X}} f(x') + \frac{1}{2\lambda}\|x - x'\|^2$ and $\hat{x}_\lambda(x) := \operatorname{argmin}_{x' \in \mathbb{X}} f(x') + \frac{1}{2\lambda}\|x - x'\|^2$. Then for any $x \in \mathbb{X}$:*
8   *(a) $\hat{x}_\lambda(x)$ is unique and $f(\hat{x}_\lambda(x)) \leq f_\lambda(x) \leq f(x)$.*
9   *(b) $f_\lambda$ is convex, differentiable, $1/\lambda$-smooth and $\nabla f_\lambda(x) = (1/\lambda)(x - \hat{x}_\lambda(x))$, and,*
10   *(c) if $f$ is $G$-Lipschitz continuous, then, $\|\hat{x}_\lambda(x) - x\| \leq G\lambda$, and $f(x) \leq f_\lambda(x) + G^2\lambda/2$.*

11   This version of Lemma 1 is $(i)$ sufficient for proving our main results (lines 577, 622, 691, 705) and $(ii)$ correct. Since
12   the reviewer's counter example uses $x \notin \mathbb{X}$, it does not contradict Lemma 1(c). We now provide a full proof below.
13   *Proof (brief due to page limit).* Denote $\phi_{\lambda,x}(x') := f(x') + (1/2\lambda)\|x - x'\|^2$. Note that $\phi_{\lambda,x}(\cdot)$ is a $1/\lambda$-strongly
14   convex function and $f_\lambda(x) = \min_{x' \in \mathbb{X}} \phi_{\lambda,x}(x')$.
15   (a) Then $f(\hat{x}_\lambda(x)) \leq \phi_{\lambda,x}(\hat{x}_\lambda(x)) = \min_{x' \in \mathbb{X}} \phi_{\lambda,x}(x') = f_\lambda(x) \leq \phi_{\lambda,x}(x) = f(x)$ and the uniqueness of $\hat{x}_\lambda(x)$
16   follows from the strong convexity of $\phi_{\lambda,x}(\cdot)$ and the fact that $f$ is a proper convex function.
17   (b) Let $x \in \mathbb{R}^d$ and $g_x := (x - \hat{x}_\lambda(x))/\lambda$. By $1/\lambda$ strong convexity of $\phi_{\lambda,x}(x')$ and $\hat{x}_\lambda(x) = \operatorname{argmin}_{x' \in \mathbb{X}} \phi_{\lambda,x}(x')$,
18   we have for any $x' \in \mathbb{X}$ that $\phi_{\lambda,x}(x') \geq \phi_{\lambda,x}(\hat{x}_\lambda(x)) + \|x' - \hat{x}_\lambda(x)\|^2/2\lambda$, which simplifies to $f(x') \geq f(\hat{x}_\lambda(x)) +$
19   $\langle g_x, x' - \hat{x}_\lambda(x)\rangle$. Using this, for any $x, y \in \mathbb{R}^d$ we get

$$f_\lambda(y) - f_\lambda(x) = f(\hat{x}_\lambda(y)) - f(\hat{x}_\lambda(x)) + (\|\hat{x}_\lambda(y) - y\|^2 - \|\hat{x}_\lambda(x) - x\|^2)/2\lambda$$
$$\geq \langle g_x, \hat{x}_\lambda(y) - \hat{x}_\lambda(x)\rangle + \lambda/2(\|g_y\|^2 - \|g_x\|^2) = \langle g_x, y - x\rangle + \lambda/2\|g_x - g_y\|^2 \quad (1)$$

20   Instantiating the above for $y \leftarrow x, x \leftarrow y$ we also get $f_\lambda(y) - f_\lambda(x) \leq \langle g_y, y - x\rangle - \lambda/2\|g_x - g_y\|^2$. Combining
21   these two inequalities

$$0 \leq \lambda/2\|g_y - g_x\|^2 \leq f_\lambda(y) - f_\lambda(x) - \langle g_x, y - x\rangle \leq -\lambda/2\|g_y - g_x\|^2 + \langle g_y - g_x, y - x\rangle \leq \|y - x\|^2/2\lambda \quad (2)$$

22   This implies that $\lim_{y \to x}(f_\lambda(y) - f_\lambda(x) - \langle g_x, y - x\rangle)/\|y - x\| = 0$. Thus $f_\lambda$ is Frechet differentiable with gradient
23   $\nabla f_\lambda(x) = g_x = (x - \hat{x}_\lambda(x))/\lambda$. The above inequality also implies $f_\lambda$ is convex and $1/\lambda$-smooth.
24   (c) Let $x \in \mathbb{X}$. Using $1/\lambda$-strong convexity of $\phi_{\lambda,x}$ and $\hat{x}_\lambda(x) \in \operatorname{argmin}_{x' \in \mathbb{X}} \phi_{\lambda,x}(x')$, and $G$-Lipschitzness of $f$,

$$\|x - \hat{x}_\lambda(x)\|^2/2\lambda \leq \phi_{\lambda,x}(x) - \phi_{\lambda,x}(\hat{x}_\lambda(x)) = f(x) - f_\lambda(x) = f(x) - f(\hat{x}_\lambda(x)) - \|x - \hat{x}_\lambda(x)\|^2/2\lambda$$
$$\leq G\|\hat{x}_\lambda(x) - x\| - \|x - \hat{x}_\lambda(x)\|^2/2\lambda \leq G^2\lambda/2 . \quad \square$$

25   We say line 557: "for simplicity...$\mathbb{X}$ is the whole vector space". This was an assumption made, in the context of Sec. A.1,
26   for ease of exposition of the failed attempt at a PO efficient algorithm (Algo. APGD).
27   **Experimental verification.** As suggested, we compared the projection-free methods using a higher-dimensional
28   ($d = 50,176$) ImageNet dataset in the same low-rank SVM problem. For achieving an optimality gap of $0.02$,
29   Randomized-FW[52] used $34717/264$ FO/LMO calls and our MOLES used $4004/241$ FO/LMO calls. We will add
30   detailed simulation results including sensitivity analysis in the next revision. We agree that our algorithms have more
31   parameters and hence harder to tune than most baselines. Overcoming this is an important direction of future research.

32   **Reviewer 2** We agree that the reviewer's definition of the stochastic subgradient oracle is more appropriate. We
33   modified the manuscript according to the additional comments.

34   **Reviewer 3** The two properties we need of the superset $\mathbb{X} \supseteq \mathcal{X}$ are that (a) it is easy to project onto $\mathbb{X}$ and (b) $f$ is
35   $G$-Lipschitz on $\mathbb{X}$. In our paper, we choose $\mathbb{X}$ to be a Euclidean ball (which is easy to project to) but any other choice of
36   $\mathbb{X}$ which satisfies the above properties works just as well. One choice for this Euclidean ball is $B(x_0, D_\mathcal{X})$, where $x_0$ is
37   the initial point and $D_\mathcal{X}$ is the diameter of $\mathcal{X}$, instead of the ball of radius $2R$ we currently use.

38   As mentioned by R3, even if $f$ is $G$-Lipschitz inside the constraint $\mathcal{X}$, it could (i) blow up or (ii) be undefined just
39   outside of $\mathcal{X}$. Thus an $\mathbb{X}$ satisfying our requirements may not exist. In our experiments, we do not explicitly project onto
40   $\mathbb{X}$ (line 3.16) but still observed that $\|x_k - x'_k\| = O(G\lambda)$ and small, which implies that the iterates $x'_k$ are close to
41   $\mathcal{X}$. This hints that we may only need Lipschitzness over a much smaller set $\mathcal{X} + B(0, O(G\lambda))$, but we do know how to
42   prove this yet. Theoretically, we can work around the issue (ii) above by minimizing the convex extension $f_\mathcal{X} : \mathbb{R}^d \to \mathbb{R}$
43   of the function $f$ from the set $\mathcal{X}$, defined as $f_\mathcal{X}(x') := \max_{x \in \mathcal{X}} \max_{g \in \partial f(x)} f(x) + \langle g, x' - x\rangle$. The extension $f_\mathcal{X}$
44   has the same value as $f$ inside $\mathcal{X}$ and is $G$-Lipschitz everywhere. Therefore the following minimization problems
45   are equivalent: $\min_{x \in \mathcal{X}} f(x)$ and $\min_{x \in \mathcal{X}} f_\mathcal{X}(x')$. However, it is not clear if we can even estimate/approximate the
46   gradients of $f_\mathcal{X}$ efficiently. We could not find any relevant prior work and leave this question for future work. We
47   modified the manuscript according to the additional comments.