

ANALYSIS OF GENERALIZED BREGMAN SURROGATE ALGORITHMS FOR NONSMOOTH NONCONVEX STATISTICAL LEARNING

BY YIYUAN SHE, ZHIFENG WANG AND JIUWU JIN

Department of Statistics, Florida State University

Modern statistical applications often involve minimizing an objective function that may be nonsmooth and/or nonconvex. This paper focuses on a broad Bregman-surrogate algorithm framework including the local linear approximation, mirror descent, iterative thresholding, DC programming and many others as particular instances. The recharacterization via generalized Bregman functions enables us to construct suitable error measures and establish global convergence rates for nonconvex and nonsmooth objectives in possibly high dimensions. For sparse learning problems with a composite objective, under some regularity conditions, the obtained estimators as the surrogate's fixed points, though not necessarily local minimizers, enjoy provable statistical guarantees, and the sequence of iterates can be shown to approach the statistical truth within the desired accuracy geometrically fast. The paper also studies how to design adaptive momentum based accelerations without assuming convexity or smoothness by carefully controlling stepsize and relaxation parameters.

1. Introduction. Many statistical learning problems can be formulated as minimizing a certain objective function. In shrinkage estimation, the objective can often be represented as the sum of a loss function and a penalty function, neither of which is necessarily smooth or convex. For example, when the number of variables is much larger than the number of observations ($p \gg n$), sparsity-inducing penalties come into play and result in nondifferentiability. Furthermore, many popular penalties are nonconvex [22, 19, 65], making the computation and analysis more challenging. Although in low dimensions there are ways to tackle nonsmooth nonconvex optimization, statisticians often prefer easy-to-implement algorithms that scale well in big data applications. Therefore, first-order methods, gradient-descent type algorithms in particular, have recently attracted a great deal of attention due to their lower complexity per iteration and better numerical stability than Newton-type algorithms.

In this work, we study a class of algorithms in a *Bregman surrogate* framework. The idea is that instead of solving the original problem $\min_{\beta} f(\beta)$, one constructs a surrogate function

$$(1) \quad g(\beta; \beta^-) = f(\beta) + \Delta_{\psi}(\beta, \beta^-),$$

and generates a sequence of iterates according to

$$(2) \quad \beta^{(t+1)} \in \arg \min_{\beta} g(\beta; \beta^{(t)}).$$

The generalized Bregman function Δ_{ψ} will be rigorously defined in Section 2.1, and we will call g a (generalized) Bregman surrogate. Note that Δ_{ψ} is not necessarily the standard Bregman divergence [9] because we do not restrict ψ to be smooth or strictly convex or even convex. Bregman divergence does not seem to have been widely used in the statistics community, but see [64]. The generalized Bregman surrogate framework has a close connection

MSC2020 subject classifications: Primary 90C26, 49J52, 68Q25.

Keywords and phrases: nonconvex optimization, nonsmooth optimization, MM algorithms, Bregman divergence, statistical algorithmic analysis, momentum-based acceleration.

to the majorization-minimization (MM) principle [28, 29]. But the surrogate here as a function of β matches $f(\beta)$ to a higher order when β^- is set to β (cf. Lemma 4) and we do not always invoke the majorization condition $g(\beta; \beta^-) \geq f(\beta)$; the benefits will be seen in step size control and acceleration.

A variety of algorithms can be recharacterized by Bregman surrogates, including DC programming [55], local linear approximation (LLA) [67] and iterative thresholding [8, 47]. In contrast to the large body of literature in convex optimization, little research has been done on the rate of convergence of nonconvex optimization algorithms when $p > n$, and there is a lack of universal methodologies. Instead of proving local convergence results for some carefully chosen initial points, this work aims to establish *global* convergence rates regardless of the specific choice of the starting point, where a crucial element is the error measure. We will see that the most natural measures are unsurprisingly problem-dependent, but can be conveniently constructed via generalized Bregman functions.

Another perhaps more intriguing question to statisticians is how the statistical accuracy improves or deteriorates as the cycles progress, and whether the finally obtained estimators can enjoy provable guarantees in a statistical sense. See, for example, [1, 20, 63]; in particular, [36], one of the main motivations of our work, showed that for a composite objective composed of a loss and a regularizer that enforces sparsity, the sequence of iterates $\beta^{(t)}$ generated by gradient-descent type algorithms can approach a minimizer β^o at a linear rate even when $p > n$, if the problem under consideration satisfies some regularity conditions. This article reveals broader conclusions when using generalized Bregman surrogate algorithms in the composite setting: the more straightforward *statistical error* between the t -th iterate $\beta^{(t)}$ and the statistical truth β^* enjoys fast convergence, and the convergent fixed points, though not necessarily local minimizers, let alone global minimizers, possess the desired statistical accuracy in a minimax sense. The studies support the practice of avoiding unnecessary over-optimization in high-dimensional sparse learning tasks. Our theory will make heavy use of the calculus of generalized Bregman functions—in fact, the proofs become readily on hand with some nice properties of Δ established. Again, a wise choice of the discrepancy measure can facilitate theoretical analysis and lead to less restrictive regularity conditions.

Finally, we would like to study and extend Nesterov’s first and second accelerations [39, 40]. Accelerated gradient algorithms [4, 57, 32] have lately gained popularity in high-dimensional convex programming because they can attain the optimal rates of convergence among first-order methods. However, since convexity is indispensable to these theories, how to adapt the momentum techniques to nonsmooth nonconvex programming is largely unknown. Ghadimi and Lan [24] studied how to accelerate gradient descent type algorithms when the objective function is nonconvex but strongly smooth; the obtained convergence rate is of the same order as gradient descent for nonconvex problems. We are interested in more general Bregman surrogates with a possible lack of smoothness and convexity, most notably in high-dimensional nonconvex sparse learning. This work will come up with two momentum-based schemes to accelerate Bregman-surrogate algorithms by carefully controlling the sequences of relaxation parameters and step sizes.

Overall, this paper aims to provide a universal tool of generalized Bregman functions in the interplay between optimization and statistics, and to demonstrate its active roles in constructing error measures, formulating less restrictive regularity conditions, characterizing strong convexity, deriving the so-called basic inequalities in nonasymptotic statistical analysis, devising line search and momentum-based updates, and so on. The rest of this paper is organized as follows. In Section 2, we introduce the generalized Bregman surrogate framework and present some examples. Section 3 gives the main theoretical results on computational accuracy and statistical accuracy. Section 4 proposes and analyzes two acceleration schemes. We conclude in Section 5. Simulation studies and all technical details are provided in the Appendices.

Notation. Throughout the paper, we use C, c to denote positive constants. They are not necessarily the same at each occurrence. The class of continuously differentiable functions is denoted by \mathcal{C}^1 . Given any matrix \mathbf{A} , we denote its (i, j) -th element by A_{ij} . The spectral norm and the Frobenius norm of \mathbf{A} are denoted by $\|\mathbf{A}\|_2$ and $\|\mathbf{A}\|_F$, respectively. The Hadamard product of two matrices \mathbf{A} and \mathbf{B} of the same dimension is denoted by $\mathbf{A} \circ \mathbf{B}$ and their inner product is $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}\{\mathbf{A}^\top \mathbf{B}\}$. If $\mathbf{A} - \mathbf{B}$ is positive semi-definite, we also write $\mathbf{A} \succeq \mathbf{B}$. Let $[p] := \{1, \dots, p\}$. Given $\mathcal{J} \subset [p]$, we use $\mathbf{A}_{\mathcal{J}}$ to denote the submatrix of \mathbf{A} formed by the columns indexed by \mathcal{J} . Given a set $A \subset \mathbb{R}^n$, we use A° , $\text{ri}(A)$, \bar{A} to denote its interior, relative interior, and closure, respectively [45]. When f is an extended real-valued function from $D \subset \mathbb{R}^p$ to $\mathbb{R} \cup \{+\infty\}$, its effective domain is defined as $\text{dom}(f) = \{\beta \in \mathbb{R}^p : f(\beta) < +\infty\}$. Let $\mathbb{R}_+ = [0, +\infty)$.

2. Basics of generalized Bregman surrogates.

2.1. Generalized Bregman functions. Bregman divergence [9], typically defined for continuously differentiable and strictly convex functions, plays an important role in convex analysis. An extension of it based on “right-hand” Gateaux differentials helps to handle nonsmooth nonconvex optimization problems. We begin with one-sided directional derivative.

DEFINITION 1. Let $\psi : D \subset \mathbb{R}^p \rightarrow \mathbb{R}$ be a function. The one-sided directional derivative of ψ at $\beta \in D$ with increment \mathbf{h} is defined as

$$(3) \quad \delta\psi(\beta; \mathbf{h}) = \lim_{\epsilon \rightarrow 0^+} \frac{\psi(\beta + \epsilon \mathbf{h}) - \psi(\beta)}{\epsilon},$$

provided \mathbf{h} is admissible in the sense that $\beta + \epsilon \mathbf{h} \in D$ for sufficiently small $\epsilon : 0 < \epsilon < \epsilon_0$. When $\psi : D \rightarrow \mathbb{R}^n$ is a vector function, $\delta\psi$ is defined componentwise.

In the following, ψ is called (one-sided) directionally differentiable at β if $\delta\psi(\beta; \mathbf{h})$ as defined in (3) exists and is finite for all admissible \mathbf{h} , and if this holds for all $\beta \in D$, we say that ψ is directionally differentiable.

When $a > 0$, $\delta\psi(\beta; a\mathbf{h}) = a\delta\psi(\beta; \mathbf{h})$, but $\delta\psi$ is not necessarily a linear operator with respect to \mathbf{h} . Definition 1 is a relaxed version of the standard Gateaux differential which studies the limit when $\epsilon \rightarrow 0$. In high-dimensional sparse problems where nonsmooth regularizers and/or losses are widely used, (3) is more convenient and useful.

DEFINITION 2 (Generalized Bregman Function (GBF)). The generalized Bregman function associated with a function ψ is defined by

$$(4) \quad \Delta_\psi(\beta, \gamma) = \psi(\beta) - \psi(\gamma) - \delta\psi(\gamma; \beta - \gamma),$$

assuming $\beta, \gamma \in \text{dom}(\psi)$ and $\delta\psi(\gamma; \beta - \gamma)$ is meaningful and finite. In particular, when ψ is differentiable and strictly convex, the generalized Bregman function Δ_ψ becomes the standard Bregman divergence:

$$(5) \quad \mathbf{D}_\psi(\beta, \gamma) := \psi(\beta) - \psi(\gamma) - \langle \nabla\psi(\gamma), \beta - \gamma \rangle.$$

When ψ is a vector function, a vector version of Δ is defined componentwise.

When $\nabla\psi$ exists at β , $\delta\psi(\beta; \mathbf{h})$ reduces to $\langle \nabla\psi(\beta), \mathbf{h} \rangle$, which is linear in \mathbf{h} . So if ψ is the restriction of a function $\varphi \in \mathcal{C}^1$ to a convex set, $\Delta_\psi(\beta, \gamma) = \Delta_\varphi(\beta, \gamma)$ for all $\beta, \gamma \in \text{dom}(\psi)$. For simplicity, all functions in our paper are assumed to be defined on a whole vector space (\mathbb{R}^p , typically) unless otherwise mentioned, although most results can be formulated in the case of extended real-valued functions under the convexity of their effective domains.

The generalized Bregman $\Delta_\psi(\cdot, \gamma)$ can be seen as the difference between the function ψ and its radial approximations made at γ . A simple but important example is $\mathbf{D}_2(\beta, \gamma) := \mathbf{D}_{\|\cdot\|_2^2/2}(\beta, \gamma) = \|\beta - \gamma\|_2^2/2$. In general, Δ_ψ or \mathbf{D}_ψ may not be symmetric. The following symmetrized version turns out to be useful:

$$(6) \quad \bar{\Delta}_\psi(\beta, \gamma) := \frac{1}{2}(\Delta_\psi + \bar{\Delta}_\psi)(\beta, \gamma) = \frac{1}{2}\{\Delta_\psi(\beta, \gamma) + \Delta_\psi(\gamma, \beta)\},$$

where $\bar{\Delta}(\beta, \gamma)$ denotes $\Delta(\gamma, \beta)$. If ψ is smooth, $\bar{\Delta}_\psi(\beta, \gamma) = \langle \nabla \psi(\beta) - \nabla \psi(\gamma), \beta - \gamma \rangle$.

To simplify the notation, we use $\Delta_\psi \geq \Delta_\phi$ to denote $\Delta_\psi(\beta, \gamma) \geq \Delta_\phi(\beta, \gamma)$ for all β, γ , and so $\Delta_\psi \geq 0$ stands for $\Delta_\psi(\beta, \gamma) \geq 0, \forall \beta, \gamma$. Some basic properties of Δ are given as follows.

LEMMA 1. *Let ψ and φ be directionally differentiable functions. Then for any α, β, γ , we have the following properties.*

- (i) $\Delta_{a\psi+b\varphi}(\beta, \gamma) = a\Delta_\psi(\beta, \gamma) + b\Delta_\varphi(\beta, \gamma), \forall a, b \in \mathbb{R}$.
- (ii) If ψ is convex, it is directionally differentiable and $\Delta_\psi \geq 0$; conversely, if ψ is directionally differentiable and $\Delta_\psi \geq 0$ then ψ is convex.
- (iii) If $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable and $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^n$ is continuous and directionally differentiable, then $\Delta_{\psi \circ \varphi}(\beta, \gamma) = \Delta_\psi(\varphi(\beta), \varphi(\gamma)) + \langle \Delta_\varphi(\beta, \gamma), \nabla \psi(\varphi(\gamma)) \rangle$. Also, if $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$ is directionally differentiable and $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^n$ is linear, then $\Delta_{\psi \circ \varphi}(\beta, \gamma) = \Delta_\psi(\varphi(\beta), \varphi(\gamma))$.
- (iv) $\Delta_\psi(\beta, \gamma) = \int_0^1 [\delta\psi(\gamma + t(\beta - \gamma); \beta - \gamma) - \delta\psi(\gamma; \beta - \gamma)] dt$, provided $\delta\psi(\gamma + t(\beta - \gamma); \beta - \gamma)$ is integrable over $t \in [0, 1]$.

The properties will be frequently used in the rest of the paper. For instance, for $\psi = \rho\|\cdot\|_2^2/2 - f$, by (i) we can write $\Delta_\psi = \rho\mathbf{D}_2 - \Delta_f$. Sometimes, though f is not necessarily convex, $f + \nu\|\cdot\|_2^2/2$ is so for some $\nu \in \mathbb{R}$, which means $\Delta_f \geq -\nu\mathbf{D}_2$, owing to (ii). For $l(\beta) = l_0(\mathbf{X}\beta + \alpha)$, commonly encountered in statistical applications, (iii) states that $\Delta_l(\beta, \gamma) = \Delta_{l_0}(\mathbf{X}\beta + \alpha, \mathbf{X}\gamma + \alpha)$. For (iv), the integrability condition is met when the directional derivative restricted to the interval $[\beta, \gamma]$ is bounded by a constant (or more generally a Lebesgue integrable function); in particular, if ψ is L -strongly smooth, that is, $\nabla \psi$ exists and is Lipschitz continuous: $\|\nabla \psi(\beta) - \nabla \psi(\gamma)\|_* \leq L\|\beta - \gamma\|$ for any β, γ , where $\|\cdot\|_*$ is the dual norm of $\|\cdot\|$, $\Delta_\psi(\beta, \gamma) \leq L\|\beta - \gamma\|^2/2$ and for the Euclidean norm, $\Delta_\psi \leq L\mathbf{D}_2$ results.

Moreover, the GBF operator satisfies some interesting “idempotence” properties under some mild assumptions, which is extremely helpful in studying iterative optimization algorithms.

LEMMA 2. (i) When ψ is convex, $\Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) \leq \Delta_\psi(\beta, \gamma)$, and when ψ is concave, $\Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) \geq \Delta_\psi(\beta, \gamma)$ for all α, β, γ .

(ii) When ψ is directionally differentiable, for all $\alpha = (1 - \theta)\gamma + \theta\beta$ with $\theta \notin (0, 1)$, $\Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) = \Delta_\psi(\beta, \gamma)$ and in particular,

$$(7) \quad \Delta_{\Delta_\psi(\cdot, \beta)}(\beta, \gamma) = \Delta_{\Delta_\psi(\cdot, \gamma)}(\beta, \gamma) = \Delta_\psi(\beta, \gamma).$$

(iii) When $\delta\psi(\cdot; \beta - \gamma)$ is bounded in a neighborhood of α and has restricted radial continuity at α : $\lim_{\epsilon \rightarrow 0+} \delta\psi(\alpha + \epsilon\mathbf{h}; \beta - \gamma) = \delta\psi(\alpha; \beta - \gamma)$ for any $\mathbf{h} \in [\beta - \alpha, \gamma - \alpha]$, or when $\delta\psi(\alpha; \cdot)$ has restricted linearity $\delta\psi(\alpha; \mathbf{h}) = \langle g(\alpha), \mathbf{h} \rangle$ for some g and all $\mathbf{h} \in [\beta - \alpha, \gamma - \alpha]$, we have

$$(8) \quad \Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) = \Delta_\psi(\beta, \gamma).$$

In particular, (8) holds when ψ is differentiable at α or $\delta\psi(\cdot; \beta - \gamma)$ is continuous at α .

We refer to (ii) as the *weak idempotence* property and (iii) as the *strong idempotence* property. When Δ_ψ becomes a legitimate Bregman divergence, (8) can be rephrased into the three-point property $\mathbf{D}_\psi(\beta, \gamma) = \mathbf{D}_\psi(\beta, \alpha) + \mathbf{D}_\psi(\alpha, \gamma) - \langle \beta - \alpha, \nabla \psi(\gamma) - \nabla \psi(\alpha) \rangle$ [14]. It is worth mentioning that although from (iii), differentiability can be used to gain strong idempotence, the weak idempotence (7) is often what we need, which always holds under just directional differentiability.

At the end of the subsection, we give some important facts of GBFs for canonical generalized linear models (GLMs) that are widely used in statistics modeling. Here, the response variable $\mathbf{y} \in \mathcal{Y}^n \subset \mathbb{R}^n$ has density $p_\eta(\cdot) = \exp\{(\langle \cdot, \eta \rangle - b(\eta))/\sigma^2 - c(\cdot, \sigma^2)\}$ with respect to measure ν_0 defined on \mathcal{Y}^n (typically the counting measure or Lebesgue measure), where $\eta \in \mathbb{R}^n$ represents the systematic component of interest, and σ is the scale parameter; see [30]. Since σ is not the parameter of interest, it is more convenient to define the density $\exp\{(\langle \cdot, \eta \rangle - b(\eta))/\sigma^2\}$ (still written as $p_\eta(\cdot)$ with a slight abuse of notation) with respect to the base measure $d\nu = \exp(-c(\cdot, \sigma^2)) d\nu_0$. The loss for η can be written as

$$(9) \quad l_0(\eta; \mathbf{y}) = \{-\langle \mathbf{y}, \eta \rangle + b(\eta)\}/\sigma^2.$$

That is, l_0 corresponds to a distribution in the exponential dispersion family with cumulant function $b(\cdot)$, dispersion σ^2 and natural parameter η . In the Gaussian case, $l_0(\eta) = -\langle \eta, \mathbf{y} \rangle/\sigma^2 + \|\mathbf{y}\|_2^2/(2\sigma^2)$.

Following [62], we define the natural parameter space $\Omega = \text{dom}(b) = \{\eta \in \mathbb{R}^n : b(\eta) < \infty\}$ (always assumed to be nonempty) and the mean parameter space $\mathcal{M} = \{\mu \in \mathbb{R}^n : \mu = \mathbb{E}\mathbf{y}, \text{ where } \mathbf{y} \sim p \text{ for some density } p \text{ defined on } \mathcal{Y}^n \text{ with respect to } \nu\}$, and call p_η minimal if $\langle \mathbf{a}, \mathbf{z} \rangle = c$ for almost every $\mathbf{z} \in \mathcal{Y}^n$ with respect to ν implies $\mathbf{a} = \mathbf{0}$. When Ω is open, p_η is called regular, and b can be shown to be differentiable to any order and convex, but not necessarily strictly convex; if, in addition, p_η is minimal, b is strictly convex and the canonical link $g = (\nabla b)^{-1}$ is well-defined on \mathcal{M}° . These can all be derived from, say, the propositions in [62].

LEMMA 3. *Assume the exponential dispersion family setup with the associated loss defined in (9). (i) If Ω is an open set or p_η is regular, then*

$$(10) \quad l_0(\eta; \mathbf{z}) = \Delta_b(\eta, \partial b^*(\mathbf{z}))/\sigma^2 - b^*(\mathbf{z})/\sigma^2$$

for all $\eta \in \Omega, \mathbf{z} \in \text{ri}(\mathcal{M})$, where b^* is the Fenchel conjugate of b , and $\partial b^*(\mathbf{z})$ can take any subgradient of b^* at \mathbf{z} . If p_η is also minimal, Δ_b becomes \mathbf{D}_b , $\partial b^*(\mathbf{z})$ becomes $g(\mathbf{z})$ (which is unique), and $\text{ri}(\mathcal{M})$ becomes \mathcal{M}° . (ii) As long as Ω is open,

$$(11) \quad l_0(\eta; \mathbf{z}) = \Delta_{b^*}(\mathbf{z}, \nabla b(\eta))/\sigma^2 - b^*(\mathbf{z})/\sigma^2$$

for all $\eta \in \Omega, \mathbf{z} \in \text{ri}(\mathcal{M})$. If p_η is also minimal, $\Delta_{b^*} = \mathbf{D}_{b^*}$ and $\text{ri}(\mathcal{M}) = \mathcal{M}^\circ$. (iii) Given any $\eta_1 \in \Omega^\circ$ and $\eta_2 \in \Omega$, the Kullback Leibler (KL) divergence of p_{η_2} from p_{η_1} relates to the GBF of l_0 or b by

$$(12) \quad KL(p_{\eta_1}, p_{\eta_2}) = \Delta_{l_0}(\eta_2, \eta_1) = \Delta_b(\eta_2, \eta_1)/\sigma^2.$$

Property (i) shows the importance of GBF in maximum likelihood estimation. A Bregman version of Property (ii) was first described in [3], while our conclusions based on Δ_b, Δ_{b^*} are more general, as they do *not* require the strict convexity of b or the differentiability of b^* . Consider for instance the multinomial GLM under a symmetric parametrization: for $[y_1, \dots, y_m] \in \mathcal{Y} = \{y_k \in \{0, 1\}, 1 \leq k \leq m, \sum y_k = 1\}$ ($n = 1$), $\mathbb{E}y_k \propto \exp(\eta_k)$ or $\mathbb{E}y_k = \exp(\eta_k) / \sum \exp(\eta_k)$ gives $b = \log \sum \exp(\eta_k)$, and thus $b^*(\mu)$ takes $\sum \mu_k \log \mu_k$ for

$[\mu_1, \dots, \mu_m] \in \mathcal{M} = \{[\mu_k] : \sum \mu_k = 1, \mu_k \geq 0\}$ and $+\infty$ otherwise. Clearly, b^* is not differentiable (given any $\mathbf{z} \in \text{ri}(\mathcal{M})$, $\partial b^*(\mathbf{z}) = \{\log \mathbf{z} + t\mathbf{1} : t \in \mathbb{R}\}$), but nicely our two GBF representations still hold. In addition, if the right-hand side of (10) or (11), as a function of \mathbf{z} , is continuous on $\overline{\mathcal{M}}$, which is the case for Bernoulli, multinomial and Poisson, (i) and (ii) hold for any $\mathbf{z} \in \overline{\mathcal{M}}$ from [62, Theorem 3.4].

Property (iii) (notice the exchange of $\boldsymbol{\eta}_1$ and $\boldsymbol{\eta}_2$ in the generalized Bregman expressions) can be used to formulate and verify model regularity conditions in minimax studies of sparse GLMs, which are of great interest in high-dimensional statistical learning [58]. More concretely, consider a general signal class

$$(13) \quad \mathcal{B}(s^*, M) = \{\boldsymbol{\beta}^* \in \mathbb{R}^p : \|\boldsymbol{\beta}^*\|_0 \leq s^*, \|\boldsymbol{\beta}^*\|_\infty \leq M\},$$

where $s^* \leq p$, $0 \leq M \leq +\infty$. Some applications limit the magnitude of the coefficients β_j via a constraint or a penalty, resulting in a finite M . Let $I(\cdot)$ be any nondecreasing function with $I(0) = 0, I \not\equiv 0$. Some particular examples are $I(t) = t$ and $I(t) = 1_{t \geq c}$. Recall the regular exponential dispersion family with systematic component $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$ and loss $l(\boldsymbol{\beta}) = l_0(\boldsymbol{\eta})$ defined by (9).

THEOREM 1. *In the regular exponential dispersion family setup (with $\text{dom}(b)$ a nonempty open set), assume $p \geq 2, 1 \leq s^* \leq p/2$. Let*

$$(14) \quad P(s^*) = s^* \log(ep/s^*).$$

(i) *If*

$$(15) \quad \Delta_{l_0}(\mathbf{0}, \mathbf{X}\boldsymbol{\beta})\sigma^2 \leq \kappa \mathbf{D}_2(\mathbf{0}, \boldsymbol{\beta}), \quad \forall \boldsymbol{\beta} \in \mathcal{B}(s^*, M)$$

where $\kappa > 0$, there exist positive constants c, \tilde{c} , depending on $I(\cdot)$ only, such that

$$\inf_{\hat{\boldsymbol{\beta}}} \sup_{\boldsymbol{\beta}^* \in \mathcal{B}(s^*, M)} \mathbb{E}\{I(\mathbf{D}_2(\boldsymbol{\beta}^*, \hat{\boldsymbol{\beta}})/[\tilde{c} \min\{\sigma^2 P(s^*)/\kappa, M^2 s^*\}])\} \geq c > 0,$$

where $\hat{\boldsymbol{\beta}}$ denotes any estimator of $\boldsymbol{\beta}^*$.

(ii) *If*

$$(16) \quad \begin{cases} \underline{\kappa} \mathbf{D}_2(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \leq \mathbf{D}_2(\mathbf{X}\boldsymbol{\beta}_1, \mathbf{X}\boldsymbol{\beta}_2) \\ \Delta_{l_0}(\mathbf{0}, \mathbf{X}\boldsymbol{\beta}_1)\sigma^2 \leq \bar{\kappa} \mathbf{D}_2(\mathbf{0}, \boldsymbol{\beta}_1), \end{cases} \quad \forall \boldsymbol{\beta}_i \in \mathcal{B}(s^*, M)$$

where $\underline{\kappa}, \bar{\kappa} \geq 0$, then there exist positive constants c, \tilde{c} depending on $I(\cdot)$ only such that

$$\inf_{\hat{\boldsymbol{\beta}}} \sup_{\boldsymbol{\beta}^* \in \mathcal{B}(s^*, M)} \mathbb{E}\{I(\mathbf{D}_2(\mathbf{X}\boldsymbol{\beta}^*, \mathbf{X}\hat{\boldsymbol{\beta}})/[c \min\{(\underline{\kappa}/\bar{\kappa})\sigma^2 P(s^*), \underline{\kappa} M^2 s^*\}])\} \geq c > 0.$$

The GBF-form conditions (15), (16) can be viewed as an extension of restricted isometry [11], and are often easy to check using the Hessian. For example, from Lemma 1, we immediately know that if l_0 is L -strongly smooth, (15) is satisfied with $\kappa = L\|X\|_2^2$ even when $M = +\infty$. This is the case for regression and logistic regression, and accordingly, no estimation algorithms can beat the minimax rate $s^* \log(ep/s^*)$ (ignoring trivial factors). The optimal lower bounds provide useful guidance in establishing sharp statistical error upper bounds of Bregman-surrogate algorithms in Section 3.2.

2.2. Examples of Bregman surrogates.

EXAMPLE 1. (Gradient descent and mirror descent). Gradient descent is a simple first-order method to minimize a function $f \in \mathcal{C}^1$ which may be nonconvex. Starting with $\beta^{(0)}$, the algorithm proceeds as follows:

$$(17) \quad \beta^{(t+1)} = \beta^{(t)} - \alpha \nabla f(\beta^{(t)}),$$

where $\alpha > 0$ is a step size parameter. Its rationale can be seen by formulating a Bregman-surrogate algorithm using $\Delta_\psi = \rho \mathbf{D}_2 - \Delta_f$:

$$(18a) \quad \beta^{(t+1)} = \arg \min_{\beta} g(\beta; \beta^{(t)}) = f(\beta) + (\rho \mathbf{D}_2 - \Delta_f)(\beta, \beta^{(t)})$$

$$(18b) \quad = \beta^{(t)} - \frac{1}{\rho} \nabla f(\beta^{(t)}),$$

where $f(\cdot) - \Delta_f(\cdot, \beta^{(t)})$ gives a linear approximation of f and $1/\rho$ amounts to the step size. We call ρ the inverse step size parameter. (The generalized Bregman surrogate in (18a) extends the class of algorithms to a directionally differentiable f , with the update given by $\beta^{(t+1)} = \beta^{(t)} + (0 \vee -\delta f(\beta^{(t)}; \mathbf{h}^\circ)) \mathbf{h}^\circ / \rho$ and $\mathbf{h}^\circ \in \arg \max_{\|\mathbf{h}\|_2=1} [\delta f(\beta^{(t)}; \mathbf{h})]_-$, where $[\cdot]_-$ denotes the negative part ($t_- = (|t| - t)/2$)).

More generally, we can use a strictly convex $\varphi \in \mathcal{C}^1$ to construct

$$(19) \quad g(\beta; \beta^{(t)}) = f(\beta) + (\rho \mathbf{D}_\varphi - \Delta_f)(\beta, \beta^{(t)}),$$

Minimizing (19) with respect to β gives the renowned mirror descent [38]: $\beta^{(t+1)} = (\nabla \varphi)^{-1}(\nabla \varphi(\beta^{(t)}) - \nabla f(\beta^{(t)})/\rho)$, where $(\nabla \varphi)^{-1}$ is the inverse of $\nabla \varphi$. Mirror descent is widely used in convex programming, but this work does *not* restrict f to be convex.

EXAMPLE 2. (Iterative thresholding). Sparsity-inducing penalties are widely used in high-dimensional problems; see, for example, ℓ_0 , ℓ_1 [56], bridge penalties [22], SCAD [19], capped- ℓ_1 [66] and MCP [65]. There is a universal connection between thresholding rules and penalty functions [48], and the mapping from penalties to thresholdings is many-to-one. This makes it possible to apply an iterative thresholding algorithm to solve a general penalized problem of the form $\min_{\beta} l(\beta) + \sum_j P(\varrho \beta_j; \lambda)$ [8, 47]:

$$(20) \quad \beta^{(t+1)} = \Theta(\varrho \beta^{(t)} - \nabla l(\beta^{(t)})/\varrho; \lambda)/\varrho,$$

where Θ is a thresholding function inducing P , and $\varrho > 0$ is an algorithm parameter for the sake of scaling and convergence control. This class of iterative algorithms is called the *Thresholding-based Iterative Selection Procedures* (TISP) in [47] and is scalable in computation. For the rigorous definition of Θ and the Θ - P coupling formula, see Section 3.1 for detail. Some examples of Θ include: (i) soft-thresholding $\Theta_S(t; \lambda) = \text{sgn}(t)(|t| - \lambda)1_{|t| > \lambda}$, which induces the ℓ_1 penalty, (ii) hard-thresholding $\Theta_H(t; \lambda) = t1_{|t| > \lambda}$, which is associated with (infinitely) many penalties, with the capped- ℓ_1 penalty, (55), and the discrete ℓ_0 penalty as particular instances. The nonconvex SCAD and MCP penalties also have their corresponding thresholding rules. In this sense, thresholdings extend proximity operators. One can regard (20) as an outcome of minimizing the following Bregman surrogate

$$(21) \quad g(\beta; \beta^{(t)}) = l(\beta) + \sum P(\varrho \beta_j; \lambda) + (\varrho^2 \mathbf{D}_2 - \Delta_l)(\beta, \beta^{(t)}).$$

Here, we linearize l only, as $\min_{\beta} g(\beta; \beta^{(t)})$ has (20) as its globally optimal solution. Interestingly, the set of fixed points under the g -mapping enjoys provable guarantees that may *not*

hold for the set of local minimizers to the original objective (Section 3.2.1). This is particularly the case when Θ has discontinuities and $P(t; \lambda)$ is given by $P_\Theta(t; \lambda) + q(t; \lambda)$, where P_Θ is defined by (48) and q is a function satisfying $q(t; \lambda) \geq 0$ for all $t \in \mathbb{R}$ and $q(t; \lambda) = 0$ if $t = \Theta(s; \lambda)$ for some $s \in \mathbb{R}$ [49].

A closely related *iterative quantile-thresholding* procedure [48, 52] proceeds by $\beta^{(t+1)} = \Theta^\#(\beta^{(t)} - \nabla l(\beta^{(t)})/\varrho^2; q)$ for the sake of feature screening: $\min l(\beta)$ s.t. $\|\beta\|_0 \leq q$, and uses a similar surrogate $g(\beta; \beta^{(t)}) = l(\beta) + (\varrho^2 \mathbf{D}_2 - \Delta_l)(\beta, \beta^{(t)})$. Here, the quantile thresholding $\Theta^\#(\alpha; q)$, as an outcome of $\min g(\beta; \beta^{(t)})$, keeps the top q elements of α_j after ordering them in magnitude, $|\alpha_{(1)}| \geq \dots \geq |\alpha_{(p)}|$, and zero out the rest. To avoid ambiguity, we assume no ties occur in performing $\Theta^\#(\alpha; q)$ throughout the paper, that is, $|\alpha_{(q)}| > |\alpha_{(q+1)}|$.

EXAMPLE 3. (Nonnegative matrix factorization). Nonnegative Matrix Factorization (NMF) [34] provides an effective tool for feature extraction and finds widespread applications in computer vision, text mining and many other areas. NMF approximates a nonnegative data matrix $\mathbf{X} \in \mathbb{R}_+^{n \times p}$ by the product of two nonnegative low-rank matrices $\mathbf{W} \in \mathbb{R}_+^{n \times r}$ and $\mathbf{H} \in \mathbb{R}_+^{r \times p}$. The KL divergence is often used to make a cost function, that is, $\min_{\mathbf{W} \in \mathbb{R}_+^{n \times r}, \mathbf{H} \in \mathbb{R}_+^{r \times p}} \text{KL}(\mathbf{X}, \mathbf{WH}) := \sum_{i,j} [X_{ij} \log(X_{ij}/(\mathbf{WH})_{ij}) - X_{ij} + (\mathbf{WH})_{ij}]$, which gives a nonconvex optimization problem. The following *multiplicative* update rule (MUR) shows good scalability in big data applications [15]:

$$(22) \quad H_{kj}^{(t+1)} = H_{kj}^{(t)} \exp \left[-\frac{1}{\rho} \sum_i \left(W_{ik} - \frac{W_{ik} X_{ij}}{(\mathbf{WH}^{(t)})_{ij}} \right) \right],$$

$$(23) \quad W_{ik}^{(t+1)} = W_{ik}^{(t)} \exp \left[-\frac{1}{\rho} \sum_j \left(H_{kj} - \frac{H_{kj} X_{ij}}{(\mathbf{W}^{(t)}\mathbf{H})_{ij}} \right) \right].$$

The update formulas can be explained from a Bregman surrogate perspective. Since the problem is symmetric in \mathbf{W} and \mathbf{H} , $\Delta_{\text{KL}}(\mathbf{X}, \mathbf{WH}) = \Delta_{\text{KL}}(\mathbf{X}^\top, \mathbf{H}^\top \mathbf{W}^\top)$, we take (22) for instance to illustrate the point. Noticing that the criterion is separable in the column vectors of \mathbf{H} , it suffices to look at $\min_{\mathbf{h} \in \mathbb{R}_+^r} f(\mathbf{h}) = \text{KL}(\mathbf{x}, \mathbf{Wh}) = \sum_i [x_i \log(x_i/(\mathbf{Wh})_i) - x_i + (\mathbf{Wh})_i]$, where \mathbf{x} can be any column of \mathbf{X} . Then it is easy to verify that the following Bregman surrogate,

$$(24) \quad g(\mathbf{h}; \mathbf{h}^{(t)}) = f(\mathbf{h}) + (\rho \mathbf{D}_\varphi - \mathbf{D}_f)(\mathbf{h}, \mathbf{h}^{(t)}), \quad \varphi(\mathbf{h}) = \sum (h_i \log h_i - h_i),$$

leads to the multiplicative update formulas.

EXAMPLE 4. (DC programming). DC programming [55] is capable of tackling a large class of nonsmooth nonconvex optimization problems; see, for example, [23, 43]. A “difference of convex” (DC) function f is defined by $f(\beta) = d_1(\beta) - d_2(\beta)$, where d_1 and d_2 are both closed convex functions. To minimize $f(\beta)$, a standard DC algorithm generates two sequences $\{\beta^{(t)}\}$ and $\{\gamma^{(t)}\}$ that obey

$$(25) \quad \gamma^{(t)} \in \partial d_2(\beta^{(t)}), \quad \beta^{(t+1)} \in \partial d_1^*(\gamma^{(t)}),$$

where $\partial d(\beta)$ is the subdifferential of $d(\cdot)$ at β , and $d_1^*(\cdot)$ is the Fenchel conjugate of $d_1(\cdot)$. (As before, d_1, d_2 are assumed to be real-valued functions defined on \mathbb{R}^p , so the sequences are well-defined and finite.) This elegant algorithm does not involve any line search and guarantees global convergence given any initial point. Many popular nonconvex algorithms can be derived from (25) [2].

Focusing on the β -update, we know that $\beta^{(t+1)}$ must be a solution to $\min_\beta d_1(\beta) - \langle \beta, \gamma^{(t)} \rangle$ or $\min_\beta d_1(\beta) - \langle \beta - \beta^{(t)}, \gamma^{(t)} \rangle$. Due to the convexity of d_2 , $\langle \beta - \beta^{(t)}, \gamma^{(t)} \rangle \leq$

$\sup_{\gamma \in \partial d_2(\beta^{(t)})} \langle \beta - \beta^{(t)}, \gamma \rangle = \delta d_2(\beta^{(t)}; \beta - \beta^{(t)})$ for all $\gamma^{(t)} \in \partial d_2(\beta^{(t)})$, $\beta \in \mathbb{R}^p$. Thus $\min_{\beta} d_1(\beta) - \langle \beta - \beta^{(t)}, \gamma^{(t)} \rangle$ should be no lower than $\min_{\beta} d_1(\beta) - \delta d_2(\beta^{(t)}; \beta - \beta^{(t)})$. Choosing $\beta^{(t+1)} \in \arg \min_{\beta} d_1(\beta) - \delta d_2(\beta^{(t)}; \beta - \beta^{(t)})$ and $\gamma^{(t)} = \delta d_2(\beta^{(t)}; \beta^{(t+1)} - \beta^{(t)}) (\beta^{(t+1)} - \beta^{(t)}) / \|\beta^{(t+1)} - \beta^{(t)}\|_2^2$ ensures (25), which simply amounts to using a Bregman surrogate

$$(26) \quad g(\beta; \beta^{(t)}) = f(\beta) + \Delta_{d_2}(\beta, \beta^{(t)}).$$

For the γ -updates, a Bregman surrogate $g(\gamma; \gamma^{(t)}) = (d_2^* - d_1^*)(\gamma) + \Delta_{d_1^*}(\gamma, \gamma^{(t)})$ can be similarly constructed.

EXAMPLE 5. (Local linear approximation). Zou and Li [67] proposed an effective local linear approximation (LLA) technique to minimize penalized negative log-likelihoods. In their paper, the loss function is assumed to be convex and smooth, and the penalty is concave on \mathbb{R}_+ . We give a new characterization of LLA by use of a Bregman surrogate.

Let l be a directionally differentiable loss function but not necessarily continuously differentiable, and P be a function that is concave and differentiable over $(0, +\infty)$, and satisfies $P(t) = P(-t)$ for any $t \in \mathbb{R}$, $P(0) = 0$. Consider the problem $\min_{\beta} l(\beta) + \sum_j P(\beta_j)$. Using the generalized Bregman notation $\Delta_{\|\cdot\|_1}(\beta, \gamma)$, or $\Delta_1(\beta, \gamma)$ for short, define

$$(27) \quad g(\beta; \beta^{(t)}) = l(\beta) + \sum P(\beta_j) + \sum [\alpha_j \Delta_1(\beta_j, \beta_j^{(t)}) - \Delta_P(\beta_j, \beta_j^{(t)})].$$

In contrast to (21), (27) linearizes P instead of l . Simple calculation shows

$$(28) \quad \Delta_1(\beta_j, \beta_j^{(t)}) = \begin{cases} |\beta_j| - \text{sgn}(\beta_j^{(t)})\beta_j, & \beta_j^{(t)} \neq 0 \\ 0, & \beta_j^{(t)} = 0, \end{cases}$$

$$(29) \quad \Delta_P(\beta_j, \beta_j^{(t)}) = \begin{cases} P(\beta_j) - P(\beta_j^{(t)}) - P'(\beta_j^{(t)})(\beta_j - \beta_j^{(t)}), & \beta_j^{(t)} \neq 0 \\ P(\beta_j) - P'_+(0)|\beta_j|, & \beta_j^{(t)} = 0, \end{cases}$$

where $\text{sgn}(\cdot)$ is the sign function and $P'_+(\beta)$ denotes the right derivative of $P(\cdot)$ at β . Interestingly, with $\alpha_j = |P'_+(\beta_j^{(t)})|$, the Δ_1 -based surrogate (27) can be shown to be

$$l(\beta) + \sum_j [P(|\beta_j^{(t)}|) + P'_+(|\beta_j^{(t)}|)(|\beta_j| - |\beta_j^{(t)}|)],$$

which is exactly the surrogate constructed by Zou and Li. To the best of our knowledge, the generalized Bregman formulation is new.

LLA requires solving a weighted lasso problem at each step. We can further linearize l as in Example 2 to improve its scalability. LLA is popular among statisticians, but to our knowledge, there is a lack of *global* convergence-rate studies in large- p applications. We will see that reformulating LLA from the generalized Bregman surrogate perspective leads to a convenient choice of the convergence measure in analyzing the algorithm.

EXAMPLE 6. (Sigmoidal regression). We use the univariate-response sigmoidal regression to illustrate this type of nonconvex problems that is commonly seen in artificial neural networks. The formulation carries over to multilayered networks and recurrent networks [51].

Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$ be the data matrix, and $\mathbf{y} = [y_1, \dots, y_n]^\top$ be the response vector. Define $\pi(\nu) = e^\nu / (1 + e^\nu)$; if ν is replaced by a vector, π is defined componentwise. The sigmoidal regression solves

$$(30) \quad \min_{\beta} f(\beta) = \frac{1}{2} \sum_{i=1}^n (y_i - \pi(\mathbf{x}_i^\top \beta))^2.$$

Then $\nabla^2 f(\beta) = \sum_{i=1}^n [(-2\mu_i^3 + 3\mu_i^2 - \mu_i)y_i + (3\mu_i^4 - 5\mu_i^3 + 2\mu_i^2)]x_i x_i^\top$, where $\mu_i = \pi(x_i^\top \beta)$. Because $\mu_i \in [0, 1]$, we get $\nabla^2 f(\beta) \preceq \mathbf{X}^\top \text{diag}\{|0.1y_i| + 0.08\}_{i=1}^n \mathbf{X}$, which motivates a Bregman surrogate

$$g(\beta; \beta^{(t)}) = f(\beta) + \mathbf{D}_{\psi-f}(\beta, \beta^{(t)}), \quad \psi(\beta) = \frac{1}{2} \beta^\top \mathbf{X}^\top \text{diag}\{|0.1y_i| + 0.08\}_{i=1}^n \mathbf{X} \beta.$$

Solving $\min_{\beta} g(\beta; \beta^{(t)})$ yields $\beta^{(t+1)} = \beta^{(t)} + \mathbf{B}^{-1} \mathbf{X}^\top (\mathbf{u}^{(t)} - \mathbf{u}^{(t)} \circ \mathbf{u}^{(t)}) \circ (\mathbf{y} - \mathbf{u}^{(t)})$, where $\mathbf{B} = \mathbf{X}^\top \text{diag}\{|0.1y_i| + 0.08\}_{i=1}^n \mathbf{X}$, $\mathbf{u}^{(t)} = \pi(\mathbf{X}^\top \beta^{(t)})$ and \circ denotes the Hadamard product. This type of surrogate functions is closely related to proximal Newton-type methods [46] and signomial programming [33].

3. Bregman-surrogate algorithm analysis. Motivated by the examples in Section 2, we study a generalized Bregman-surrogate algorithm family for solving $\min_{\beta} f(\beta)$, with the sequence of iterates defined by

$$(31) \quad \beta^{(t+1)} \in \arg \min_{\beta} g(\beta; \beta^{(t)}) := f(\beta) + \Delta_{\psi}(\beta, \beta^{(t)}), \quad t \geq 0$$

The objective function f and the auxiliary function ψ are assumed to be directionally differentiable but need not be smooth or convex. ψ has flexible options as seen from the previous examples.

Equation (31) does not necessarily give an MM procedure, as the majorization condition $g(\beta; \beta^-) \geq f(\beta)$ may not hold. But we have the following zeroth-order and first-order degeneracies when $\beta^- = \beta$, which provides rationality of investigating the accuracy of *fixed points* under the g -mapping (31).

LEMMA 4. *Let $g(\beta; \beta^-) = f(\beta) + \Delta_{\psi}(\beta, \beta^-)$ with f and ψ directionally differentiable. Then (i) $g(\beta; \beta) = f(\beta)$, and (ii) $\delta g(\beta; \beta^-, \mathbf{h})|_{\beta^- = \beta} = \delta f(\beta; \mathbf{h}), \forall \beta, \mathbf{h}$, where $\delta g(\beta; \beta^-, \mathbf{h})$ is the directional derivative of $g(\cdot; \beta^-)$ at β with increment \mathbf{h} .*

The lemma relates the set of fixed points of the algorithm mapping,

$$(32) \quad \{\beta : \beta \in \arg \min_{\beta} g(\beta; \beta^-)|_{\beta^- = \beta}\},$$

which we will call the fixed points of g for short, to the set of directional stationary points of f (under directional differentiability),

$$(33) \quad \{\beta : \delta f(\beta; \mathbf{h}) \geq 0 \text{ for any admissible } \mathbf{h}\},$$

which becomes the set of stationary points when $f \in \mathcal{C}^1$. The link is general for any generalized Bregman surrogate in (31) *regardless* of the specific form of ψ . An important implication is that in studying convergence it is legitimate to measure how $\beta^{(t+1)}$ and $\beta^{(t)}$ differ, as widely used in practice. Later we will see that it is indeed possible to provide provable guarantees for the fixed points of this type of surrogates. In contrast, a general MM algorithm does not always have the first-order degeneracy and so attaining $\beta^{(t+1)} = \beta^{(t)}$ does not necessarily ensure a good-quality solution, especially in nonconvex scenarios.

3.1. Computational accuracy. We first study the optimization error of (31), then turn to its statistical error in Section 3.2. This subsection aims to derive universal rates of convergence under no regularity conditions.

• *General setting.* In this part, the objective $f(\beta)$ does not have any known structure. To better connect with some conventional results in convex optimization, we first present two propositions for (31) on the function-value convergence and iterate convergence. While the resultant rates are encouraging, the error bounds are most informative under certain smoothness and convexity assumptions. This suggests the necessity of choosing a proper convergence measure in order to avoid stringent or awkward technical conditions in nonconvex optimization.

PROPOSITION 1. *Given an arbitrary initial point $\beta^{(0)}$, let $\beta^{(t)}$ be the sequence generated according to (31) where ψ is differentiable. Then*

$$(34) \quad \text{avg}_{0 \leq t \leq T} f(\beta^{(t+1)}) - f(\bar{\beta}) \leq \frac{1}{T+1} [\Delta_\psi(\bar{\beta}, \beta^{(0)}) - \Delta_\psi(\bar{\beta}, \beta^{(T+1)})]$$

for any $\bar{\beta}$ satisfying

$$(35) \quad \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + \Delta_f(\bar{\beta}, \beta^{(t+1)}) \geq 0, \quad 0 \leq t \leq T.$$

Here, $\text{avg}_{0 \leq t \leq T} f(\beta^{(t+1)})$ denotes the average of $f(\beta^{(1)}), \dots, f(\beta^{(T+1)})$.

In particular, if both f and ψ are convex, then $f(\beta^{(t)})$ is nonincreasing and

$$(36) \quad f(\beta^{(T+1)}) - f(\beta) \leq \frac{\Delta_\psi(\beta, \beta^{(0)})}{T+1}, \quad \forall \beta.$$

Equation (34) shows a convergence rate of $\mathcal{O}(1/T)$ under (35) that amounts to step size control. For example, for $\Delta_\psi = \rho \mathbf{D}_\varphi - \Delta_f$ in mirror descent, (35) shows that ρ should be sufficiently large, which in turns gives a small stepsize $1/\rho$:

$$\rho \geq (\Delta_f(\beta^{(t+1)}, \beta^{(t)}) - \Delta_f(\bar{\beta}, \beta^{(t+1)})) / \mathbf{D}_\varphi(\beta^{(t+1)}, \beta^{(t)}),$$

or $\rho \geq \Delta_f(\beta^{(t+1)}, \beta^{(t)}) / \mathbf{D}_\varphi(\beta^{(t+1)}, \beta^{(t)})$ when f is convex. In nonconvex scenarios, the condition may be hard to verify, but one has reason to believe that with a properly small step size, a generalized Bregman-surrogate algorithm should not be much slower than gradient descent.

Actually, a faster rate of convergence may be obtained under some GBF comparison conditions, (37) and (39) below, which can be viewed as substitutes for conventional strong convexity in a more general sense. (The corresponding geometric decay of the errors is motivating in high dimensional statistical learning, in light of the “restricted” strongly convexity often possessed by such a type of problems [36].)

PROPOSITION 2. *Consider the iterative algorithm defined by (31) starting at an arbitrary point $\beta^{(0)}$ with ψ differentiable, and let β^o be a minimizer of $f(\beta)$. (i) If for some $\kappa > 1$, $\Delta_\phi = \Delta_\psi + \Delta_f$ satisfies*

$$(37) \quad \bar{\Delta}_\phi \geq \frac{\kappa}{\kappa - 1} \Delta_\psi,$$

then for any $T \geq 0$, we have

$$(38) \quad \bar{\Delta}_\phi(\beta^o, \beta^{(T+1)}) \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^{T+1} \bar{\Delta}_\phi(\beta^o, \beta^{(0)}) - \frac{\kappa}{2} \min_{0 \leq t \leq T} \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}).$$

(ii) Alternatively, if

$$(39) \quad 2\bar{\Delta}_f \geq \varepsilon \Delta_\psi$$

for some $\varepsilon > 0$, then

$$(40) \quad \Delta_\psi(\beta^o, \beta^{(T+1)}) \leq \left(\frac{1}{1+\varepsilon}\right)^{T+1} \Delta_\psi(\beta^o, \beta^{(0)}) - \frac{1}{\varepsilon} \min_{0 \leq t \leq T} \Delta_\psi(\beta^{(t+1)}, \beta^{(t)})$$

for any $T \geq 0$.

REMARK 1. We give an illustration of (i) and (ii) to compare their assumptions and conclusions. In gradient descent with $\Delta_\phi = \rho \mathbf{D}_2$, (37) becomes $\rho \mathbf{D}_2 \geq (\rho \mathbf{D}_2 - \Delta_f) \kappa / (\kappa - 1)$ or $\Delta_f \geq (\rho / \kappa) \mathbf{D}_2$ and when f is μ -strongly convex and ρ -strongly smooth, $\kappa = \rho / \mu$. Then (38) reads

$$(41) \quad \mathbf{D}_2(\beta^o, \beta^{(T+1)}) \leq \left(\frac{\rho - \mu}{\rho + \mu}\right)^{T+1} \mathbf{D}_2(\beta^o, \beta^{(0)}).$$

The \mathbf{D}_2 -form bound is classical for problems with strong convexity; see, for example, Theorem 2.1.15 in [41]. Yet it is worth mentioning that our Bregman comparison conditions do not require ψ to be strongly convex to attain the linear rate. (40) gives a linear convergence result, too, in terms of yet another measure. In the same setup, (39) holds for $\varepsilon : \varepsilon \rho / (2 + \varepsilon) = \mu$ and similarly

$$(42) \quad \Delta_\psi(\beta^o, \beta^{(T+1)}) \leq \left(\frac{\rho - \mu}{\rho + \mu}\right)^{T+1} \Delta_\psi(\beta^o, \beta^{(0)}).$$

A careful examination of the proof in Section A.8 shows that (39) is applied once, while (37) is applied twice on both sides of (A.13), and so (ii) appears less technically demanding. Picking a suitable error function can assist analysis and relax regularity assumptions. The same Δ_ψ will be used in studying the statistical error convergence in Theorem 5.

Instead of naively comparing $f(\beta^{(t)})$ with f^o , or $\beta^{(t)}$ with β^o , which may be unattainable or nonunique in nonconvex optimization, one can measure the algorithm convergence in a wiser manner. Ben-Tal and Nemirovski [5] pointed out that with an inappropriate measure of discrepancy, the convergence rate of gradient descent for minimizing a nonconvex objective can be arbitrarily slow, and a common choice is to bound

$$(43) \quad \min_{t \leq T} \|\nabla f(\beta^{(t)})\|^2.$$

This is reasonable since when $\nabla f(\beta^{(t)}) = 0$, gradient descent stops iterating and delivers a stationary point. (43) can be rewritten as ρ^2 times

$$(44) \quad \min_{t \leq T} \mathbf{D}_2(\beta^{(t+1)}, \beta^{(t)})$$

as $\beta^{(t+1)} - \beta^{(t)} = -\nabla f(\beta^{(t)}) / \rho$. The idea of checking stationarity by the difference between two successive iterates generalizes, thanks to Lemma 4, and eventually leads to an error bound that can get rid of condition (35).

THEOREM 2. Any generalized Bregman surrogate algorithm defined by (31) satisfies the following bound for all $T \geq 1$,

$$(45) \quad \text{avg}_{0 \leq t \leq T} (2\bar{\Delta}_\psi + \Delta_f)(\beta^{(t)}, \beta^{(t+1)}) \leq \frac{1}{T+1} [f(\beta^{(0)}) - f(\beta^{(T+1)})].$$

(45) obtains the same rate of convergence as Proposition 1, but is free of any conditions other than directional differentiability, because only the weak idempotence is needed to derive

the bound. A proper stepsize control can often make the GBF error nonnegative (e.g., (50)). But even when $\beta^{(t)}$ diverges, (45) still applies.

Notice the factor ‘2’ proceeding the symmetrized Bregman $\bar{\Delta}_\psi$ on the left-hand side of (45). This gives a relaxed stepsize control than MM. We use mirror descent $\Delta_\psi = \rho \bar{\mathbf{D}}_\varphi - \Delta_f$ to exemplify the point without requiring f to be convex, cf. Example 1.

COROLLARY 1. *In the mirror descent setup with a possibly nonconvex objective, suppose that $\Delta_f \leq L \bar{\mathbf{D}}_\varphi$ for some $L > 0$, $\inf_{\beta} f(\beta) \geq 0$, and the inverse stepsize parameter ρ is taken such that $\rho > L/2$. Then any accumulation point of $\beta^{(t)}$ is a fixed point of g and*

$$(46) \quad \text{avg}_{0 \leq t \leq T} \bar{\mathbf{D}}_\varphi(\beta^{(t)}, \beta^{(t+1)}) \leq \frac{f(\beta^{(0)})}{(T+1)(2\rho - L)}.$$

Hence in the special case of gradient descent, (46) recovers $\min_{0 \leq t \leq T} \|\nabla f(\beta^{(t)})\|_2^2 = \mathcal{O}(1/T)$ [5] when $\rho > L/2$. In comparison, MM algorithms always require $\Delta_\psi \geq 0$, or $\rho \geq L$. A smaller value of ρ means a larger step size with which the algorithm converges faster.

• *Composite setting.* High-dimensional statistical learning often has an additive objective $f(\beta) = l_0(\mathbf{X}\beta) + P(\varrho\beta; \lambda)$, where $\mathbf{X} \in \mathbb{R}^{n \times p}$ is the predictor or feature matrix, $l_0(\cdot)$ is the loss defined on $\mathbf{X}\beta$ (and so $l(\beta) = l_0(\mathbf{X}\beta)$), $P(\cdot; \lambda)$ is a sparsity-inducing regularizer and ϱ is a controllable parameter, typically taking $\|\mathbf{X}\|_2$ to match the scale. Unless otherwise mentioned, $P(\beta; \lambda)$ denotes $\sum_j P(\beta_j; \lambda)$ with a little abuse of notation.

Such a composite setup is widely assumed in convex optimization [57, 18]. But among the abundant choices of l_0 and P in the literature, many of them are nonconvex. The good news is that the main theorem proved in the previous subsection adapts to the composite setting and we give some results for iterative thresholding and LLA as an illustration (cf. Examples 2, 5).

Iterative thresholding. Many popularly used penalty functions are associated with thresholdings rigorously defined as follows.

DEFINITION 3 (Thresholding function). *A threshold function is a real-valued function $\Theta(t; \lambda)$ defined for $-\infty < t < \infty$ and $0 \leq \lambda < \infty$ such that (i) $\Theta(-t; \lambda) = -\Theta(t; \lambda)$; (ii) $\Theta(t; \lambda) \leq \Theta(t'; \lambda)$ for $t \leq t'$; (iii) $\lim_{t \rightarrow \infty} \Theta(t; \lambda) = \infty$; (iv) $0 \leq \Theta(t; \lambda) \leq t$ for $0 \leq t < \infty$.*

Given Θ , a critical concavity number $\mathcal{L}_\Theta \leq 1$ can be introduced such that $d\Theta^{-1}(u; \lambda) du \geq 1 - \mathcal{L}_\Theta$ for almost every $u \geq 0$, or

$$(47) \quad \mathcal{L}_\Theta = 1 - \text{ess inf} \{d\Theta^{-1}(u; \lambda)/du : u \geq 0\},$$

with ess inf the essential infimum and $\Theta^{-1}(u; \lambda) := \sup\{t : \Theta(t; \lambda) \leq u\}, \forall u > 0$. For the widely used soft-thresholding $\Theta_S(t; \lambda) = \text{sgn}(t)(|t| - \lambda)1_{|t| > \lambda}$ and hard-thresholding $\Theta_H(t; \lambda) = t1_{|t| > \lambda}$, \mathcal{L}_Θ equals 0 and 1, respectively. In fact, when $\mathcal{L}_\Theta > 0$, the penalty induced by Θ via (48) is nonconvex, and \mathcal{L}_Θ gives a concavity measure of it according to Lemma A.3. The Bregman surrogate characterization of iterative thresholding in (21) yields a general conclusion for any Θ in possibly high dimensions.

PROPOSITION 3. *Given any thresholding Θ and directionally differentiable $l(\cdot)$, consider the iterative thresholding procedure (20): $\beta^{(t+1)} = \Theta(\varrho\beta^{(t)} - \nabla l(\beta^{(t)})/\varrho; \lambda)/\varrho$ with $\varrho > 0$. Construct*

$$(48) \quad P_\Theta(t; \lambda) = \int_0^{|t|} (\Theta^{-1}(u; \lambda) - u) du, \quad \forall t \in \mathbb{R},$$

and define $f(\beta) = l(\beta) + P_\Theta(\varrho\beta; \lambda)$, $g(\beta, \beta^-) = l(\beta) + P_\Theta(\varrho\beta; \lambda) + (\varrho^2 \mathbf{D}_2 - \Delta_l)(\beta, \beta^-)$. Then $\beta^{(t)} \in \arg \min_\beta g(\beta, \beta^{(t-1)})$ and for all $T \geq 1$

$$(49) \quad \text{avg}_{0 \leq t \leq T} (\varrho^2(2 - \mathcal{L}_\Theta) \mathbf{D}_2 - \Delta_l)(\beta^{(t)}, \beta^{(t+1)}) \leq \frac{1}{T+1} [f(\beta^{(0)}) - f(\beta^{(T+1)})].$$

When the loss satisfies $\Delta_l \leq L \mathbf{D}_2$, a reasonable choice of ϱ is

$$(50) \quad \varrho^2 > L/(2 - \mathcal{L}_\Theta).$$

So when $\mathcal{L}_\Theta > 0$, the step size upper bound will be smaller than that as $\mathcal{L}_\Theta = 0$. This is often the price to pay for nonconvex optimization. On the other hand, (49) still ensures the universal rate of convergence of $\mathcal{O}(1/T)$, in spite of the high dimensionality and nonconvexity.

Local linear approximation. Next, we study the computational convergence of LLA for solving the penalized estimation problem $\min f(\beta) = l(\beta) + P(\varrho\beta)$, assuming l is directionally differentiable, $P(0) = 0$, $P'_+(0) < +\infty$, $P(t) = P(-t) \geq 0$ and $P(t)$ is differentiable for any $t > 0$. Recall its Bregman form surrogate

$$(51) \quad g_{\text{LLA}}^{(t)}(\beta; \beta^{(t)}) = l(\beta) + P(\varrho\beta) + \Delta_{\|\alpha^{(t)} \circ (\cdot)\|_1 - P(\cdot)}(\varrho\beta, \varrho\beta^{(t)}),$$

where $\alpha^{(t)} = [\alpha_j^{(t)}]$ with $\alpha_j^{(t)} = |P'_+(\beta_j^{(t)})|$, $1 \leq j \leq p$. We abbreviate $\Delta_{\|\alpha^{(t)} \circ (\cdot)\|_1 - P(\cdot)}$ to $\Delta_{\text{LLA}}^{(t)}$, which does not satisfy strong idempotence. By combining $\bar{\Delta}_{\text{LLA}}^{(t)}$ and Δ_f to evaluate LLA's optimization error, we obtain a convergence result without any additional assumptions.

PROPOSITION 4. *Given any starting point $\beta^{(0)}$, the LLA iterates satisfy the following bound for all $T \geq 1$:*

$$\text{avg}_{0 \leq t \leq T} [2\bar{\Delta}_{\text{LLA}}^{(t)}(\varrho\beta^{(t)}, \varrho\beta^{(t+1)}) + \Delta_f(\beta^{(t)}, \beta^{(t+1)})] \leq \frac{1}{T+1} [f(\beta^{(0)}) - f(\beta^{(T+1)})].$$

Ignoring the cost difference per iteration, the convergence rate of LLA is no slower than that of gradient descent. If l is a negative log-likelihood function associated with a log-concave density and P is concave on \mathbb{R}_+ , as assumed in [67], $2\bar{\Delta}_{\text{LLA}}^{(t)}(\varrho\beta, \varrho\beta') + \Delta_f(\beta, \beta') = \Delta_l(\beta, \beta') + \Delta_{-P}(\varrho\beta', \varrho\beta) + 2 \sum_j \alpha_j^{(t)} \bar{\Delta}_1(\varrho\beta_j, \varrho\beta'_j) \geq 0, \forall \beta, \beta'$. But Proposition 4 holds even when P is nonconcave on \mathbb{R}_+ and l is nonconvex.

The global convergence-rate results presented in this subsection are free of any regularity conditions on sparsity, sample size, initial point and design incoherence. High-dimensional learning algorithms may however show a better convergence rate when the problems under consideration are “regular” in a certain sense.

3.2. Statistical accuracy. To statisticians, the statistical accuracy of Bregman-surrogate algorithms with respect to a statistical truth (denoted by β^*) is perhaps more meaningful than the optimization error to a certain local or global minimizer, since real world data are always noisy. Section 3.2.1 and Section 3.2.2 will study the statistical error of the final estimate $\hat{\beta}$ and the t -th iterate $\beta^{(t)}$, respectively, where combining the generalized Bregman calculus and the empirical process theory eases the treatment of a nonquadratic loss.

The techniques based on GBFs apply to a general problem (see, e.g., Theorem A.1 in Section A.18), but here we focus on the aforementioned sparse learning in the composite setting: $\min_\beta l(\beta) + P_\Theta(\varrho\beta; \lambda)$, where $l(\beta) = l_0(\eta) = l_0(\mathbf{X}\beta)$ is directionally differentiable and $P_\Theta(\cdot; \lambda)$ is induced by a thresholding Θ via (48). Since l_0 is placed on $\mathbf{X}\beta$, we include here a scaling parameter ϱ (often $\|\mathbf{X}\|_2$) in the penalty; this will yield a universal choice of

the regularization parameter λ that does not vary with the sample size. Throughout Section 3.2, we assume that ϱ satisfies $\varrho \geq \|\mathbf{X}\|_2$. Note that neither the loss nor the penalty needs to be convex or smooth.

Give any directionally differentiable ψ , the sequence of iterates is generated by

$$(52) \quad \beta^{(t+1)} \in \arg \min_{\beta} g(\beta; \beta^{(t)}) := l(\beta) + P_{\Theta}(\varrho\beta; \lambda) + \Delta_{\psi}(\beta, \beta^{(t)}).$$

Nonconvex iterative thresholding and LLA are particular instances.

First, we must characterize the notion of noise in this nonlikelihood setting, to take into account the randomness of samples. Assume l_0 is differentiable at point $\mathbf{X}\beta^*$ (but not necessarily differentiable on all of \mathbb{R}^n) and define the *effective noise* by

$$(53) \quad \epsilon = -\nabla l_0(\mathbf{X}\beta^*).$$

(An alternative assumption is that $\delta l_0(\mathbf{X}\beta^*; \mathbf{h})$ is a sub-Gaussian random variable with mean 0 and scale bounded by $c\sigma$ for any unit vector \mathbf{h} , but we will not pursue further in the current paper.)

Typically, $\mathbb{E}[\epsilon]$ should be 0, and so $\nabla\{\mathbb{E}[l_0(\mathbf{X}\beta^*)]\} = 0$ assuming the differentiation and expectation are exchangeable, which means the statistical truth makes the gradient of its risk vanish. For a GLM with y_i ($1 \leq i \leq n$) following a distribution in the exponential family that has cumulant function b and canonical link function $g = (b')^{-1}$, the loss is then $l(\beta) = l_0(\mathbf{X}\beta) = -\langle \mathbf{y}, \mathbf{X}\beta \rangle + \langle \mathbf{1}, b(\mathbf{X}\beta) \rangle$ (cf. (9) with $\sigma = 1$), and so

$$(54) \quad \epsilon = \mathbf{y} - g^{-1}(\mathbf{X}\beta^*) = \mathbf{y} - \mathbb{E}(\mathbf{y}).$$

Our effective noise, as a joint outcome of the loss and the response, does not depend on the regularizer, and may differ from the raw noise. For example, under $\mathbf{y} = \mathbf{X}\beta^* + \epsilon^{\text{raw}}$, $l(\beta) = l_{\text{Huber}}(\mathbf{r}) = \sum_{i: |r_i| \leq a\sigma} r_i^2/2 + \sum_{i: |r_i| > a\sigma} (a|r_i| - a^2\sigma^2/2)$ with $\mathbf{r} = \mathbf{y} - \mathbf{X}\beta$ [27], simple calculation gives $\epsilon_i = \epsilon_i^{\text{raw}} 1_{|\epsilon_i^{\text{raw}}| \leq a\sigma} + a\sigma 1_{|\epsilon_i^{\text{raw}}| > a\sigma}$, which is bounded by $a\sigma$, thereby sub-Gaussian, no matter what distribution the raw noise follows. This nonparametricness is apparent for any l_0 that is (globally) Lipschitz, for example, the logistic deviance and hinge loss for classification.

In this section, we assume that ϵ is a sub-Gaussian random vector with mean zero and scale bounded by σ , cf. Definition A.1, where ϵ_i are not required to be independent. Examples include Gaussian random variables and bounded random variables such as Bernoulli.

The support of β is denoted by $\mathcal{J}(\beta) = \{j : \beta_j \neq 0\}$, and its cardinality is $J(\beta) = |\mathcal{J}(\beta)| = \|\beta\|_0$. We abbreviate $J(\beta^*)$ to J^* and $J(\hat{\beta})$ to \hat{J} . In sparse learning, $J^* \ll n \ll p$ is typically true. The sparsity suggests the possibility of obtaining a fast rate of convergence in statistical error. The following penalty induced by the hard-thresholding $\Theta_H(t; \lambda) = t 1_{|t| > \lambda}$ by (48) turns out to play a key role in the analysis

$$(55) \quad P_H(t; \lambda) = (-t^2/2 + \lambda|t|) 1_{|t| < \lambda} + (\lambda^2/2) 1_{|t| \geq \lambda}.$$

An important fact is that $P_{\Theta}(t; \lambda) \geq P_H(t; \lambda)$ for any $t \in \mathbb{R}$ and any thresholding rule Θ . This is simply because in shrinkage estimation, any $\Theta(t; \lambda)$ with λ as the threshold is identical to zero as $t \in [0, \lambda)$ and is bounded above by the identity line for $t \geq \lambda$.

3.2.1. Statistical accuracy of fixed-point solutions. The finally obtained solutions from a Bregman surrogate algorithm can be described as the fixed points of g (recall (32)),

$$(56) \quad \hat{\beta} \in \arg \min_{\beta} g(\beta; \hat{\beta}).$$

We denote the set by \mathcal{F} , and call such solutions the *F-estimators*. When the objective function is convex, an F-estimator is necessarily a globally optimal solution to the original problem by

Lemma 4, thus an M-estimator. In general, however, the lack of convexity and smoothness may make $\hat{\beta}$ neither an M-estimator nor a Z-estimator [60], which poses new and intriguing challenges to statistical algorithmic analysis. It is also worth mentioning that another important class of “A-estimators” that have *alternative* optimality, typically arising from block coordinate descent (BCD) algorithms like in Example 3, can often be converted to F-estimators; see Section A.17.

Nicely, if the problem is regular, all F-estimators defined through g can achieve essentially the best statistical precision in possibly high dimensions. This is nontrivial since even f ’s locally optimal solutions do not all have the provable guarantee (cf. Remark 4). Theorem 3 and Theorem 4 below only make use of the weak idempotence property; another notable feature is that the conditions and conclusions below are *regardless* of the form of Δ_ψ .

THEOREM 3. *Suppose there exist $\delta > 0$, $\vartheta > 0$ and large enough $K \geq 0$ so that the following inequality holds for any $\beta \in \mathbb{R}^p$:*

$$(57) \quad \begin{aligned} & \varrho^2 \mathcal{L}_\Theta \mathbf{D}_2(\beta, \beta^*) + \delta \mathbf{D}_2(\mathbf{X}\beta, \mathbf{X}\beta^*) + \vartheta P_H(\varrho(\beta - \beta^*); \lambda) + P_\Theta(\varrho\beta^*; \lambda) \\ & \leq 2\bar{\Delta}_l(\beta, \beta^*) + P_\Theta(\varrho\beta; \lambda) + K\lambda^2 J(\beta^*), \end{aligned}$$

where $\lambda = A\sigma\sqrt{\log(ep)}/\sqrt{(\delta \wedge \vartheta)\vartheta}$ with A a sufficiently large constant. Then

$$(58) \quad \mathbf{D}_2(\mathbf{X}\hat{\beta}, \mathbf{X}\beta^*) \leq \frac{2KA^2}{(\delta \wedge \vartheta)\delta\vartheta} \sigma^2 J^* \log(ep),$$

$$(59) \quad P_H(\varrho(\hat{\beta} - \beta^*); \lambda) \leq \frac{4KA^2}{(\delta \wedge \vartheta)\vartheta^2} \sigma^2 J^* \log(ep),$$

with probability at least $1 - Cp^{-cA^2}$, where C, c are positive constants.

Moreover, an *oracle inequality* [17, 31] can be built to justify the estimators even when β^* is not exactly sparse. Toward this goal, recall the notion of a pseudo-metric d (cf. Definition A.2), that is, d is nonnegative, symmetric, and satisfies the triangle inequality, and suppose without loss of generality that

$$\alpha d^2(\eta, \eta') \leq \Delta_{l_0}(\eta, \eta') \leq L d^2(\eta, \eta'), \forall \eta, \eta'$$

for some pseudo-metric d with $-\infty \leq \alpha \leq L \leq +\infty$. For regression $l(\beta) = l_0(\eta) = \|\mathbf{y} - \eta\|_2^2/2$, $\alpha = L = 1 > 0$.

THEOREM 4. *Assume for given $\beta \in \mathbb{R}^p$, there exist $r: 0 \leq r < 1, \alpha r/L \geq 0$, positive δ , ϑ , and a large enough $K \geq 0$ so that*

$$(60) \quad \begin{aligned} & \varrho^2 \mathcal{L}_\Theta \mathbf{D}_2(\beta, \gamma) + \delta \mathbf{D}_2(\mathbf{X}\beta, \mathbf{X}\gamma) + \vartheta P_H(\varrho(\beta - \gamma); \lambda) + P_\Theta(\varrho\beta; \lambda) \\ & \leq (1 + \frac{\alpha}{L}r) \Delta_l(\beta, \gamma) + P_\Theta(\varrho\gamma; \lambda) + K\lambda^2 J(\beta) \end{aligned}$$

for any $\gamma \in \mathbb{R}^p$, where $\lambda = A\sigma\sqrt{\log(ep)}/\sqrt{(\delta \wedge \vartheta)\vartheta}$ with A a sufficiently large constant. The oracle inequality below holds for some constant $C > 0$,

$$(61) \quad \mathbb{E} \Delta_l(\hat{\beta}, \beta^*) \leq \mathbb{E} \left\{ \left(\frac{1+r}{1-r} \right)^2 \Delta_l(\beta, \beta^*) + \frac{(1+r)KA^2}{(1-r)(2\vartheta \wedge \delta)\vartheta} \sigma^2 J(\beta) \log(ep) \right\} + \frac{C(1+r)}{(1-r)(2\vartheta \wedge \delta)} \sigma^2.$$

Compared with (57) which fixes γ at β^* , (60) has $(1 + \frac{\alpha}{L}r) \Delta_l$ in place of $2\bar{\Delta}_l$ as the first term on the right-hand side. Nonrigorously, these conditions ask $2\bar{\Delta}_l$ or $(1 + \frac{\alpha}{L}r) \Delta_l$ to dominate $\varrho^2 \mathcal{L}_\Theta \mathbf{D}_2$ in a restricted sense; Remark 2 argues that (60) is not technically demanding compared with many other regularity conditions in the literature.

When $r = 0$, the multiplicative constant preceding $\Delta_l(\beta, \beta^*)$ in (61) is as small as 1, resulting in a sharp oracle inequality [31]. If one sets $\beta = \beta^*$ in (61), the Bregman error $\Delta_l(\hat{\beta}, \beta^*)$ is of the order $\sigma^2 J^* \log(ep)$ for any thresholding (when δ, ϑ, K are treated as constants). But the bias term $\Delta_l(\beta, \beta^*)$ or $\Delta_{l_0}(\mathbf{X}\beta, \mathbf{X}\beta^*)$ helps to handle *approximately* sparse signals: when β^* contains a number of small nonzero elements, rather than taking $\beta = \beta^*$, a reference β with a reduced support will yield an even smaller error bound benefiting from the bias-variance tradeoff.

Unlike the optimization error bounds, the statistical error bounds never vanish (unless $\sigma \rightarrow 0$). We can similarly analyze the set of global minimizers, in which case the term $\varrho^2 \mathcal{L}_\Theta \mathbf{D}_2(\beta, \beta^*)$ is dropped from the regularity conditions, but the error bounds remain of the same order (cf. Remark A.1 in Section A.12). In fact, for sparse GLMs, by Theorem 1, the rate $\sigma^2 J^* \log(ep)$ is essentially minimax optimal (thus unbeatable) up to a logarithmic factor.

REMARK 2 (Regularity condition comparison). *The GBF-based regularity conditions (57), (60) are no more demanding than some commonly used regularity conditions. Assume that P_Θ is subadditive: $P_\Theta(t + s) \leq P_\Theta(t) + P_\Theta(s)$, which holds when it is concave on \mathbb{R}_+ . Let $\mathcal{J} = \mathcal{J}(\beta)$, $J = |\mathcal{J}(\beta)|$, $\gamma = \beta' - \beta$. Then, from $P_\Theta(\varrho \beta'_{\mathcal{J}}; \lambda) - P_\Theta(\varrho \beta_{\mathcal{J}}; \lambda) \leq P_\Theta(\varrho(\beta' - \beta)_{\mathcal{J}}; \lambda)$ and $P_\Theta(\varrho \beta'_{\mathcal{J}^c}; \lambda) = P_\Theta(\varrho(\beta' - \beta)_{\mathcal{J}^c}; \lambda)$, (60) is implied by $P_\Theta(\varrho \gamma_{\mathcal{J}}; \lambda) + \vartheta P_H(\varrho \gamma_{\mathcal{J}}; \lambda) + \mathcal{L}_\Theta \mathbf{D}_2(\varrho \beta, \varrho \beta') + \delta \|\mathbf{X}\gamma\|_2^2/2 \leq (2 - \varepsilon) \Delta_l(\beta, \beta') + K \lambda^2 J + P_\Theta(\varrho \gamma_{\mathcal{J}^c}; \lambda) - \vartheta P_H(\varrho \gamma_{\mathcal{J}^c}; \lambda)$, or $(1 + \vartheta) P_\Theta(\varrho \gamma_{\mathcal{J}}; \lambda) + \mathcal{L}_\Theta \mathbf{D}_2(\varrho \beta, \varrho \beta') + \delta \|\mathbf{X}\gamma\|_2^2/2 \leq (2 - \varepsilon) \Delta_l(\beta, \beta') + K \lambda^2 J + (1 - \vartheta) P_\Theta(\varrho \gamma_{\mathcal{J}^c}; \lambda)$ since $P_H \leq P_\Theta$.*

To get more intuition, let $l(\beta) = \|\mathbf{X}\beta - \mathbf{y}\|_2^2/2$. Then the above condition simplifies to $(1 + \vartheta) P_\Theta(\varrho \gamma_{\mathcal{J}}; \lambda) + \mathcal{L}_\Theta \|\varrho \gamma\|_2^2/2 \leq (2 - \varepsilon') \|\mathbf{X}\gamma\|_2^2/2 + K \lambda^2 J + (1 - \vartheta) P_\Theta(\varrho \gamma_{\mathcal{J}^c}; \lambda)$ with $\varepsilon' = \varepsilon + \delta$, or the following sufficient condition (with K redefined) for all $\gamma \in \mathbb{R}^p$:

$$(62) \quad (1 + \vartheta) P_\Theta(\varrho \gamma_{\mathcal{J}}; \lambda) + \frac{\mathcal{L}_\Theta}{2} \|\varrho \gamma\|_2^2 \leq K \sqrt{J} \lambda \|\mathbf{X}\gamma\|_2 + (1 - \vartheta) P_\Theta(\varrho \gamma_{\mathcal{J}^c}; \lambda).$$

For lasso, where $P_\Theta(\beta; \lambda) = \lambda \|\beta\|_1$, there is a rich collection of regularity conditions in the literature. In this convex case, $\mathcal{L}_\Theta = 0$ and ϱ can be arbitrarily large. (62) reduces to (with ϑ and K redefined and λ canceled)

$$(63) \quad (1 + \vartheta) \varrho \|\gamma_{\mathcal{J}}\|_1 \leq K \sqrt{J} \|\mathbf{X}\gamma\|_2 + \varrho \|\gamma_{\mathcal{J}^c}\|_1, \forall \gamma$$

for some $K \geq 0, \vartheta > 0$. Taking $\varrho = c \|\mathbf{X}\|_2$ results in scale invariance with respect to \mathbf{X} . Let's compare (63) with the restricted eigenvalue (RE) condition and the compatibility condition [7, 59]. For given \mathcal{J} , the two conditions assume that there exist positive numbers κ, ϑ_{RE} such that $J \|\mathbf{X}\gamma\|_2^2 \geq \kappa \|\gamma_{\mathcal{J}}\|_1^2$ (compatibility) or more restrictively, $\|\mathbf{X}\gamma\|_2^2 \geq \kappa \|\gamma_{\mathcal{J}}\|_2^2$ (RE), for all $\gamma : (1 + \vartheta_{RE}) \|\gamma_{\mathcal{J}}\|_1 \geq \|\gamma_{\mathcal{J}^c}\|_1$. Therefore, $(1 + \vartheta) \varrho \|\gamma_{\mathcal{J}}\|_1 \leq K \sqrt{J} \|\mathbf{X}\gamma\|_2 \vee \varrho \|\gamma_{\mathcal{J}^c}\|_1$ with $K = (1 + \vartheta_{RE})/(\varrho \sqrt{\kappa})$, $\vartheta = \vartheta_{RE}$. That is, the RE-type conditions are more demanding than (63) (and (60)). Another popular set of regularity conditions is based on restricted strong convexity (RSC). Under a version of RSC condition (and assuming f is differentiable), [36, Theorem 1] showed that $\|\tilde{\beta} - \beta^\|_2^2$ has a bound of order $\sigma^2 (J^* \log p)/n$ for any stationary point $\tilde{\beta}$. In the lasso case, the condition becomes $\|\mathbf{X}\gamma\|_2^2 \geq \alpha \|\gamma\|_2^2 - \tau \log p \|\gamma\|_1^2/n$ for some constant $\alpha > 0$ and $\tau \geq 0$, from which it follows that for any $\gamma : (1 + \vartheta_{RE}) \|\gamma_{\mathcal{J}}\|_1 \geq \|\gamma_{\mathcal{J}^c}\|_1$, $\|\mathbf{X}\gamma\|_2^2 \geq \alpha \|\gamma\|_2^2 - \tau (2 + \vartheta_{RE})^2 \frac{\log p}{n} \|\gamma_{\mathcal{J}}\|_1^2 \geq \alpha \|\gamma\|_2^2 - \tau (2 + \vartheta_{RE})^2 \frac{J \log p}{n} \|\gamma_{\mathcal{J}}\|_2^2 \geq \kappa' \|\gamma_{\mathcal{J}}\|_2^2$, where $\kappa' = \alpha - \tau (2 + \vartheta_{RE})^2 (J \log p)/n$. Therefore, when $n \gg J \log p$, RSC implies RE and so is more restrictive than (63). See Remark A.1 in Section A.12 for an extension to general penalties.*

REMARK 3 (Technical treatment). *A big difference between our work and [36] is that the latter enforces an ℓ_1 -type side constraint, for example, $\|\beta\|_1 \leq R$, in addition to the sparsity-inducing penalty P . The use of the constraint is a necessary ingredient of the proofs and the constraint parameter R appears in the minimum sample size condition and the error bounds implicitly. However, few practically used algorithms seem to include such an additional ℓ_1 constraint.*

Our analysis does not need any side constraint, and the resulting error bounds and the oracle inequality hold with no minimum sample size requirement. In fact, in dealing with a general penalty that may be nonconvex, our treatment of the stochastic term is distinctive from the conventional “ ℓ_1 fashion” via Hölder’s inequality: $\langle \epsilon, \mathbf{X}\beta \rangle \leq \|\mathbf{X}^\top \epsilon\|_\infty \|\beta\|_1$ (see, e.g., [10, 7, 37]). More concretely, applying the union bound to $\|\mathbf{X}^\top \epsilon\|_\infty$ will lead to a further upper bound $\|\beta\|_2^2 + P(\beta; \lambda)$ up to multiplicative factors [36], while we can bound $\langle \epsilon, \mathbf{X}\beta \rangle$ by the sum of $\|\mathbf{X}\beta\|_2^2/a$ and a light penalty $P_H(\beta; \lambda)/b$ for any $a, b > 0$, with a proper choice of λ .

REMARK 4 (Fixed points vs. local minimizers). *Targeting at the fixed points of the Bregman surrogate instead of the local minimizers of the original objective seems more reasonable from a statistical perspective. Certainly, if f is smooth, \mathcal{F} contains more valid solutions (cf. Lemma 4). But a more important reason is that \mathcal{F} can adaptively exclude bad local solutions for some statistical learning problems with severe nonsmoothness and nonconvexity.*

For instance, each bridge ℓ_q -penalty ($q : 0 \leq q < 1$) [22] determines a thresholding Θ_q , which is however the solution for infinitely many penalties; picking the particular one constructed from (48) that is the lowest and directionally differentiable [49], one can repeat the analysis in Theorems 3, 4 to show provable guarantees for all the fixed points of the iterative Θ_q procedure. In contrast, as pointed out by [36], the original optimization problem may contain “faulty” local minimizers. In fact, when $q = 0$, the ℓ_0 -penalized problem $\min_\beta \|\mathbf{X}\beta - \mathbf{y}\|_2^2/2 + (\lambda^2/2)\|\beta\|_0$ (not directionally differentiable) always has $\mathbf{0}$ as a local minimizer which is however a poor estimator as β^ is large. Switching to the surrogate’s fixed points successfully addresses the issue: $\hat{\beta} = \mathbf{0}$ is a valid fixed point only when $\mathbf{X}^\top \mathbf{y}$ is properly small: $\|\mathbf{X}^\top \mathbf{y}\|_\infty \leq \lambda$, or the true signal is inconsequential relative to the maximum noise level.*

3.2.2. *Statistical analysis of the iterates from Bregman surrogates.* We show a nice result for (52) in the composite setting: under a regularity condition similar to those in Section 3.2.1, with high probability, the t -th iterate can approach the statistical target within the desired precision geometrically fast, even when $p > n$. Specifically, we add a mild multiple of Δ_ψ to the left-hand side of (57) and assume that for some $\delta > 0$, $\varepsilon > 0$, $\vartheta > 0$ and large $K \geq 0$,

$$(64) \quad \begin{aligned} & \varepsilon \Delta_\psi(\beta^*, \beta) + \delta \mathbf{D}_2(\mathbf{X}\beta, \mathbf{X}\beta^*) + \vartheta P_H(\varrho(\beta - \beta^*); \lambda) + P_\Theta(\varrho\beta^*; \lambda) \\ & \leq (2\bar{\Delta}_l - \varrho^2 \mathcal{L}_\Theta \mathbf{D}_2)(\beta, \beta^*) + P_\Theta(\varrho\beta; \lambda) + K\lambda^2 J(\beta^*), \forall \beta \end{aligned}$$

and ψ is differentiable for simplicity. Recall that (39) in Proposition 2 requires $2\bar{\Delta}_f$ to dominate $\varepsilon \Delta_\psi$; (64) gives a large- p extension of it.

THEOREM 5. *Under the above regularity condition, for $\lambda = A\sigma\sqrt{\log(ep)}/\sqrt{(\delta \wedge \vartheta)\vartheta}$ with A sufficiently large and $\kappa = 1/(1 + \varepsilon)$, we have*

$$(65) \quad \Delta_\psi(\beta^*, \beta^{(t)}) \leq \kappa^t \Delta_\psi(\beta^*, \beta^{(0)}) + \frac{\kappa}{1 - \kappa} (K\lambda^2 J^* - \min_{1 \leq s \leq t} \Delta_\psi(\beta^{(s)}, \beta^{(s-1)}))$$

for any $t \geq 1$ with probability at least $1 - Cp^{-cA^2}$, where C, c are universal positive constants.

The error measure $\Delta_\psi(\beta^*, \beta^{(t)})$ in (65) has β^* as its first argument and differs from the $\Delta_t(\hat{\beta}, \beta^*)$ used in (61). According to the proof, (64) only needs to hold for $\beta = \beta^{(s)}$ ($0 \leq s \leq t$), and so different starting values may give different values of κ . With $\Delta_\psi \geq 0$ (which can be realized by stepsize control), the fast converging statistical error to $\mathcal{O}(\sigma^2 J^* \log(ep))$ implies that over-optimization may be unnecessary. As an example, consider the iterative thresholding procedures with $\Delta_t \leq LD_2$ and $\varrho^2 > L$. Then (65) yields

$$\|\beta^* - \beta^{(t)}\|_2^2 \leq \kappa^t \frac{\varrho^2}{\varrho^2 - L} \|\beta^* - \beta^{(0)}\|_2^2 + \frac{2\kappa K}{(1 - \kappa)(\varrho^2 - L)} \lambda^2 J^*.$$

So it is possible to terminate the iterative algorithm before full computational convergence without sacrificing much statistical accuracy. The simulations in Section C.2 support this point.

REMARK 5. *Theorem 5 reveals the fast decay of the direct statistical error between $\beta^{(t)}$ and β^* . [1] and [36] argued a similar point for gradient descent type algorithms, in a somewhat indirect manner: (i) $\beta^{(t)}$ can approach any globally optimal solution $\hat{\beta}$ geometrically fast in computation under a combination of an RSC condition and an RSM condition, and (ii) under some regularity conditions, every local minimum point is close enough to the authentic β^* . In the RSC condition for (i), the factor preceding the dominant term $\bar{\Delta}_t$ is 1 (there are two different sets of RSC conditions used in Theorem 1 and Theorem 3 of [36], the factor α_1 in the second set corresponding to half of the α_1 used in the first set). But (64) allows it to be 2. Moreover, Theorem 5 does not need the extra RSM condition and applies to a broader class of algorithms. For example, we can show that the statistical error of the LLA algorithm reduces at a linear rate to the desired precision under some regularity conditions; see Proposition 5 and Lemma A.7 in Section A.16.*

4. Two acceleration schemes for generalized Bregman surrogates. How to accelerate first-order algorithms without incurring much additional cost per iteration has lately attracted lots of attention in big data applications. In convex optimization, Nesterov's momentum techniques prove to be quite effective in that the rate of convergence can be improved from $\mathcal{O}(1/t)$ to $\mathcal{O}(1/t^2)$, which is optimal when using first-order methods on smooth problems [41, 57, 4, 32]. This section attempts to extend Nesterov's *first* and *second* accelerations [39, 40] to Bregman-surrogate algorithms. With a possible lack of smoothness or convexity, carefully choosing the relaxation parameters and step sizes is the key, and we will see the benefit of maximizing a quantity $R_t/(\theta_t^2 \rho_t)$ at the t -th iteration, with R_t appropriately defined via generalized Bregman notation. We consider the following two broad scenarios to devise the acceleration schemes.

Scenario 1. $g(\beta; \gamma) = f(\beta) - \Delta_{\psi_0}(\beta, \gamma) + \rho D_2(\beta, \gamma)$. This surrogate family includes gradient descent type algorithms. Often, if $\min_\beta f(\beta) + \Delta_\psi(\beta, \gamma)$ is easy to solve, so is $\min_\beta f(\beta) + \Delta_\psi(\beta, \gamma) + \rho D_2(\beta, \gamma)$, in which case $\psi_0 = -\psi$.

Scenario 2. $g(\beta; \gamma) = f(\beta) - \Delta_{\psi_0}(\beta, \gamma) + \rho \Delta_\phi(\beta, \gamma)$. This gives a more general class than the first one.

This section assumes that $f, \psi_0, \phi, \Delta_{\psi_0}(\cdot, \gamma), \Delta_{\psi_0}(\cdot, \gamma)$ are directionally differentiable given any γ . We introduce a convenient notation \mathbf{C}_ψ defined for any ψ as follows

$$(66) \quad \mathbf{C}_\psi(\alpha, \beta, \theta) = \theta\psi(\alpha) + (1 - \theta)\psi(\beta) - \psi(\theta\alpha + (1 - \theta)\beta),$$

where $0 \leq \theta \leq 1$. Like Δ , \mathbf{C} is a linear operator of ψ and its nonnegativity means convexity. Some connections between Δ and \mathbf{C} are given below.

LEMMA 5. *Let ψ be directionally differentiable. (i) $\mathbf{C}_\psi(\alpha, \beta, \theta) = (1 - \theta)\Delta_\psi(\beta, \alpha) - \Delta_\psi(\theta\alpha + (1 - \theta)\beta, \alpha)$ for any α, β and $\theta \in [0, 1]$. (ii) $\mathbf{C}_{\Delta_\psi(\cdot, \alpha)} = \mathbf{C}_\psi$ if ψ is differentiable at α .*

An acceleration scheme of the second kind. Scenario 2 is of our primary interest since it applies more broadly. Below, we modify the surrogate and define an iterative algorithm (not a descent method) that involves three sequences $\alpha^{(t)}, \beta^{(t)}, \gamma^{(t)}$ starting at $\alpha^{(0)} = \beta^{(0)}$:

$$(67a) \quad \gamma^{(t)} = (1 - \theta_t)\beta^{(t)} + \theta_t\alpha^{(t)}$$

$$(67b) \quad \alpha^{(t+1)} = \operatorname{argmin} f(\beta) - \Delta_{\psi_0}(\beta, \gamma^{(t)}) + \mu_0\Delta_\phi(\beta, \gamma^{(t)}) + \theta_t\rho_t\Delta_\phi(\beta, \alpha^{(t)})$$

$$(67c) \quad \beta^{(t+1)} = (1 - \theta_t)\beta^{(t)} + \theta_t\alpha^{(t+1)},$$

for some $\mu_0 \geq 0$, $\theta_t \in (0, 1]$, $\rho_t > 0$ ($\forall t \geq 0$), to be chosen later. Notice the extra GBF term $\mu_0\Delta_\phi(\cdot, \gamma^{(t)})$ in (67b) in addition to $\Delta_\phi(\cdot, \alpha^{(t)})$. The design of relaxation parameters θ_t and inverse step size parameters ρ_t, μ_0 holds the key to acceleration. Let

$$(68) \quad \bar{\psi}_0 = \psi_0 - \mu_0\phi.$$

We advocate the following line search criterion

$$(69a) \quad R_t := \theta_t^2\rho_t\Delta_\phi(\alpha^{(t+1)}, \alpha^{(t)}) - \Delta_{\bar{\psi}_0}(\beta^{(t+1)}, \gamma^{(t)}) + (1 - \theta_t)\Delta_{\bar{\psi}_0}(\beta^{(t)}, \gamma^{(t)}) \\ + \mathbf{C}_{f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})}(\alpha^{(t+1)}, \beta^{(t)}, \theta_t) \geq 0,$$

$$(69b) \quad \frac{\theta_t^2}{1 - \theta_t} = \frac{\theta_{t-1}(\rho_{t-1}\theta_{t-1} + \mu_0)}{\rho_t}, \quad t \geq 1.$$

The update of the relaxation parameter involves ρ and μ as well.

Theorem 6 presents two error bounds without assuming convexity or smoothness, and shows in general the reasonability of (69a).

THEOREM 6. *Let ρ_t be any positive sequence. Consider the algorithm defined by (67a)–(67c) and (69b). Let $\mathcal{E}_t(\beta) = \Delta_{\bar{\psi}_0}(\beta, \gamma^{(t)}) + \Delta_{f(\cdot) - \Delta_{\psi_0}(\cdot, \gamma^{(t)})}(\beta, \alpha^{(t+1)}) + (\mu_0\Delta_{\Delta_\phi(\cdot, \gamma^{(t)}) - \phi(\cdot)} + \theta_t\rho_t\Delta_{\Delta_\phi(\cdot, \alpha^{(t)}) - \phi(\cdot)})(\beta, \alpha^{(t+1)})$.*

(i) *When $\mu_0 = 0$, for any β and $T \geq 0$,*

$$(70) \quad \frac{f(\beta^{(T+1)}) - f(\beta)}{\theta_T^2\rho_T} + T \cdot \operatorname{avg}_{0 \leq t \leq T} \frac{\mathcal{E}_t(\beta)}{\theta_t\rho_t} + T \cdot \operatorname{avg}_{0 \leq t \leq T} \frac{R_t}{\theta_t^2\rho_t} \\ \leq \Delta_\phi(\beta, \alpha^{(0)}) - \Delta_\phi(\beta, \alpha^{(T+1)}) + \frac{1 - \theta_0}{\theta_0^2\rho_0} [f(\beta^{(0)}) - f(\beta)].$$

(ii) *Moreover, given any $\mu_0 \geq 0$,*

$$(71) \quad f(\beta^{(T+1)}) - f(\beta) + \theta_T^2(\rho_T + \frac{\mu_0}{\theta_T})\Delta_\phi(\beta, \alpha^{(T+1)}) \\ + \sum_{t=0}^T (\Pi_{s=t+1}^T (1 - \theta_s))(R_t + \theta_t\mathcal{E}_t(\beta)) \\ \leq \left(\prod_{t=1}^T (1 - \theta_t) \right) [(1 - \theta_0)(f(\beta^{(0)}) - f(\beta)) + \theta_0^2\rho_0\Delta_\phi(\beta, \beta^{(0)})]$$

for all β and $T \geq 0$, where by convention, $\prod_{s=l}^u a_s = 1$ as $l > u$.

First, we make a discussion of the results for convex optimization. Assume $\Delta_\phi \geq \sigma \mathbf{D}_2$ for some $\sigma > 0$. With the additional knowledge that $f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})$ is convex and $\Delta_{\bar{\psi}_0} \leq L_{\bar{\psi}_0} \mathbf{D}_2$ for some $L_{\bar{\psi}_0} \geq 0$, (69a) is implied by

$$(72) \quad \theta_t^2(\rho_t - L_{\bar{\psi}_0}/\sigma)\Delta_\phi(\alpha^{(t+1)}, \alpha^{(t)}) + (1 - \theta_t)\Delta_{\bar{\psi}_0}(\beta^{(t)}, \gamma^{(t)}) \geq 0.$$

So when f is convex, criterion (69) is satisfied by $\rho_t = \rho \geq L_{\bar{\psi}_0}/\sigma$, $\psi_0 = f$, $\mu_0 = 0$ and $\theta_{t+1} = (\sqrt{\theta_t^4 + 4\theta_t^2} - \theta_t^2)/2$, degenerating to Nesterov's second method [40, 57], and the convergence rate is of order $\mathcal{O}(1/T^2)$ according to (70) and (75). The second conclusion tells more when strong convexity (or restricted strong convexity) arises. Given a convex f satisfying $\mu \mathbf{D}_\phi \leq \Delta_f \leq L \mathbf{D}_\phi$ with $0 < \mu \leq L$ and ϕ differentiable, taking $\psi_0 = f$, $\mu_0 = \mu$, and $\rho_t = L - \mu$ ensures $\mathcal{E}_t(\beta) = \Delta_{f-\mu\phi}(\beta, \gamma^{(t)}) + \Delta_{f(\cdot) - \Delta_f(\cdot, \gamma^{(t)})}(\beta, \alpha^{(t+1)}) \geq \Delta_{f-\mu\phi}(\beta, \gamma^{(t)}) \geq 0$ and $R_t \geq \theta_t^2 \rho_t \mathbf{D}_\phi(\alpha^{(t+1)}, \alpha^{(t)}) - \Delta_{\bar{\psi}_0}(\beta^{(t+1)}, \gamma^{(t)}) \geq \theta_t^2(\rho_t + \mu_0 - L)\sigma \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) = 0$. According to (69b), the following choice

$$(73) \quad \theta_t = \theta_0 = \frac{2}{\sqrt{4\kappa - 3} + 1} \text{ with } \kappa = L/\mu$$

suffices, and the optimization problem to solve in (67b) becomes

$$(74) \quad \min f(\gamma^{(t)}) + \delta f(\beta; \beta - \gamma^{(t)}) + \mu \mathbf{D}_\phi(\beta, \gamma^{(t)}) + \frac{2(L - \mu)}{\sqrt{4\kappa - 3} + 1} \mathbf{D}_\phi(\beta, \alpha^{(t)}).$$

From (71), both $f(\beta^{(T+1)}) - f(\beta)$ and $\mathbf{D}_\phi(\beta, \alpha^{(T+1)})$ enjoy a linear convergence with rate parameter $\frac{\sqrt{4\kappa - 3} - 1}{\sqrt{4\kappa - 3} + 1}$, or an iteration complexity of $\mathcal{O}(\sqrt{\kappa} \log(1/\epsilon))$, significantly faster than $\mathcal{O}(\kappa \log(1/\epsilon))$ in Proposition 2. Hence (67), (69) can achieve rate-optimality in various convex scenarios. To the best of our knowledge, this is the first “all-in-one” form of the *second* acceleration that adapts.

The proposed algorithm can even go beyond convexity. As a demonstration, let us apply the acceleration to the iterative quantile-thresholding procedure (cf. Example 2) for solving the feature screening problem: $\min l(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2/2$ s.t. $\|\beta\|_0 \leq q$, which is nonconvex. Here, q is bounded above by p but may be larger than n . Take $\phi = \|\cdot\|_2^2/2$, $\mu_0 = 0$ and $\psi_0(\beta) = l(\beta) - \mathcal{L}\phi(\beta)$ for some $\mathcal{L} \geq 0$. Given any $s \leq p$ and \mathbf{X} , define the restricted isometry number $\rho_+(s)$ [12] that satisfies $\|\mathbf{X}\beta\|_2^2 \leq \rho_+(s)\|\beta\|_2^2$, $\forall \beta : \|\beta\|_0 \leq s$, which can be much smaller than $\|\mathbf{X}\|_2^2$ as s is small.

COROLLARY 2. *Assume q is set larger than the target $\|\beta^*\|_0$ with the ratio denoted by r . Then for any $\mathcal{L} \geq \rho_+(2q)/\sqrt{r}$, there exists a universal ρ_t ($\rho_t = \rho_+(2q)(1 - 1/\sqrt{r})$, say), thereby $\theta_{t+1} = (\sqrt{\theta_t^4 + 4\theta_t^2} - \theta_t^2)/2$, such that the accelerated iterative quantile-thresholding according to (67a)–(67c) satisfies $l(\beta^{(T+1)}) - l(\beta^*) + \min_{0 \leq t \leq T} \Delta_{\psi_0}(\beta^*, \gamma^{(t)}) \leq A/T^2$ for all $T \geq 0$, where A is independent of T .*

The proof of the corollary shows the power of an *accumulative* R_t -control, and applies more generally: if the objective function $f(\beta)$, possibly nonconvex, can be written as the sum of a convex function $l(\beta)$ with $\Delta_l \leq L \mathbf{D}_2$ and a function $P(\beta)$ that can be lifted: $\Delta_P + \mathcal{L}_0 \mathbf{D}_2 \geq 0$ for some finite $\mathcal{L}_0 \geq 0$, then one can utilize a ψ_0 as $l - 0.6\mathcal{L}_0 \|\cdot\|_2^2$ and a universal ρ_t to fulfill $T \cdot \text{avg}_{t \leq T} R_t / (\theta_t^2 \rho_t) \geq 0$ in (70) (although not every R_t is necessarily nonnegative) so as to attain an $\mathcal{O}(1/T^2)$ error bound. See Remark A.3 in Section A.14.

Of course, a time-varying ρ_t can provide finer control, and the theorem does not limit ρ_t to be constant. In fact, under $\mu_0 = 0$, as long as $\rho_t/\rho_{t-1} \geq 1 - (at + ab + 1)/(t + b - 1)^2$ ($t \geq 1$) for some constants a, b : $a > -2, b \geq a + 1$, induction based on (69b) gives $\theta_t \leq (a + 2)/(t + b)$

and $\sum_{t=0}^T \rho_t / (\rho_t \theta_t) \geq (T + c_1)^2 / (a + 2)^2 + c_2$ (with constants c_i dependent on a, b) for any $t \geq 1$, from which it follows that

$$(75) \quad \theta_T^2 = \mathcal{O}(1/T^2) \quad \text{and} \quad T \cdot \text{avg}_{0 \leq t \leq T} (1/(\rho_t \theta_t)) \geq \mathcal{O}(T^2/\rho_T).$$

Now, under $R_t \geq 0$ or just $\sum_{t=0}^T R_t / (\theta_t^2 \rho_t) \geq 0$, (70) gives $f(\beta^{(T+1)}) - f(\beta) + \min_{0 \leq t \leq T} \mathcal{E}_t(\beta) \leq \mathcal{O}(\rho_T/T^2)$ for any β . Typically, (69a) involves a line search. If the condition fails for the current value of ρ_t , one can set $\rho_t = \alpha \rho_t$ for some $\alpha > 1$ and recalculate θ_t , $\gamma^{(t)}$, $\alpha^{(t+1)}$ and $\beta^{(t+1)}$ according to (69b) and (67) to verify it again. In implementation, it is wise to limit the number of searches at each iteration (denoted by M) to control the per-iteration complexity. If (69a) does not hold after m times of search, we simply pick the ρ_t that gives the largest $R_t / (\theta_t^2 \rho_t)$ based on Theorem 6. Some details are in Algorithm B.1. In simulation studies, letting $M = 3$, $\alpha = 2$ already shows excellent performance; see Figure C.5 and Figure C.6.

An acceleration scheme of the first kind. For the algorithms falling into Scenario 1, we can alternatively consider two sequences of iterates generated by

$$(76a) \quad \gamma^{(t)} = \beta^{(t)} + \{\rho_{t-1}\theta_t(1 - \theta_{t-1})/(\rho_{t-1}\theta_{t-1} + \mu_0)\}(\beta^{(t)} - \beta^{(t-1)}),$$

$$(76b) \quad \beta^{(t+1)} = \arg \min f(\beta) - \Delta_{\psi_0}(\beta, \gamma^{(t)}) + \mu_0 \mathbf{D}_2(\beta, \gamma^{(t)}) + \rho_t \mathbf{D}_2(\beta, \gamma^{(t)}),$$

for some $\mu_0 \geq 0$, $\theta_t \in (0, 1]$, $\rho_t > 0$ for all $t \geq 0$, and we force $\gamma^{(0)} = \beta^{(0)}$. (76a), (76b) give a new first type acceleration, and notably, the novel update of $\gamma^{(t)}$ involves ρ_{t-1} . When $\beta^{(t+1)} = \gamma^{(t)}$ one stops the algorithm and obtains a fixed point with provable statistical guarantees as shown in Section 3.2.1.

Similar to (68), let $\bar{\psi}_0 = \psi_0 - \mu_0 \|\cdot\|_2^2/2$. Define the line search criterion

$$(77a) \quad R_t := (\rho_t \mathbf{D}_2 - \Delta_{\bar{\psi}_0})(\beta^{(t+1)}, \gamma^{(t)}) + (1 - \theta_t) \Delta_{\bar{\psi}_0}(\beta^{(t)}, \gamma^{(t)}) \geq 0,$$

$$(77b) \quad \frac{\theta_t^2}{1 - \theta_t} = \frac{\theta_{t-1}(\rho_{t-1}\theta_{t-1} + \mu_0)}{\rho_t}, \quad \theta_t \geq 0, \quad \rho_t > 0, \quad t \geq 1.$$

Note that R_t is defined differently from (69a). The following theorem reveals the importance of maximizing R_t in each iteration step when performing possibly nonconvex optimization.

THEOREM 7. *Given any $\rho_t > 0$ ($t \geq 0$), consider the algorithm defined by (76a), (76b), and (77b). Let $\mathcal{E}_t(\beta) = \Delta_{\bar{\psi}_0}(\beta, \gamma^{(t)}) + \{\mathbf{C}_{f(\cdot) - \Delta_{\psi_0}(\cdot, \gamma^{(t)})}(\beta, \beta^{(t)}, \theta_t) + \Delta_{f(\cdot) - \Delta_{\psi_0}(\cdot, \gamma^{(t)})}(\theta_t \beta + (1 - \theta_t)\beta^{(t)}, \beta^{(t+1)})\}/\theta_t$.*

(i) *When $\mu_0 = 0$, we have*

$$\begin{aligned} & \frac{f(\beta^{(T+1)}) - f(\beta)}{\theta_T^2 \rho_T} + T \cdot \text{avg}_{0 \leq t \leq T} \frac{\mathcal{E}_t(\beta)}{\theta_t \rho_t} + T \cdot \text{avg}_{0 \leq t \leq T} \frac{R_t}{\theta_t^2 \rho_t} \\ & \leq \mathbf{D}_2(\beta, \beta^{(0)}) + \frac{1 - \theta_0}{\theta_0^2 \rho_0} [f(\beta^{(0)}) - f(\beta)] \quad \text{for any } \beta \text{ and } T \geq 0. \end{aligned}$$

(ii) *Moreover, given any $\mu_0 \geq 0$, for all β and $T \geq 0$,*

$$\begin{aligned} & f(\beta^{(T+1)}) - f(\beta) + \theta_T^2 (\rho_T + \frac{\mu_0}{\theta_T}) \mathbf{D}_2(\beta, (\gamma^{(T+1)} - (1 - \theta_{T+1})\beta^{(T+1)})/\theta_{T+1}) \\ & + \sum_{t=0}^T (\Pi_{s=t+1}^T (1 - \theta_s)) (R_t + \theta_t \mathcal{E}_t(\beta)) \\ & \leq \left(\prod_{t=1}^T (1 - \theta_t) \right) [(1 - \theta_0)(f(\beta^{(0)}) - f(\beta)) + \theta_0^2 \rho_0 \mathbf{D}_2(\beta, \beta^{(0)})]. \end{aligned}$$

Again, the new proposal of the iterate and parameter updates adapts to various situations, with μ_0 (which can be a sequence μ_t , cf. Remark A.2) measuring the degree of convexity (or restricted convexity in a nonconvex composite problem). For example, when f is convex and L -strongly smooth, $\mu_0 = 0$, $\rho_t = L$, $\psi_0 = f$, and $\theta_{t+1} = (\sqrt{\theta_t^4 + 4\theta_t^2} - \theta_t^2)/2$ make (77) hold, corresponding to Nesterov's first method. Interestingly, if f is μ -strongly convex, the associated standard momentum update $\gamma^{(t)} = \beta^{(t)} + \theta_t(\theta_{t-1}^{-1} - 1)(\beta^{(t)} - \beta^{(t-1)})$ only attains a linear rate at $1 - 1/\kappa$ ($\kappa = L/\mu$) (cf. Remark A.4), showing *no* theoretical advantage over the plain gradient descent. (76) fixes the issue: with $\mu_0 = \mu$, $\rho_t = L - \mu$, $\theta_t = 2/(\sqrt{4\kappa - 3} + 1)$, an accelerated linear rate parameter is obtained as $(\sqrt{4\kappa - 3} - 1)/(\sqrt{4\kappa - 3} + 1) (\leq 1 - \sqrt{3/(4\kappa)})$. (When μ_0 is unknown, (76b) based on the split $L = \rho_t + \mu_t$ is still advantageous over the classical acceleration with $\rho_t = L$.) We proved these error bounds by use of GBFs, which is perhaps more straightforward than Nesterov's ingenious proof based on the notion of estimate sequence, and more importantly, (76), (77) provide a universal “all-in-one” form, instead of separate schemes in different situations [41].

Theorem 7 accommodates diverse choices of the parameters $\psi_0, \mu_0, \rho_t, \theta_t$ and is motivating in the nonconvex composite setup. Consider, for example, $\min f(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2/2 + P_\Theta(\varrho\beta; \lambda)$. Because the objective is nonconvex when $p > n$ and $\mathcal{L}_\Theta > 0$, how to accelerate the associated iterative thresholding procedure is an unconventional problem. From the studies in Section 3.2, we have learned that a sparsity-inducing penalty with a properly large threshold to suppress the noise can result in strong convexity in a restricted sense. We can then use a surrogate $f(\beta) + (\rho\mathbf{D}_2 - \Delta_{\psi_0})(\beta, \beta^-)$ where $\psi_0(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2/2 - \varrho^2\mathcal{L}_\Theta\|\beta\|_2^2/2$ and $\mu_0 = 0$. Since $f(\cdot) - \Delta_{\psi_0}(\cdot, \gamma)$ is convex (cf. Lemma A.3), $\mathcal{E}_t(\beta) \geq \Delta_{\psi_0}(\beta, \gamma^{(t)})$. Moreover, thanks to the sparsity in $\beta^{(t)}$, and thus $\gamma^{(t)}$, $\mathbf{X}(\beta^{(t)} - \gamma^{(t)})$ involves just a small number of features. So with an incoherent design, a properly small ϱ can make $\Delta_{\psi_0}(\beta^{(t)}, \gamma^{(t)}) \geq 0$. Now, taking a constant ρ_t as large as, for instance, $\|\mathbf{X}\|_2^2 - \varrho^2\mathcal{L}_\Theta$, may yield a convergence rate of order $\mathcal{O}(1/t^2)$. (Actually, linear convergence may result from the restricted strong convexity under some regularity conditions.) More generally, different ρ_t 's are allowed in the theorem: (75) is still secured with just, say, $\rho_t/\rho_{t-1} \geq 1 - (t+3)/(t+1)^2$. A line search can be used to determine a proper sequence ρ_t ; see Algorithm B.2 for more details.

The proposed accelerations of the first kind and of the second kind can be utilized in a wide range of problems. Because they are momentum based, the original algorithms need not be substantially modified to have an improved iteration complexity, and the two theorems proved in this section apply in any dimensions with no design coherence restrictions. Another delightful fact is that our “all-in-one” forms update the iterates adaptively according to the degree of convexity $\mu_0 \geq 0$, which can be relaxed to a sequence of local measures μ_t (Remark A.2). With a line search to get properly large μ_t , this could be helpful in high dimensional sparse learning problems which may or may *not* have restricted strong convexity (the associated parameter often hard to determine in theory).

5. Summary. This paper studied the class of iterative algorithms derived from GBF-defined surrogates with a possible lack of convexity and/or smoothness. These surrogates differ from the MM surrogates frequently used in statistical computation, in that they gain additional first-order degeneracy and may drop the majorization requirement. GBFs have interesting connections to the densities in the exponential family and possess some idempotence properties that are useful for studying iterative algorithms.

The GBF calculus built by the lemmas not only facilitates optimization error analysis but can be bound to the empirical process theory for nonasymptotic statistical analysis (cf. Sections 3.2 and A.18). In addition to obtaining some insightful results in the realm of convex optimization, we were able to build universal global convergence rates for a broad class of

Bregman-surrogate algorithms for nonsmooth nonconvex optimization. Moreover, in the non-convex composite setting that is of great interest in high dimensional statistics, we found that the sequence of iterates generated by Bregman surrogates can approach the statistical truth at a linear rate even when $p > n$, and the obtained fixed points enjoy oracle inequalities with essentially the optimal order of statistical accuracy, under some regularity conditions less demanding than those used in the literature. Finally, we devised two “all-in-one” acceleration schemes with novel updates of the iterates and relaxation and stepsize parameters, and some sharp theoretical bounds were shown without assuming smoothness or convexity.

APPENDIX A: PROOFS

We list some notation and symbols that are used in the proofs. Given a directionally differentiable function ψ , $\Delta_\psi(\beta, \gamma) = \psi(\beta) - \psi(\gamma) - \delta\psi(\gamma; \beta - \gamma)$, $\bar{\Delta}_\psi(\beta, \gamma) = (\Delta_\psi(\beta, \gamma) + \Delta_\psi(\gamma, \beta))/2$, and $\tilde{\Delta}_\psi(\beta, \gamma) = \Delta_\psi(\gamma, \beta)$. We occasionally denote $\Delta_\psi(\beta, \gamma)$ by $\Delta(\beta, \gamma)$ when there is no ambiguity. The classes of continuous functions and continuously differentiable functions are denoted by \mathcal{C}^0 and \mathcal{C}^1 , respectively. Recall that all functions are assumed to be defined on a vector space unless otherwise mentioned.

DEFINITION A.1. We call ξ a sub-Gaussian random variable if and only if there exist constants $C, c > 0$ such that $\mathbb{P}\{|\xi| \geq t\} \leq Ce^{-ct^2}$, $\forall t > 0$. The scale (or ψ_2 -norm) of ξ is defined by $\sigma(\xi) = \inf\{\sigma > 0 : \mathbb{E} \exp(\xi^2/\sigma^2) \leq 2\}$. More generally, $\xi \in \mathbb{R}^p$ is called a sub-Gaussian random vector with scale bounded by σ if all one-dimensional marginals $\langle \xi, \alpha \rangle$ are sub-Gaussian satisfying $\|\langle \xi, \alpha \rangle\|_{\psi_2} \leq \sigma \|\alpha\|_2$, $\forall \alpha \in \mathbb{R}^p$.

DEFINITION A.2. We call d a pseudo-metric if it satisfies $d(\eta_1, \eta_2) = d(\eta_2, \eta_1) \geq 0$ and $d(\eta_1, \eta_2) \leq d(\eta_1, \eta_3) + d(\eta_2, \eta_3)$, for all η_1, η_2, η_3 .

We state a first-order optimality condition satisfied by all local minimizers of f that is directionally differentiable. The result is basic and we omit the proof. It holds the key to deriving the so-called “basic inequality” in a variety of statistical learning problems.

LEMMA A.1. Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a real-valued function and $C \subset \mathbb{R}^p$ be a convex set. Suppose that f is directionally differentiable at β^o that is a local minimizer to the problem $\min_{\beta \in C} f(\beta)$. Then $\delta f(\beta^o; \mathbf{h}) \geq 0$ with $\mathbf{h} = \beta - \beta^o$ or $f(\beta) - f(\beta^o) \geq \Delta_f(\beta, \beta^o)$ for all $\beta \in C$.

A.1. Proof of Lemma 1. (i) This property is straightforward by definition:

$$\begin{aligned} & \Delta_{a\psi+b\varphi}(\beta, \gamma) \\ &= (a\psi + b\varphi)(\beta) - (a\psi + b\varphi)(\gamma) - \delta(a\psi + b\varphi)(\gamma; \beta - \gamma) \\ &= a[\psi(\beta) - \psi(\gamma) - \delta\psi(\gamma; \beta - \gamma)] + b[\varphi(\beta) - \varphi(\gamma) - \delta\varphi(\gamma; \beta - \gamma)] \\ &= a\Delta_\psi(\beta, \gamma) + b\Delta_\varphi(\beta, \gamma). \end{aligned}$$

(ii) From [45, Theorem 23.1], the convexity of ψ implies the directional differentiability of ψ and the positively homogenous convexity of $\delta\psi(\beta; \cdot)$ for any given β , and we can write

$$(A.1) \quad \delta\psi(\beta; \mathbf{h}) = \inf_{\epsilon > 0} \frac{\psi(\beta + \epsilon \mathbf{h}) - \psi(\beta)}{\epsilon}.$$

Putting $\epsilon = 1$ and $\mathbf{h} = \gamma - \beta$ in (A.1) gives $\delta\psi(\beta; \gamma - \beta) \leq \psi(\gamma) - \psi(\beta)$, thus $\Delta_\psi(\gamma, \beta) \geq 0$.

Conversely, suppose that ψ defined on \mathbb{R}^n is directionally differentiable ($\delta\psi$ exists and is finite), thus radially continuous, and $\Delta_\psi \geq 0$. For any $s_\beta : \langle s_\beta, \mathbf{h} \rangle \leq \delta\psi(\beta; \mathbf{h})$, $s_\gamma : \langle s_\gamma, \mathbf{h} \rangle \leq \delta\psi(\gamma; \mathbf{h}) \forall \mathbf{h}$,

$$(A.2) \quad \psi(\beta) - \psi(\gamma) - \langle s_\gamma, \beta - \gamma \rangle \geq \Delta_\psi(\beta, \gamma) \geq 0,$$

$$(A.3) \quad \psi(\gamma) - \psi(\beta) - \langle s_\beta, \gamma - \beta \rangle \geq \Delta_\psi(\gamma, \beta) \geq 0.$$

Adding them together gives $\langle s_\beta - s_\gamma, \beta - \gamma \rangle \geq 0$. This indicates the monotone property of the Clarke-Rockafellar subdifferential of ψ , thereby its convexity according to [16].

(iii) To show the first result, notice that $\Delta_{\psi \circ \varphi}(\beta, \gamma) - \Delta_\psi(\varphi(\beta), \varphi(\gamma)) = \lim_{\epsilon \rightarrow 0+} \{\psi(\varphi(\gamma) + \epsilon(\varphi(\beta) - \varphi(\gamma))) - \psi(\varphi(\gamma))\} / \epsilon = \langle \nabla \psi(\varphi(\gamma)), \varphi(\beta) - \varphi(\gamma) \rangle - \delta(\psi \circ \varphi)(\gamma; \beta - \gamma)$. From $\psi \in \mathcal{C}^1$ and $\varphi \in \mathcal{C}^0$,

$$\begin{aligned} \delta(\psi \circ \varphi)(\gamma; \beta - \gamma) &= \lim_{\epsilon \rightarrow 0+} \{\psi(\varphi(\gamma) + (\varphi(\gamma + \epsilon(\beta - \gamma)) - \varphi(\gamma))) - \psi(\varphi(\gamma))\} / \epsilon \\ &= \lim_{\epsilon \rightarrow 0+} \langle \nabla \psi(\varphi(\gamma)), \varphi(\gamma + \epsilon(\beta - \gamma)) - \varphi(\gamma) \rangle / \epsilon \\ &= \langle \nabla \psi(\varphi(\gamma)), \delta\varphi(\gamma; \beta - \gamma) \rangle. \end{aligned}$$

Using the definition of $\Delta_\varphi(\beta, \gamma)$ (the componentwise extension), we obtain the conclusion.

Next, we prove the second result. Let $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^n$ be the linear function $\varphi(\beta) = \mathbf{X}\beta + \alpha$ with its Jacobian matrix $D\varphi(\beta) := [D_j\varphi_i(x)] = \mathbf{X} \in \mathbb{R}^{n \times p}$. By definition, $\Delta_{\psi \circ \varphi}(\beta, \gamma) = \psi(\varphi(\beta)) - \psi(\varphi(\gamma)) - \delta(\psi \circ \varphi)(\gamma; \beta - \gamma)$ and

$$\begin{aligned} \delta(\psi \circ \varphi)(\gamma; \beta - \gamma) &= \lim_{\epsilon \rightarrow 0+} \{\psi(\varphi(\gamma + \epsilon(\beta - \gamma))) - \psi(\varphi(\gamma))\} / \epsilon \\ &= \lim_{\epsilon \rightarrow 0+} \{\psi(\varphi(\gamma) + \epsilon D\varphi(\gamma)(\beta - \gamma)) - \psi(\varphi(\gamma))\} / \epsilon \\ &= \delta\psi(\varphi(\gamma); D\varphi(\gamma)(\beta - \gamma)) = \delta\psi(\varphi(\gamma); \varphi(\beta) - \varphi(\gamma)), \end{aligned}$$

from which it follows that $\Delta_{\psi \circ \varphi}(\beta, \gamma) = \Delta_\psi(\varphi(\beta), \varphi(\gamma))$.

(iv) From Theorem 11 in [25], for any continuous function f with finite Dini derivative $D^+f(x) := \limsup_{\epsilon \rightarrow 0+} (f(x + \epsilon) - f(x)) / \epsilon$, if $D^+f(x)$ is integrable over $[a, b]$, $f(b) - f(a) = \int_a^b D^+f(x) dx$. By definition, ψ is continuous when restricted to the line segment $[\beta, \gamma]$ (radial continuity). It follows that

$$\begin{aligned} \psi(\beta) - \psi(\gamma) &= \psi(\gamma + t(\beta - \gamma)) \Big|_{t=0}^1 \\ &= \int_0^1 \lim_{\epsilon \rightarrow 0+} \frac{1}{\epsilon} [\psi(\gamma + (t + \epsilon)(\beta - \gamma)) - \psi(\gamma + t(\beta - \gamma))] dt \\ &= \int_0^1 \lim_{\epsilon \rightarrow 0+} \frac{1}{\epsilon} [\psi(\gamma + t(\beta - \gamma) + \epsilon(\beta - \gamma)) - \psi(\gamma + t(\beta - \gamma))] dt \\ &= \int_0^1 \delta\psi(\gamma + t(\beta - \gamma); \beta - \gamma) dt. \end{aligned}$$

Hence, Δ_ψ can be formulated by

$$\Delta_\psi(\beta, \gamma) = \int_0^1 [\delta\psi(\gamma + t(\beta - \gamma); \beta - \gamma) - \delta\psi(\gamma; \beta - \gamma)] dt.$$

A.2. Proof of Lemma 2. (i) First, if $\delta\psi(\alpha; \cdot - \alpha)$ is directionally differentiable, then

$$(A.4) \quad \Delta_\psi(\beta, \gamma) - \Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) = \Delta_{\delta\psi(\alpha; \cdot - \alpha)}(\beta, \gamma)$$

for any α, β, γ . In fact, $\Delta_\psi(\beta, \gamma) - \Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) = \Delta_{\psi(\cdot) - \Delta_\psi(\cdot, \alpha)}(\beta, \gamma) = \Delta_{\psi(\alpha) + \delta\psi(\alpha; \cdot - \alpha)}(\beta, \gamma) = \Delta_{\delta\psi(\alpha; \cdot - \alpha)}(\beta, \gamma)$.

Accordingly, when ψ is convex, which means $\delta\psi(\alpha; \cdot - \alpha)$ is convex as well (cf. Section A.1), $\Delta_{\delta\psi(\alpha; \cdot - \alpha)}(\beta, \gamma) \geq 0$ by Lemma 1. The result under concavity can be similarly proved.

(ii) Let

$$(A.5) \quad q(\cdot; \alpha) = \delta\psi(\alpha; \cdot - \alpha).$$

We want to show for $\alpha = \theta\beta + (1 - \theta)\gamma$ with $\theta \leq 0$ or $\theta \geq 1$, $\Delta_{q(\cdot; \alpha)}(\beta, \gamma)$ is well-defined and equals 0. This is intuitive due to the linearity of q when restricted to $[\beta, \gamma]$, assuming $\beta - \alpha$ and $\gamma - \alpha$ are positively collinear.

To verify it, by definition,

$$\begin{aligned} \delta q(\cdot; \alpha)(\gamma; \beta - \gamma) &= \lim_{\epsilon \rightarrow 0+} [q(\gamma + \epsilon(\beta - \gamma); \alpha) - q(\gamma; \alpha)]/\epsilon \\ &= \lim_{\epsilon \rightarrow 0+} [\delta\psi(\alpha; \gamma + \epsilon(\beta - \gamma) - \alpha) - \delta\psi(\alpha; \gamma - \alpha)]/\epsilon \\ &= \lim_{\epsilon \rightarrow 0+} [\delta\psi(\alpha; (\theta - \epsilon)(\gamma - \beta)) - \delta\psi(\alpha; \theta(\gamma - \beta))]/\epsilon \\ &= \lim_{\epsilon \rightarrow 0+} [\delta\psi(\alpha; (\epsilon - \theta)(\beta - \gamma)) - \delta\psi(\alpha; (-\theta)(\beta - \gamma))]/\epsilon, \end{aligned}$$

and so with $\theta > 0$,

$$\begin{aligned} \delta q(\cdot; \alpha)(\gamma; \beta - \gamma) &= \lim_{\epsilon \rightarrow 0+} [(\theta - \epsilon)\delta\psi(\alpha; \gamma - \beta) - \theta\delta\psi(\alpha; \gamma - \beta)]/\epsilon \\ &= -\delta\psi(\alpha; \gamma - \beta), \end{aligned}$$

and with $\theta \leq 0$,

$$\begin{aligned} \delta q(\cdot; \alpha)(\gamma; \beta - \gamma) &= \lim_{\epsilon \rightarrow 0+} [(\epsilon - \theta)\delta\psi(\alpha; (\beta - \gamma)) - (-\theta)\delta\psi(\alpha; (\beta - \gamma))]/\epsilon \\ &= \delta\psi(\alpha; \beta - \gamma). \end{aligned}$$

The above derivation also guarantees the existence of $\Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma)$. Now, as $\theta \geq 1$, $\langle \beta - \alpha, \gamma - \beta \rangle \geq 0$ and so $q(\beta; \alpha) - q(\gamma; \alpha) - \delta q(\cdot; \alpha)(\gamma; \beta - \gamma) = \delta\psi(\alpha; \beta - \alpha) - \delta\psi(\alpha; \gamma - \alpha) + \delta\psi(\alpha; \gamma - \beta) = 0$. As $\theta \leq 0$, $\langle \beta - \alpha, \gamma - \beta \rangle \leq 0$ and $q(\beta; \alpha) - q(\gamma; \alpha) - \delta q(\cdot; \alpha)(\gamma; \beta - \gamma) = \delta\psi(\alpha; \beta - \alpha) - \delta\psi(\alpha; \gamma - \alpha) - \delta\psi(\alpha; \beta - \gamma) = 0$.

(iii) By definition, we have

$$\begin{aligned} &\delta q(\cdot; \alpha)(z; \beta - \gamma) \\ &= \lim_{\epsilon_2 \rightarrow 0+} \frac{1}{\epsilon_2} \{q(z + \epsilon_2(\beta - \gamma); \alpha) - q(z; \alpha)\} \\ &= \lim_{\epsilon_2 \rightarrow 0+} \frac{1}{\epsilon_2} \{\delta\psi(\alpha; z + \epsilon_2(\beta - \gamma) - \alpha) - \delta\psi(\alpha; z - \alpha)\}. \end{aligned}$$

Under the restricted linearity condition $\delta\psi(\alpha; h) = \langle g(\alpha), h \rangle, \forall h \in [\beta - \alpha, \gamma - \alpha]$, for $z = \gamma + t(\beta - \gamma)$ with $t \in [0, 1]$,

$$\begin{aligned} \delta q(\cdot; \alpha)(z; \beta - \gamma) &= \lim_{\epsilon_2 \rightarrow 0+} \frac{1}{\epsilon_2} \langle g(\alpha), z + \epsilon_2(\beta - \gamma) - \alpha - z + \alpha \rangle \\ &= \langle g(\alpha), \beta - \gamma \rangle. \end{aligned}$$

Under the restricted continuity condition $\lim_{\epsilon \rightarrow 0+} \delta\psi(\alpha + \epsilon h; \beta - \gamma) = \delta\psi(\alpha; \beta - \gamma)$, $\forall h \in [\beta - \alpha, \gamma - \alpha]$, for $z = \gamma + t(\beta - \gamma)$ with $t \in [0, 1)$,

$$\begin{aligned}
& \delta q(\cdot; \alpha)(z; \beta - \gamma) \\
&= \lim_{\epsilon_2 \rightarrow 0+} \frac{1}{\epsilon_2} \left\{ \lim_{\epsilon_1 \rightarrow 0+} \frac{1}{\epsilon_1} \left[\psi(\alpha + \epsilon_1[z + \epsilon_2(\beta - \gamma) - \alpha]) - \psi(\alpha) \right] \right. \\
&\quad \left. - \lim_{\epsilon_1 \rightarrow 0+} \frac{1}{\epsilon_1} \left[\psi(\alpha + \epsilon_1(z - \alpha)) - \psi(\alpha) \right] \right\} \\
&= \lim_{\epsilon_2 \rightarrow 0+} \lim_{\epsilon_1 \rightarrow 0+} \frac{1}{\epsilon_1 \epsilon_2} \left[\psi((1 - \epsilon_1)\alpha + \epsilon_1 z - \epsilon_1 \epsilon_2 \gamma + \epsilon_1 \epsilon_2 \beta) - \psi((1 - \epsilon_1)\alpha + \epsilon_1 z) \right] \\
&= \lim_{\epsilon_2 \rightarrow 0+} \lim_{\epsilon_1 \rightarrow 0+} \frac{1}{\epsilon_1 \epsilon_2} \int_0^1 \delta\psi((1 - \epsilon_1)\alpha + \epsilon_1 z + \epsilon_1 \epsilon_2 s(\beta - \gamma); \epsilon_1 \epsilon_2(\beta - \gamma)) \, ds \\
&= \lim_{\epsilon_2 \rightarrow 0+} \lim_{\epsilon_1 \rightarrow 0+} \int_0^1 \delta\psi((1 - \epsilon_1)\alpha + \epsilon_1 z + \epsilon_1 \epsilon_2 s(\beta - \gamma); \beta - \gamma) \, ds \\
&= \lim_{\epsilon_2 \rightarrow 0+} \int_0^1 \lim_{\epsilon_1 \rightarrow 0+} \delta\psi(\alpha + \epsilon_1(z + \epsilon_2 s(\beta - \gamma) - \alpha); \beta - \gamma) \, ds \\
&= \delta\psi(\alpha; \beta - \gamma),
\end{aligned}$$

where we used the positive homogeneity of $\delta\psi(\alpha; \cdot)$ and the dominated convergence theorem. (The integral is well-defined due to the boundedness and Lebesgue measurability of the integrand.)

The two sets of conditions are not equivalent in multiple dimensions. But in either case, $\delta q(\cdot; \alpha)(z; \beta - \gamma)$ is a term independent of z . Hence by Lemma 1 (iv),

$$\Delta_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) = \int_0^1 \left[\delta q(\cdot; \alpha)(\gamma + t(\beta - \gamma); \beta - \gamma) - \delta q(\cdot; \alpha)(\gamma; \beta - \gamma) \right] dt = 0.$$

A.3. Proof of Lemma 3. (i) Let $\varphi = b^*$. Then for all subgradient $g \in \partial\varphi(z)$, (g, z) makes a conjugate pair and so $\langle g, z \rangle = b(g) + \varphi(z)$ (see, e.g., [45]). Using the shorthand notation $(\partial\varphi(z), z)$, we represent it as $\langle \partial\varphi(z), z \rangle = b(\partial\varphi(z)) + \varphi(z)$. Therefore,

$$\begin{aligned}
\sigma^2 l_0(\eta; z) + b^*(z) &= -\langle z, \eta \rangle + b(\eta) + \varphi(z) \\
&= -\langle z, \eta \rangle + b(\eta) + \langle z, \partial\varphi(z) \rangle - b(\partial\varphi(z)) \\
&= b(\eta) - b(\partial\varphi(z)) - \langle z, \eta - \partial\varphi(z) \rangle \\
&= b(\eta) - b(\partial\varphi(z)) - \langle \nabla b(\partial\varphi(z)), \eta - \partial\varphi(z) \rangle \\
&= \Delta_b(\eta, \partial\varphi(z)).
\end{aligned}$$

When p_η is minimal, \mathcal{M} is full-dimensional and the canonical link $g = (\nabla b)^{-1}$ is well-defined on \mathcal{M}° (Proposition 3.1 and Proposition 3.2 in [62] can be slightly modified to include the dispersion parameter), and so $(g(z), \nabla b(g(z)))$ or $(g(z), z)$ makes a conjugate pair.

(ii) Let $\mu(\eta) = \nabla b(\eta)$ or μ for brevity. It follows that $\eta \in \partial\varphi(\mu)$ and so

$$\begin{aligned}
-\langle z, \eta \rangle + b(\eta) + b^*(z) &= -\langle z, \eta \rangle + \langle \mu, \eta \rangle - \varphi(\mu) + \varphi(z) \\
&= -\langle z - \mu, \eta \rangle - \varphi(\mu) + \varphi(z) \\
&\geq -\delta\varphi(\mu; z - \mu) - \varphi(\mu) + \varphi(z) = \Delta_\varphi(z, \mu),
\end{aligned}$$

where the inequality is due to [45, Theorem 23.2]. We claim that the inequality is actually an equality.

Indeed, if there exist $\eta_1, \eta_2 \in \partial\varphi(\mu)$ with $\eta_1 \neq \eta_2$, then $\bar{\Delta}_b(\eta_2, \eta_1) = \langle \nabla b(\eta_2) - \nabla b(\eta_1), \eta_2 - \eta_1 \rangle = 0$ and so $\Delta_b(\eta_2, \eta_1) = 0$ since b is convex. Therefore, for any random vector \mathbf{y} following p_η in the exponential family, where $\eta = t\eta_1 + (1-t)\eta_2$, $t \in (0, 1)$,

$$\text{Var}((\eta_2 - \eta_1)^T \mathbf{y}) = 0,$$

which can be obtained from Proposition 3.1 of [62]. Because $\exp((\langle \cdot, \eta \rangle - b(\eta))/\sigma^2) > 0$ for any $\eta \in \Omega$, we have $\langle \eta_2 - \eta_1, \mathbf{z} \rangle = c$ for almost every $\mathbf{z} \in \mathcal{Y}^n$ with respect to ν (i.e., p_η is not minimal). It follows that

$$\langle \mathbf{z} - \mu, \eta_1 - \eta_2 \rangle = 0.$$

Finally, from $\delta\varphi(\mu; \mathbf{h}) = \sup\{\langle \mathbf{g}, \mathbf{h} \rangle : \mathbf{g} \in \partial\varphi(\mu)\}$ [45, Theorem 23.4], the claim is true.

In the case that p_η is also minimal, φ can be shown to be strictly convex and differentiable on \mathcal{M}° [45, Theorem 26.4].

(iii) Let $dP_\eta = p_\eta d\nu_0$. By definition,

$$\text{KL}(p_{\eta_1}, p_{\eta_2}) = \int \log(dP_{\eta_1}/dP_{\eta_2}) dP_{\eta_1} = \int p_{\eta_1} \log(p_{\eta_1}/dp_{\eta_2}) d\nu_0$$

and so

$$\begin{aligned} \text{KL}(p_{\eta_1}, p_{\eta_2}) &= \int \log\{e^{(\langle \mathbf{y}, \eta_1 \rangle - b(\eta_1))/\sigma^2 - c(\mathbf{y}, \sigma^2)} / e^{(\langle \mathbf{y}, \eta_2 \rangle - b(\eta_2))/\sigma^2 - c(\mathbf{y}, \sigma^2)}\} dP_{\eta_1} \\ &= \frac{1}{\sigma^2} \int \langle \mathbf{y}, \eta_1 - \eta_2 \rangle - b(\eta_1) + b(\eta_2) dP_{\eta_1} \\ &= \frac{1}{\sigma^2} \{b(\eta_2) - b(\eta_1) + \int \langle \mathbf{y}, \eta_1 - \eta_2 \rangle dP_{\eta_1}\} \\ &= \frac{1}{\sigma^2} \{b(\eta_2) - b(\eta_1) + \langle \nabla b(\eta_1), \eta_1 - \eta_2 \rangle\} \\ &= \Delta_b(\eta_2, \eta_1)/\sigma^2, \end{aligned}$$

where the third equality is due to $\mathbb{E}_{\mathbf{y} \sim p_{\eta_1}} \mathbf{y} = \nabla b(\eta_1)$ under $\eta_1 \in \Omega^\circ$ (which can be derived from Proposition 3.1 of [62]). Moreover, from Lemma 1, $\sigma^2 \Delta_{l_0}(\eta_2, \eta_1) = \Delta_b(\eta_2, \eta_1)$.

A.4. Proof of Lemma 4. In this proof, all directional derivatives are with respect with β . The result of (i) is trivial from the construction of g . For (ii), by definition, we have $\delta g(\beta; \beta^-, \mathbf{h}) = \delta f(\beta; \mathbf{h}) + \delta \psi(\beta; \mathbf{h}) - \delta q(\beta; \beta^-, \mathbf{h})$ with $q(\beta; \beta^-) = \delta \psi(\beta^-; \beta - \beta^-)$. It follows from $q(\beta; \beta^-) = \lim_{\epsilon \rightarrow 0+} [\psi(\beta^- + \epsilon(\beta - \beta^-)) - \psi(\beta^-)]/\epsilon$ that

$$\begin{aligned} \delta q(\beta; \beta^-, \mathbf{h}) &= \lim_{\epsilon' \rightarrow 0+} [q(\beta + \epsilon' \mathbf{h}; \beta^-) - q(\beta; \beta^-)]/\epsilon' \\ &= \lim_{\epsilon' \rightarrow 0+} \left\{ (1/\epsilon') \lim_{\epsilon \rightarrow 0+} [\psi(\beta^- + \epsilon(\beta + \epsilon' \mathbf{h} - \beta^-)) - \psi(\beta^-)]/\epsilon \right\} \\ &\quad - \lim_{\epsilon' \rightarrow 0+} \delta \psi(\beta^-; \beta - \beta^-)/\epsilon'. \end{aligned}$$

When $\beta^- = \beta$, $\delta \psi(\beta^-; \beta - \beta^-) = 0$ and so

$$\begin{aligned} \delta q(\beta; \beta^-, \mathbf{h})|_{\beta^- = \beta} &= \lim_{\epsilon' \rightarrow 0+} \lim_{\epsilon \rightarrow 0+} [\psi(\beta + \epsilon(\epsilon' \mathbf{h})) - \psi(\beta)]/(\epsilon \epsilon') \\ &= \lim_{\epsilon'' \rightarrow 0+} [\psi(\beta + \epsilon'' \mathbf{h}) - \psi(\beta)]/\epsilon'' \\ &= \delta \psi(\beta; \mathbf{h}). \end{aligned}$$

The above argument also guarantees the existence of $\delta g(\beta; \beta^-, \mathbf{h})|_{\beta^- = \beta}$. Therefore, $\delta g(\beta; \beta^-, \mathbf{h})|_{\beta^- = \beta} = \delta f(\beta; \mathbf{h})$ for any β and \mathbf{h} .

A.5. Proof of Lemma 5. All results in Lemma 1 and Lemma 2 can be formulated for \mathbf{C} . For example, ψ is convex if and only if $\mathbf{C}_\psi \geq 0$, $\mathbf{C}_{a\phi+b\varphi} = a\mathbf{C}_\phi + b\mathbf{C}_\varphi$, $\Delta_\psi \geq \mu\mathbf{D}_2$ implies $\mathbf{C}_\psi \geq \mu\mathbf{C}_2$ since $\mathbf{C}_2(\alpha, \beta, \theta) := \mathbf{C}_{\|\cdot\|_2^2/2}(\alpha, \beta, \theta) = \theta(1-\theta)\mathbf{D}_2(\alpha, \beta)$, and so on. To show (i), we have

$$\begin{aligned} & \mathbf{C}_\psi(\alpha, \beta, \theta) + \Delta_\psi(\theta\alpha + (1-\theta)\beta, \alpha) \\ &= \theta\psi(\alpha) + (1-\theta)\psi(\beta) - \psi(\theta\alpha + (1-\theta)\beta) \\ & \quad + \psi(\theta\alpha + (1-\theta)\beta) - \psi(\alpha) - \delta\psi(\alpha; \theta\alpha + (1-\theta)\beta - \alpha) \\ &= (\theta-1)\psi(\alpha) + (1-\theta)\psi(\beta) - \delta\psi(\alpha; (1-\theta)(\beta-\alpha)) \\ &= (1-\theta)\psi(\beta) - (1-\theta)\psi(\alpha) - (1-\theta)\delta\psi(\alpha; \beta-\alpha) \\ &= (1-\theta)\Delta_\psi(\beta, \alpha). \end{aligned}$$

Similar to the proof of Lemma 2, let $q(\cdot; \alpha) = \delta\psi(\alpha; \cdot - \alpha)$. Then

$$\mathbf{C}_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma, \theta) - \mathbf{C}_\psi(\beta, \gamma, \theta) = \mathbf{C}_{\psi(\alpha) + q(\cdot; \alpha)}(\beta, \gamma, \theta) = \mathbf{C}_{q(\cdot; \alpha)}(\beta, \gamma, \theta),$$

without requiring the directional differentiability of $q(\cdot; \alpha)$. We can show analogous results to Lemma 2. For example, for any convex ψ , from the positively homogenous convexity of q ,

$$\mathbf{C}_{\Delta_\psi(\cdot, \alpha)} \leq \mathbf{C}_\psi$$

holds for any α , and for $\alpha = (1-\theta)\gamma + \theta\beta$ with $\theta \notin (0, 1)$,

$$\mathbf{C}_{\Delta_\psi(\cdot, \alpha)}(\beta, \gamma) = \mathbf{C}_\psi(\beta, \gamma)$$

follows from the restricted linearity of q . In particular, when $\nabla\psi(\alpha)$ exists, q is linear and so $\mathbf{C}_{q(\cdot; \alpha)} \equiv 0$ which gives the result in (ii).

A.6. Proof of Theorem 1. The theorem can be proved based on Theorem 6.1 of [37] and property (iii) of Lemma 3 in Section 2.1. We give some details for the second conclusion; the proof of the first follows similar lines and is easier. Consider a signal subclass

$$\mathcal{B}^1 = \{\beta : \beta_j \in \{0, \tau R\}, \|\beta\|_0 \leq s^*\},$$

where

$$R = [\sigma(\log(ep/s^*))^{1/2}/\bar{\kappa}^{1/2}] \wedge M$$

and $1 > \tau > 0$ is a small constant to be chosen later. Clearly, $\mathcal{B}^1 \in \mathcal{B}(s^*, M)$. By Stirling's approximation, $\log |\mathcal{B}^1| \geq \log \binom{p}{s^*} \geq s^* \log(p/s^*) \geq cs^* \log(ep/s^*)$ for some universal constant c .

Let $\rho(\beta_1, \beta_2) = \|\beta_1 - \beta_2\|_0$, the Hamming distance between β_1 and β_2 . By Lemma A.3 in [44], there exists a subset $\mathcal{B}^{10} \subset \mathcal{B}^1$ such that $\mathbf{0} \in \mathcal{B}^{10}$ and

$$\log |\mathcal{B}^{10}| \geq c_1 s^* \log(ep/s^*), \rho(\beta_1, \beta_2) \geq c_2 s^*, \forall \beta_1, \beta_2 \in \mathcal{B}^{10}, \beta_1 \neq \beta_2$$

for some universal constants $c_1, c_2 > 0$. Then

$$(A.6) \quad \|\mathbf{X}\beta_1 - \mathbf{X}\beta_2\|_2^2 \geq \underline{\kappa} \|\beta_1 - \beta_2\|_2^2 = \underline{\kappa} \tau^2 R^2 \rho(\beta_1, \beta_2) \geq c_2 \underline{\kappa} \tau^2 R^2 s^*$$

for any $\beta_1, \beta_2 \in \mathcal{B}^{10}, \beta_1 \neq \beta_2$.

By Lemma 3 (iii), since Ω is open, for any $\beta \in \mathcal{B}^{10}$, we have

$$\text{KL}(p_\beta, p_0) = \Delta_{l_0}(\mathbf{0}, \mathbf{X}\beta) \leq \tau^2 \bar{\kappa} R^2 s^* / (2\sigma^2).$$

Therefore,

$$(A.7) \quad \frac{1}{|\mathcal{B}^{10}| - 1} \sum_{\beta \in \mathcal{B}^{10} \setminus \{\mathbf{0}\}} \text{KL}(p_\beta, p_0) \leq \tau^2 s^* \log(ep/s^*).$$

Combining (A.6) and (A.7) and choosing a sufficiently small value for τ , we can apply Theorem 2.7 of [58] to get the desired lower bound.

A.7. Proof of Proposition 1. We first introduce a lemma.

LEMMA A.2. *For the sequence of iterates $\{\beta^{(t)}\}$ defined by (31) starting from an arbitrary point $\beta^{(0)}$, if $f(\cdot)$ and $g(\cdot; \beta^{(t)})$ are directionally differentiable, the following inequality holds for any β and $t \geq 0$*

$$(A.8) \quad \begin{aligned} & f(\beta) + \Delta_\psi(\beta, \beta^{(t)}) \\ & \geq f(\beta^{(t+1)}) + \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + (\Delta_{\Delta_\psi(\cdot; \beta^{(t)})} + \Delta_f)(\beta, \beta^{(t+1)}). \end{aligned}$$

It can be proved by Lemma A.1 and Lemma 1 (details omitted). Rearranging (A.8) gives

$$\begin{aligned} & f(\beta^{(t+1)}) - f(\beta) + \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + \Delta_f(\beta, \beta^{(t+1)}) \\ & \leq \Delta_\psi(\beta, \beta^{(t)}) - \Delta_{\Delta_\psi(\cdot; \beta^{(t)})}(\beta, \beta^{(t+1)}). \end{aligned}$$

Under $\Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + \Delta_f(\beta, \beta^{(t+1)}) \geq 0$, we have

$$(A.9) \quad f(\beta^{(t+1)}) - f(\beta) \leq \Delta_\psi(\beta, \beta^{(t)}) - \Delta_{\Delta_\psi(\cdot; \beta^{(t)})}(\beta, \beta^{(t+1)}).$$

By Lemma 2, when ψ is differentiable, $\Delta_{\Delta_\psi(\cdot; \beta^{(t)})}$ is well-defined and equals Δ_ψ . Adding up the corresponding inequality for $t = 0, 1, \dots, T$ leads to

$$\sum_{t=0}^T [f(\beta^{(t+1)}) - f(\beta)] \leq \Delta_\psi(\beta, \beta^{(0)}) - \Delta_\psi(\beta, \beta^{(T+1)}).$$

Therefore,

$$\text{avg}_{0 \leq t \leq T} f(\beta^{(t+1)}) - f(\beta) \leq \frac{1}{T+1} [\Delta_\psi(\beta, \beta^{(0)}) - \Delta_\psi(\beta, \beta^{(T+1)})].$$

Note that under just the directional differentiability of $\Delta_\psi(\cdot; \beta^{(t)})$, (35) can be replaced by $\Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + (\Delta_f + \Delta_{\Delta_\psi(\cdot; \beta^{(t)})} - \Delta_\psi)(\beta, \beta^{(t+1)}) \geq 0$, $0 \leq t \leq T$.

In the specific case that both f and ψ are convex, (35) is always satisfied by Lemma 1 and letting $\beta = \beta^{(t)}$ in (A.9) gives

$$f(\beta^{(t+1)}) - f(\beta^{(t)}) \leq -\Delta_\psi(\beta^{(t)}, \beta^{(t+1)}) \leq 0.$$

Hence $f(\beta^{(T+1)}) - f(\beta) = \min_{0 \leq t \leq T} f(\beta^{(t+1)}) - f(\beta) \leq \text{avg}_{0 \leq t \leq T} f(\beta^{(t+1)}) - f(\beta)$. The proof is complete.

A.8. Proof of Proposition 2. Substituting β^o for β in Lemma A.2 gives

$$(A.10) \quad \begin{aligned} & f(\beta^{(t+1)}) - f(\beta^o) + \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + \Delta_{\Delta_\psi(\cdot; \beta^{(t)})}(\beta^o, \beta^{(t+1)}) \\ & \leq \Delta_\psi(\beta^o, \beta^{(t)}) - \Delta_f(\beta^o, \beta^{(t+1)}). \end{aligned}$$

By Lemma A.1, we get

$$(A.11) \quad f(\beta^{(t+1)}) - f(\beta^o) \geq \Delta_f(\beta^{(t+1)}, \beta^o).$$

Combining (A.10) and (A.11) yields

$$(A.12) \quad (2\bar{\Delta}_f + \Delta_{\Delta_\psi(\cdot; \beta^{(t)})})(\beta^o, \beta^{(t+1)}) + \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) \leq \Delta_\psi(\beta^o, \beta^{(t)}).$$

It follows from the strong idempotence property that

$$(A.13) \quad (2\bar{\Delta}_f + \Delta_\psi)(\beta^o, \beta^{(t+1)}) \leq \Delta_\psi(\beta^o, \beta^{(t)}) - \min_{0 \leq t \leq T} \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}),$$

for any $0 \leq t \leq T$, and so (40) can be obtained under $2\bar{\Delta}_f \geq \varepsilon \Delta_\psi$.

To show the first result, since $\Delta_\phi = \Delta_\psi + \Delta_f$, (A.13) becomes

$$(2\bar{\Delta}_f + \Delta_\phi - \Delta_f)(\beta^o, \beta^{(t+1)}) \leq (\Delta_\phi - \Delta_f)(\beta^o, \beta^{(t)}) - \min_{0 \leq t \leq T} \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}).$$

Because $\kappa > 1$, (37) implies that

$$\Delta_f \geq (\kappa + 1)\bar{\Delta}_\phi/\kappa - \bar{\Delta}_\phi.$$

Applying the inequality twice, we obtain $((\kappa + 1)/\kappa)\bar{\Delta}_\phi(\beta^o, \beta^{(t+1)}) \leq (2 - (\kappa + 1)/\kappa)\bar{\Delta}_\phi(\beta^o, \beta^{(t)}) - \min_{0 \leq t \leq T} \Delta_\psi(\beta^{(t+1)}, \beta^{(t)})$, or

$$(A.14) \quad \bar{\Delta}_\phi(\beta^o, \beta^{(t+1)}) \leq \frac{\kappa - 1}{\kappa + 1} \bar{\Delta}_\phi(\beta^o, \beta^{(t)}) - \frac{\kappa}{\kappa + 1} \min_{0 \leq t \leq T} \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}).$$

The final conclusion can be obtained by applying (A.14) iteratively for $t = 0, 1, \dots, T$.

A.9. Proofs of Theorem 2 and Corollary 1. The proof of the theorem follows from Section A.7. In fact, setting $\beta = \beta^{(t)}$ in (A.8) gives

$$(\bar{\Delta}_\psi + \Delta_{\Delta_\psi(\cdot; \beta^{(t)})} + \Delta_f)(\beta^{(t)}, \beta^{(t+1)}) \leq f(\beta^{(t)}) - f(\beta^{(t+1)}),$$

which, by the weak idempotence property (with $\alpha = \beta^{(t)}$), reduces to

$$(A.15) \quad (2\bar{\Delta}_\psi + \Delta_f)(\beta^{(t)}, \beta^{(t+1)}) \leq f(\beta^{(t)}) - f(\beta^{(t+1)}).$$

Summing up (A.15) over $t = 0, 1, \dots, T$ gives the conclusion.

Next, we prove a result slightly more general than Corollary 1. Recall the surrogate

$$g(\beta; \beta^-) = f(\beta) + (\rho \mathbf{D}_\varphi - \Delta_f)(\beta, \beta^-)$$

where $\varphi \in \mathcal{C}^1$ is a strictly convex function, and f is continuous and directionally differentiable but not necessarily convex or differentiable. Denote $\arg \min g(\beta; \beta^-)$ by $\mathcal{T}(\beta^-)$.

COROLLARY 1'. Suppose that $\Delta_f \leq L\bar{\mathbf{D}}_\varphi$ for some $L > 0$ and the inverse stepsize parameter ρ satisfies $\rho > L/2$. Then $\text{avg}_{0 \leq t \leq T} (2\rho\bar{\mathbf{D}}_\varphi - \bar{\Delta}_f)(\beta^{(t)}, \beta^{(t+1)}) \leq \frac{f(\beta^{(0)})}{(T+1)}$ and so $\text{avg}_{0 \leq t \leq T} \bar{\mathbf{D}}_\varphi(\beta^{(t)}, \beta^{(t+1)}) \leq \frac{f(\beta^{(0)})}{(T+1)(2\rho-L)}$.

Moreover, for any accumulation point of $\beta^{(t)}$ at which \mathcal{T} is continuous, it must be a fixed point of \mathcal{T} . This is particularly true when $f \in \mathcal{C}^1$.

PROOF. Observe from (A.15) that

$$\begin{aligned} f(\beta^{(t)}) - f(\beta^{(t+1)}) &\geq (2\rho\bar{\mathbf{D}}_\varphi - 2\bar{\Delta}_f + \Delta_f)(\beta^{(t)}, \beta^{(t+1)}) \\ &\geq (2\rho\bar{\mathbf{D}}_\varphi - \Delta_f)(\beta^{(t+1)}, \beta^{(t)}) \\ &\geq (2\rho - L)\bar{\mathbf{D}}_\varphi(\beta^{(t+1)}, \beta^{(t)}) \geq 0. \end{aligned}$$

The error bounds can be obtained.

Let β^o be the limit point of some subsequence β^{t_l} as $l \rightarrow \infty$. Hence $f(\beta^{(t)})$ converges monotonically to $\lim_{l \rightarrow \infty} f(\beta^{t_l}) = f(\beta^o)$. It follows that

$$\lim_{t \rightarrow +\infty} \bar{\mathbf{D}}_\varphi(\beta^{(t+1)}, \beta^{(t)}) = 0.$$

\mathcal{T} is a well-defined function because of the strict convexity of the g -optimization problem. From the continuity assumptions,

$$0 = \lim_{l \rightarrow +\infty} \bar{\mathbf{D}}_\varphi(\beta^{(t_l+1)}, \beta^{(t_l)}) = \bar{\mathbf{D}}_\varphi(\mathcal{T}(\beta^o), \beta^o)$$

and thus $\mathcal{T}(\beta^o) = \beta^o$, i.e., β^o is a fixed point of \mathcal{T} . \square

A.10. Proof of Proposition 3. First, we show a result when using the Bregman surrogate $g(\beta; \beta^-) = l(\beta) + P(\varrho\beta) + \Delta_\psi(\beta, \beta^-)$ for solving $\min_\beta f(\beta) = l(\beta) + P(\varrho\beta)$ where l and P directionally differentiable and can be nonconvex. Define

$$(A.16) \quad \mathcal{L}_P := \inf\{\mathcal{L} \in \mathbb{R} : \Delta_P + \mathcal{L}\mathbf{D}_2 \geq 0\},$$

which provides an index to characterize the degree of nonconvexity of P , c.f. [36]. Assume $\mathcal{L}_P > -\infty$. Then for $\beta^{(t+1)} \in \arg \min_\beta g(\beta; \beta^{(t)})$, the following inequality holds for all $T \geq 1$

$$\text{avg}_{0 \leq t \leq T} (2\bar{\Delta}_\psi + \Delta_l - \varrho^2 \mathcal{L}_P \mathbf{D}_2)(\beta^{(t)}, \beta^{(t+1)}) \leq \frac{1}{T+1} [f(\beta^{(0)}) - f(\beta^{(T+1)})].$$

The result can be proved from Theorem 2, noticing the fact that $\Delta_f(\beta, \beta^-) = \Delta_l(\beta, \beta^-) + \Delta_P(\varrho\beta, \varrho\beta^-) \geq \Delta_l(\beta, \beta^-) - \mathcal{L}_P \mathbf{D}_2(\varrho\beta, \varrho\beta^-) = \Delta_l(\beta, \beta^-) - \varrho^2 \mathcal{L}_P \mathbf{D}_2(\beta, \beta^-)$ for any β, β^- . The details are omitted.

It suffices to proving the following lemma to complete the proof of Proposition 3.

LEMMA A.3. *Given any thresholding function Θ satisfying Definition 3, let P_Θ be the Θ -induced penalty in (48). Then \mathcal{L}_Θ as defined in (47) equals \mathcal{L}_{P_Θ} that is given in (A.16).*

PROOF. Since $\Delta_{P_\Theta}(\beta, \gamma) = \sum_j \Delta_{P_\Theta}(\beta_j, \gamma_j)$, it suffices to show the result in the univariate case. Recall that $\Theta^{-1}(u; \lambda) := \sup\{t : \Theta(t; \lambda) \leq u\}, \forall u > 0$. Since $P_\Theta(\gamma) = P_\Theta(|\gamma|) = \int_0^{|\gamma|} (\Theta^{-1}(u; \lambda) - u) du$, we assume $\gamma \geq 0$ without loss of generality. We define $s(u; \lambda) = \Theta^{-1}(u; \lambda) - u$ for $u \geq 0$, and extend $s(\cdot)$ to $(-\infty, 0)$ by $s(-u) = -s(u), u > 0$. Clearly, $s'(u) = s'(|u|)$ a.e., and so $-\mathcal{L}_\Theta = \text{ess inf}\{s'(u; \lambda) : u \neq 0\}$. By definition,

$$\delta P_\Theta(\gamma; \beta - \gamma) = \begin{cases} s(\gamma)(\beta - \gamma), & \text{if } \gamma \geq 0, \\ s(0)|\beta|, & \text{if } \gamma = 0. \end{cases}$$

When $\beta \geq 0$ and $\gamma \neq 0$, we get

$$\begin{aligned}
& (\Delta_{P_\Theta} + \mathcal{L}\mathbf{D}_2)(\beta, \gamma) \\
&= P_\Theta(\beta) - P_\Theta(\gamma) - \delta P_\Theta(\gamma; \beta - \gamma) + \mathcal{L}\mathbf{D}_2(\beta, \gamma) \\
&= \int_\gamma^\beta s(u) \, du - s(\gamma)(\beta - \gamma) + \frac{1}{2}\mathcal{L}(\beta - \gamma)^2 \\
&= \int_\gamma^\beta s(u) \, du - \int_\gamma^\beta s(\gamma) \, du + \mathcal{L} \int_\gamma^\beta (u - \gamma) \, du \\
&= \int_\gamma^\beta [s(u) - s(\gamma) + \mathcal{L}(u - \gamma)] \, du \\
&= \int_\gamma^\beta \int_\gamma^u [s'(v) + \mathcal{L}] \, dv \, du.
\end{aligned}$$

When $\beta < 0$ and $\gamma \neq 0$,

$$\begin{aligned}
& (\Delta_{P_\Theta} + \mathcal{L}\mathbf{D}_2)(\beta, \gamma) \\
&= P_\Theta(\beta) - P_\Theta(\gamma) - \delta P_\Theta(\gamma; \beta - \gamma) + \mathcal{L}\mathbf{D}_2(\beta, \gamma) \\
&= \int_\gamma^{-\beta} s(u) \, du - s(\gamma)(\beta - \gamma) + \frac{1}{2}\mathcal{L}(\beta - \gamma)^2 \\
&= \int_{-\gamma}^{-\beta} s(u) \, du - \int_{-\gamma}^{-\beta} s(-\gamma) \, du + \mathcal{L} \int_{-\gamma}^{-\beta} (u + \gamma) \, du \\
&= \int_{-\gamma}^{-\beta} [s(u) - s(-\gamma) + \mathcal{L}(u + \gamma)] \, du \\
&= \int_{-\gamma}^{-\beta} \int_{-\gamma}^u [s'(v) + \mathcal{L}] \, dv \, du.
\end{aligned}$$

Similarly, when $\gamma = 0$, $(\Delta_{P_\Theta} + \mathcal{L}\mathbf{D}_2)(\beta, 0) = \int_0^{|\beta|} \int_0^u [s(v) + \mathcal{L}] \, dv \, du$. It is then easy to verify that $\mathcal{L}_\Theta = \mathcal{L}_{P_\Theta}$. \square

A.11. Proof of Proposition 4. Let $f(\beta) = l(\beta) + P(\varrho\beta)$ and recall

$$\begin{aligned}
g_{\text{LLA}}^{(t)}(\beta; \beta^{(t)}) &= f(\beta) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta, \varrho\beta^{(t)}) \\
&= f(\beta) + \sum_j (\alpha_j^{(t)} \Delta_1 - \Delta_P)(\varrho\beta_j, \varrho\beta_j^{(t)}).
\end{aligned}$$

The proof is similar to that of Theorem 2 and we give some details for completeness. The important fact $\beta^{(t+1)} \in \arg \min_{\beta} g_{\text{LLA}}^{(t)}(\beta; \beta^{(t)})$ as shown in Example 5 implies

$$\Delta_{g_{\text{LLA}}^{(t)}(\cdot; \beta^{(t)})}(\beta^{(t)}, \beta^{(t+1)}) \leq f(\beta^{(t)}) - f(\beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t)}, \varrho\beta^{(t)}) - \Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t+1)}, \varrho\beta^{(t)})$$

or

$$\begin{aligned}
& f(\beta^{(t)}) - f(\beta^{(t+1)}) \\
& \geq \Delta_f(\beta^{(t)}, \beta^{(t+1)}) + \Delta_{\Delta_{\text{LLA}}^{(t)}(\varrho, \varrho\beta^{(t)})}(\beta^{(t)}, \beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t+1)}, \varrho\beta^{(t)}) \\
& = \Delta_f(\beta^{(t)}, \beta^{(t+1)}) + \Delta_{\Delta_{\text{LLA}}^{(t)}(\cdot, \varrho\beta^{(t)})}(\varrho\beta^{(t)}, \varrho\beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t+1)}, \varrho\beta^{(t)}) \\
& = \Delta_f(\beta^{(t)}, \beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t)}, \varrho\beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t+1)}, \varrho\beta^{(t)}) \\
& = \Delta_f(\beta^{(t)}, \beta^{(t+1)}) + 2\bar{\Delta}_{\text{LLA}}^{(t)}(\varrho\beta^{(t)}, \varrho\beta^{(t+1)}).
\end{aligned}$$

The conclusion follows from summing up this inequality for $t = 0, 1, \dots, T$.

A.12. Proofs of Theorem 3 and Theorem 4. Let $f(\beta) = l(\beta) + P_{\Theta}(\varrho\beta; \lambda)$ and recall $g(\beta; \beta^-) = f(\beta) + \Delta_{\psi}(\beta, \beta^-)$. We first introduce a lemma.

LEMMA A.4. *Let $\hat{\beta} \in \mathcal{F}$. Then for any $\beta \in \mathbb{R}^p$, we have the following inequality regardless of the specific form of ψ*

$$\begin{aligned}
& (\Delta_l - \varrho^2 \mathcal{L}_{\Theta} \mathbf{D}_2)(\beta, \hat{\beta}) + \Delta_l(\hat{\beta}, \beta^*) + P_{\Theta}(\varrho\hat{\beta}; \lambda) \\
& \leq \Delta_l(\beta, \beta^*) + \langle \epsilon, \mathbf{X}\hat{\beta} - \mathbf{X}\beta \rangle + P_{\Theta}(\varrho\beta; \lambda).
\end{aligned}
\tag{A.17}$$

PROOF. Denote $\hat{g}(\beta) := g(\beta; \hat{\beta}) = l(\beta) + P_{\Theta}(\varrho\beta; \lambda) + \Delta_{\psi}(\beta, \hat{\beta})$. Since $\hat{\beta}$ is a minimizer of $\hat{g}(\cdot)$, Lemma A.1 shows that for any β , $\Delta_{\hat{g}}(\beta, \hat{\beta}) \leq \hat{g}(\beta) - \hat{g}(\hat{\beta})$. On the one hand,

$$\begin{aligned}
& \hat{g}(\beta) - \hat{g}(\hat{\beta}) \\
& = l(\beta) - l(\hat{\beta}) + P_{\Theta}(\varrho\beta; \lambda) - P_{\Theta}(\varrho\hat{\beta}; \lambda) + \Delta_{\psi}(\beta, \hat{\beta}) \\
& = l(\beta) - l(\beta^*) - (l(\hat{\beta}) - l(\beta^*)) + P_{\Theta}(\varrho\beta; \lambda) - P_{\Theta}(\varrho\hat{\beta}; \lambda) + \Delta_{\psi}(\beta, \hat{\beta}) \\
& = \Delta_l(\beta, \beta^*) + \langle \nabla l(\beta^*), \beta - \beta^* \rangle - (\Delta_l(\hat{\beta}, \beta^*) + \langle \nabla l(\beta^*), \hat{\beta} - \beta^* \rangle) \\
& \quad + P_{\Theta}(\varrho\beta; \lambda) - P_{\Theta}(\varrho\hat{\beta}; \lambda) + \Delta_{\psi}(\beta, \hat{\beta}) \\
& = \Delta_l(\beta, \beta^*) - \Delta_l(\hat{\beta}, \beta^*) + \langle \epsilon, \mathbf{X}\hat{\beta} - \mathbf{X}\beta \rangle + P_{\Theta}(\varrho\beta; \lambda) - P_{\Theta}(\varrho\hat{\beta}; \lambda) + \Delta_{\psi}(\beta, \hat{\beta}).
\end{aligned}$$

On the other hand, by Lemma 1, Lemma 2, and Lemma A.3,

$$\begin{aligned}
\Delta_{\hat{g}}(\beta, \hat{\beta}) & = \Delta_l(\beta, \hat{\beta}) + \Delta_{P_{\Theta}(\varrho\cdot)}(\beta, \hat{\beta}) + \Delta_{\Delta_{\psi}(\cdot, \hat{\beta})}(\beta, \hat{\beta}) \\
& = \Delta_l(\beta, \hat{\beta}) + \Delta_{P_{\Theta}(\varrho\cdot)}(\beta, \hat{\beta}) + \Delta_{\psi}(\beta, \hat{\beta}) \\
& \geq \Delta_l(\beta, \hat{\beta}) - \varrho^2 \mathcal{L}_{\Theta} \mathbf{D}_2(\beta, \hat{\beta}) + \Delta_{\psi}(\beta, \hat{\beta}).
\end{aligned}$$

The conclusion follows. \square

To handle the stochastic term $\langle \epsilon, \mathbf{X}\hat{\beta} - \mathbf{X}\beta \rangle$ in (A.17), we introduce the following result.

LEMMA A.5. *Let $\mathbf{X} \in \mathbb{R}^{n \times p}$, ϵ be a sub-Gaussian random vector with mean 0 and scale bounded by σ , and $\lambda^o = \sigma \sqrt{\log(ep)}$. Suppose that $\varrho \geq \|\mathbf{X}\|_2$. Then there exist universal constants $A_0, C, c > 0$ such that for any $a \geq 2b > 0$ and $A_1 \geq A_0$, the following event*

$$\sup_{\beta \in \mathbb{R}^p} \left\{ 2\langle \epsilon, \mathbf{X}\beta \rangle - \frac{1}{a} \|\mathbf{X}\beta\|_2^2 - \frac{1}{b} [P_H(\varrho\beta; \sqrt{ab}A_1\lambda^o)] \right\} \geq a\sigma^2 t$$

occurs with probability at most $C \exp(-ct)p^{-cA_1^2}$.

The lemma can be proved by Lemma 4 of [49] based on a scaling argument.

Let $R = \sup_{\beta, \hat{\beta} \in \mathbb{R}^p} \{ \langle \epsilon, \mathbf{X}\hat{\beta} - \mathbf{X}\beta \rangle - \frac{1}{2a} \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta\|_2^2 - \frac{1}{2b} [P_H(\varrho(\hat{\beta} - \beta); \sqrt{ab}A\lambda^o)] \}$ with $\lambda^o = \sigma \sqrt{\log(ep)}$. Plugging the bound into (A.17) gives

$$(A.18) \quad \begin{aligned} & (\Delta_l - \varrho^2 \mathcal{L}_\Theta \mathbf{D}_2)(\beta, \hat{\beta}) + \Delta_l(\hat{\beta}, \beta^*) + P_\Theta(\varrho\hat{\beta}; \lambda) - P_\Theta(\varrho\beta; \lambda) \\ & \leq \Delta_l(\beta, \beta^*) + \frac{1}{2a} \|\mathbf{X}\hat{\beta} - \mathbf{X}\beta\|_2^2 + \frac{1}{2b} [P_H(\varrho(\hat{\beta} - \beta); \sqrt{ab}A\lambda^o)] + R \end{aligned}$$

with $\mathbb{P}(2R \geq a\sigma^2 t) \leq C \exp(-ct)p^{-cA^2}$ for any $a \geq 2b > 0$ and A large.

To prove Theorem 3, substitute β^* for β in (A.18) and combine it with the regularity condition (57), resulting in

$$(\delta - \frac{1}{a}) \mathbf{D}_2(\mathbf{X}\hat{\beta}, \mathbf{X}\beta^*) + (\vartheta - \frac{1}{2b}) P_H(\varrho(\hat{\beta} - \beta^*); \lambda) \leq K\lambda^2 J(\beta^*) + R$$

where $\lambda = \sqrt{ab}A\lambda^o$, $a \geq 2b > 0$, and $A \geq A_0$ with A_0 given in Lemma A.5. Setting $a = 2/(\delta \wedge (2\vartheta))$, $b = 1/(2\vartheta)$ or $a = 2/((2\delta) \wedge \vartheta)$, $b = 1/\vartheta$ bounds $\mathbf{D}_2(\mathbf{X}\hat{\beta}, \mathbf{X}\beta^*)$ or $P_H(\varrho(\hat{\beta} - \beta^*); \lambda)$. Finally, by Lemma A.5, $\mathbb{P}(R \leq 0) \geq 1 - Cp^{-cA^2}$.

For Theorem 4, we combine (A.18) and (60) with $\gamma = \hat{\beta}$:

$$\Delta_l(\hat{\beta}, \beta^*) + (\delta - \frac{1}{a}) \mathbf{D}_2(\mathbf{X}\beta, \mathbf{X}\hat{\beta}) \leq \frac{\alpha}{L} r \Delta_l(\beta, \hat{\beta}) + \Delta_l(\beta, \beta^*) + K\lambda^2 J(\beta) + R.$$

Take the same choice for λ and set $a = 2/(\delta \wedge (2\vartheta))$, $b = 1/(2\vartheta)$.

Case (i): $\alpha r/L = 0$. The conclusion follows easily, and does not need any restriction on α or L .

Case (ii): $\alpha r/L > 0$. Then $0 < r < 1$ and $\alpha/L > 0$. If $\alpha < 0$, $(\alpha r/L) \Delta_{l_0}(\mathbf{X}\beta, \mathbf{X}\hat{\beta}) \leq \alpha r d^2(\mathbf{X}\beta, \mathbf{X}\hat{\beta}) \leq 0$, reducing to the first case. Assume $\alpha > 0$. Then

$$\begin{aligned} \frac{\alpha}{L} r \Delta_l(\beta, \hat{\beta}) &= \frac{\alpha}{L} r \Delta_{l_0}(\mathbf{X}\beta, \mathbf{X}\hat{\beta}) \leq \alpha r d^2(\mathbf{X}\beta, \mathbf{X}\hat{\beta}) \\ &\leq \alpha r (d(\mathbf{X}\beta, \mathbf{X}\beta^*) + d(\mathbf{X}\hat{\beta}, \mathbf{X}\beta^*))^2 \\ &\leq \alpha r (1 + 1/M) d^2(\mathbf{X}\beta, \mathbf{X}\beta^*) + \alpha r (1 + M) d^2(\mathbf{X}\hat{\beta}, \mathbf{X}\beta^*) \\ &\leq r(1 + 1/M) \Delta_l(\beta, \beta^*) + r(1 + M) \Delta_l(\hat{\beta}, \beta^*), \end{aligned}$$

for any $M > 0$. Take $M = (1 - r)/(1 + r)$. Then

$$1/\{1 - r(1 + M)\} = 1 + r(1 + 1/M) = (1 + r)/(1 - r).$$

So we obtain

$$\begin{aligned} & (\frac{1-r}{1+r} \Delta_{l_0} + \frac{\delta}{2} \mathbf{D}_2)(\mathbf{X}\hat{\beta}, \mathbf{X}\beta^*) \\ & \leq \frac{1+r}{1-r} \Delta_{l_0}(\mathbf{X}\beta, \mathbf{X}\beta^*) + \frac{KA^2}{((2\vartheta) \wedge \delta)\vartheta} \sigma^2 J(\beta) \log(ep) + R. \end{aligned}$$

Finally, from $\mathbb{P}(2R \geq a\sigma^2 t) \leq C \exp(-ct)p^{-cA^2} \leq C \exp(-ct)$, we have $\mathbb{E}R \leq Ca\sigma^2$. The oracle inequality is proved. In fact, we also get

$$\begin{aligned} \mathbb{E}[\mathbf{D}_2(\mathbf{X}\hat{\beta}, \mathbf{X}\beta^*)] &\leq \frac{2(1+r)}{(1-r)\delta} \mathbb{E}[\Delta_{l_0}(\mathbf{X}\beta, \mathbf{X}\beta^*)] \\ &\quad + \frac{2KA^2}{(\vartheta \wedge \delta)\vartheta\delta} \sigma^2 J(\beta) \log(ep) + \frac{C}{(\vartheta \wedge \delta)\delta} \sigma^2 \end{aligned}$$

under the same condition.

REMARK A.1. Recall $J^* = J(\beta^*)$, $\mathcal{J}^* = \mathcal{J}(\beta^*)$ and $P(\beta; \lambda) = \sum_j P(\beta_j; \lambda)$. When $\hat{\beta}$ is a global minimizer, applying the bound of the stochastic term proved in Lemma A.5 gives the same conclusions (58), (59), under

$$(A.19) \quad \begin{aligned} & \vartheta P_H(\varrho(\beta - \beta^*); \lambda) + P_\Theta(\varrho\beta^*; \lambda) \\ & \leq (2\bar{\Delta}_{l_0} - \delta\mathbf{D}_2)(\mathbf{X}\beta, \mathbf{X}\beta^*) + P_\Theta(\varrho\beta; \lambda) + K\lambda^2 J^*, \forall \beta \end{aligned}$$

for some $\delta > 0$, $\vartheta > 0$ and large enough $K \geq 0$. Assuming P_Θ is subadditive, we can follow the arguments in Remark 2 to show that (A.19) is implied by $(1 + \vartheta)P_\Theta(\varrho(\beta - \beta^*)_{\mathcal{J}^*}; \lambda) \leq (2\bar{\Delta}_{l_0} - \delta\mathbf{D}_2)(\mathbf{X}\beta, \mathbf{X}\beta^*) + K\lambda^2 J^* + (1 - \vartheta)P_\Theta(\varrho(\beta - \beta^*)_{\mathcal{J}^{*c}}; \lambda)$. Furthermore, when l_0 is μ -strongly convex as in regression, one can take $\delta = \mu$ and the regularity condition is implied by

$$(A.20) \quad (1 + \bar{\vartheta})P_\Theta(\varrho\gamma_{\mathcal{J}^*}; \lambda) \leq \bar{K}\lambda\sqrt{J^*}\|\mathbf{X}\gamma\|_2 + P_\Theta(\varrho\gamma_{\mathcal{J}^{*c}}; \lambda), \forall \gamma$$

for some $\bar{\vartheta} (= \frac{2\vartheta}{1-\vartheta}) > 0$ and $\bar{K} (= \frac{\sqrt{2(2\mu-\delta)K}}{1-\vartheta}) \geq 0$, or the constrained forms

$$(A.21) \quad [P_\Theta(\varrho\gamma_{\mathcal{J}^*}; \lambda)]^2 \leq \tilde{K}\lambda^2 J^* \|\mathbf{X}\gamma\|_2^2, \quad \forall \gamma : P_\Theta(\varrho\gamma_{\mathcal{J}^{*c}}; \lambda) \leq (1 + \bar{\vartheta})P_\Theta(\varrho\gamma_{\mathcal{J}^*}; \lambda)$$

$$(A.22) \quad \sum_{j \in \mathcal{J}^*} P_\Theta^2(\varrho\gamma_j; \lambda) \leq \tilde{K}\lambda^2 \|\mathbf{X}\gamma\|_2^2, \quad \forall \gamma : P_\Theta(\varrho\gamma_{\mathcal{J}^{*c}}; \lambda) \leq (1 + \bar{\vartheta})P_\Theta(\varrho\gamma_{\mathcal{J}^*}; \lambda)$$

for some $\bar{\vartheta} > 0$ and $\tilde{K} \geq 0$. The conclusions and conditions can also be formulated in the oracle inequality setup of Theorem 4. (A.20), (A.21) and (A.22) extend the comparison condition (62), compatibility condition and RE condition to a more general penalty.

A.13. Proof of Theorem 5. We prove the result under a more relaxed assumption: l, g are merely directionally differentiable, and (64) is replaced by

$$\begin{aligned} & (\varrho^2 \mathcal{L}_\Theta \mathbf{D}_2 + \varepsilon \Delta_\psi)(\beta^*, \beta) + (\Delta_\psi - \Delta_{\Delta_\psi(\cdot; \alpha)})(\beta^*, \beta) + \vartheta P_H(\varrho(\beta - \beta^*); \lambda) + P_\Theta(\varrho\beta^*; \lambda) \\ & \leq (2\bar{\Delta}_{l_0} - \delta\mathbf{D}_2)(\mathbf{X}\beta, \mathbf{X}\beta^*) + P_\Theta(\varrho\beta; \lambda) + K\lambda^2 J(\beta^*), \forall \beta, \alpha \end{aligned}$$

for some $\delta > 0$, $\varepsilon > 0$, $\vartheta > 0$ and large $K \geq 0$.

Recall the objective function $f(\beta) = l(\beta) + P_\Theta(\varrho\beta; \lambda)$ and the surrogate function $g(\beta; \beta^-) = f(\beta) + \Delta_\psi(\beta, \beta^-)$. From Lemma A.2 and Lemma A.3, we obtain

$$\begin{aligned} & f(\beta) + \Delta_\psi(\beta, \beta^{(t)}) \\ & \geq f(\beta^{(t+1)}) + \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + (\Delta_{\Delta_\psi(\cdot; \beta^{(t)})} + \Delta_l + \Delta_{P_\Theta(\varrho\cdot)})(\beta, \beta^{(t+1)}) \\ & \geq f(\beta^{(t+1)}) + \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + (\Delta_{\Delta_\psi(\cdot; \beta^{(t)})} + \Delta_l - \varrho^2 \mathcal{L}_\Theta \mathbf{D}_2)(\beta, \beta^{(t+1)}). \end{aligned}$$

Substituting β^* for β yields

$$\begin{aligned} & \Delta_\psi(\beta^*, \beta^{(t)}) \\ & \geq \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + (\Delta_{\Delta_\psi(\cdot; \beta^{(t)})} + \Delta_l)(\beta^*, \beta^{(t+1)}) - \mathcal{L}_\Theta \mathbf{D}_2(\varrho\beta^*, \varrho\beta^{(t+1)}) \\ & \quad + l(\beta^{(t+1)}) - l(\beta^*) + P_\Theta(\varrho\beta^{(t+1)}; \lambda) - P_\Theta(\varrho\beta^*; \lambda) \\ & = \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) + \Delta_{\Delta_\psi(\cdot; \beta^{(t)})}(\beta^*, \beta^{(t+1)}) - \mathcal{L}_\Theta \mathbf{D}_2(\varrho\beta^*, \varrho\beta^{(t+1)}) \\ & \quad 2\bar{\Delta}_l(\beta^*, \beta^{(t+1)}) - \langle \epsilon, \mathbf{X}\beta^{(t+1)} - \mathbf{X}\beta^* \rangle + P_\Theta(\varrho\beta^{(t+1)}; \lambda) - P_\Theta(\varrho\beta^*; \lambda). \end{aligned}$$

From the above regularity condition,

$$\begin{aligned} & \varepsilon \Delta_\psi(\beta^*, \beta^{(t+1)}) + \mathcal{L}_\Theta \mathbf{D}_2(\varrho \beta^*, \varrho \beta^{(t+1)}) + \delta \mathbf{D}_2(\mathbf{X} \beta^*, \mathbf{X} \beta^{(t+1)}) \\ & + \vartheta P_H(\varrho(\beta^{(t+1)} - \beta^*); \lambda) + P_\Theta(\varrho \beta^*; \lambda) + (\Delta_\psi - \Delta_{\Delta_\psi(\cdot; \beta^{(t)})})(\beta^*, \beta^{(t+1)}) \\ & \leq 2\bar{\Delta}_l(\beta^*, \beta^{(t+1)}) + P_\Theta(\varrho \beta^{(t+1)}; \lambda) + K\lambda^2 J(\beta^*), \end{aligned}$$

and so we obtain

$$\begin{aligned} & (\varepsilon + 1) \Delta_\psi(\beta^*, \beta^{(t+1)}) + \Delta_\psi(\beta^{(t+1)}, \beta^{(t)}) \\ (A.23) \quad & + \delta \mathbf{D}_2(\mathbf{X} \beta^*, \mathbf{X} \beta^{(t+1)}) + \vartheta P_H(\varrho(\beta^{(t+1)} - \beta^*); \lambda) \\ & \leq \Delta_\psi(\beta^*, \beta^{(t)}) + \langle \epsilon, \mathbf{X} \beta^{(t+1)} - \mathbf{X} \beta^* \rangle + K\lambda^2 J^*. \end{aligned}$$

According to Lemma A.5, as long as $A \geq A_0/\sqrt{2}$, choosing $\lambda = \sqrt{2ab}A\lambda^o$, $b = 1/(2\vartheta)$, $a = 1/(\delta \wedge \vartheta)$ guarantees that the probability of the following inequality occurring for all t ,

$$\langle \epsilon, \mathbf{X} \beta^{(t+1)} - \mathbf{X} \beta^* \rangle - \delta \mathbf{D}_2(\mathbf{X} \beta^*, \mathbf{X} \beta^{(t+1)}) - \vartheta P_H(\varrho(\beta^{(t+1)} - \beta^*); \lambda) \geq 0,$$

is no greater than Cp^{-cA^2} . Together with (A.23), with probability $1 - Cp^{-cA^2}$

$$(A.24) \quad \Delta_\psi(\beta^*, \beta^{(t+1)}) \leq \frac{1}{\varepsilon + 1} [\Delta_\psi(\beta^*, \beta^{(t)}) - \Delta(\beta^{(t+1)}, \beta^{(t)}) + K\lambda^2 J^*]$$

for all t . The desired inequality can be shown by iteratively applying (A.24) for $t = 0, 1, 2, \dots$

A.14. Proofs of Theorem 6 and Corollary 2. First we prove Theorem 6. Note that (67b), (69b) have additional terms involving μ_0 . The first result for $\mu_0 = 0$ can be shown based on a GBF translation of the proof of Proposition 1 in [57]. For convenience, let $h_t(\beta) = f(\beta) - \Delta_{\bar{\psi}_0}(\beta, \gamma^{(t)}) = f(\beta) - \Delta_{\psi_0}(\beta, \gamma^{(t)}) + \mu_0 \Delta_\phi(\beta, \gamma^{(t)})$. Applying Lemma A.2 to (67b) yields $(\Delta_f - \Delta_{\Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})} + \theta_t \rho_t \Delta_{\Delta_\phi(\cdot, \alpha^{(t)})})(\beta, \alpha^{(t+1)}) \leq h_t(\beta) + \theta_t \rho_t \Delta_\phi(\beta, \alpha^{(t)}) - h_t(\alpha^{(t+1)}) - \theta_t \rho_t \Delta_\phi(\alpha^{(t+1)}, \alpha^{(t)}) \forall \beta$, or

$$\begin{aligned} & h_t(\alpha^{(t+1)}) - h_t(\beta) + \theta_t \rho_t \Delta_\phi(\alpha^{(t+1)}, \alpha^{(t)}) \\ (A.25) \quad & \leq \theta_t \rho_t \Delta_\phi(\beta, \alpha^{(t)}) - (\theta_t \rho_t \Delta_{\Delta_\phi(\cdot, \alpha^{(t)})} + \Delta_{f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})})(\beta, \alpha^{(t+1)}). \end{aligned}$$

By definition, $\mathbf{C}_{h_t}(\alpha^{(t+1)}, \beta^{(t)}, \theta_t) = \theta_t h_t(\alpha^{(t+1)}) + (1 - \theta_t) h_t(\beta^{(t)}) - h_t(\beta^{(t+1)})$; adding it to (A.25) multiplied by θ_t gives

$$\begin{aligned} & h_t(\beta^{(t+1)}) - (1 - \theta_t) h_t(\beta^{(t)}) - \theta_t h_t(\beta) \\ & + \theta_t^2 \rho_t \Delta_\phi(\alpha^{(t+1)}, \alpha^{(t)}) + \mathbf{C}_{h_t}(\alpha^{(t+1)}, \beta^{(t)}, \theta_t) \\ & + \theta_t^2 \rho_t \Delta_{\Delta_\phi(\cdot, \alpha^{(t)}) - \phi(\cdot)}(\beta, \alpha^{(t+1)}) \\ & \leq \theta_t^2 \rho_t \Delta_\phi(\beta, \alpha^{(t)}) - (\theta_t^2 \rho_t \Delta_\phi + \theta_t \Delta_{f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})})(\beta, \alpha^{(t+1)}), \end{aligned}$$

and so

$$\begin{aligned} & f(\beta^{(t+1)}) - f(\beta) - (1 - \theta_t)[f(\beta^{(t)}) - f(\beta)] + \theta_t \Delta_{\bar{\psi}_0}(\beta, \gamma^{(t)}) \\ (A.26) \quad & + \theta_t \{(\Delta_{f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})} + \theta_t \rho_t \Delta_{\Delta_\phi(\cdot, \alpha^{(t)}) - \phi(\cdot)})(\beta, \alpha^{(t+1)})\} + R_t \\ & \leq \theta_t^2 \rho_t (\Delta_\phi(\beta, \alpha^{(t)}) - \Delta_\phi(\beta, \alpha^{(t+1)})), \forall t \geq 0 \end{aligned}$$

where R_t is given by $\theta_t^2 \rho_t \Delta_\phi(\alpha^{(t+1)}, \alpha^{(t)}) - \Delta_{\bar{\psi}_0}(\beta^{(t+1)}, \gamma^{(t)}) + (1 - \theta_t) \Delta_{\bar{\psi}_0}(\beta^{(t)}, \gamma^{(t)}) + \mathbf{C}_{f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})}(\alpha^{(t+1)}, \beta^{(t)}, \theta_t)$.

Under $\mu_0 = 0$, (A.26) implies that

$$(A.27) \quad \begin{aligned} & \frac{1}{\theta_t^2 \rho_t} [f(\beta^{(t+1)}) - f(\beta)] - \frac{1 - \theta_t}{\theta_t^2 \rho_t} [f(\beta^{(t)}) - f(\beta)] + \frac{\mathcal{E}_t(\beta)}{\theta_t \rho_t} + \frac{R_t}{\theta_t^2 \rho_t} \\ & \leq \Delta_\phi(\beta, \alpha^{(t)}) - \Delta_\phi(\beta, \alpha^{(t+1)}). \end{aligned}$$

Since in this case (69b) gives $(1 - \theta_t)/\theta_t^2 \rho_t = 1/\theta_{t-1}^2 \rho_{t-1}$ for any $t \geq 1$, we obtain the first conclusion

$$\begin{aligned} & \frac{1}{\theta_T^2 \rho_T} [f(\beta^{(T+1)}) - f(\beta)] - \frac{1 - \theta_0}{\theta_0^2 \rho_0} [f(\beta^{(0)}) - f(\beta)] + \sum_{t=0}^T \left(\frac{\mathcal{E}_t(\beta)}{\theta_t \rho_t} + \frac{R_t}{\theta_t^2 \rho_t} \right) \\ & \leq \Delta_\phi(\beta, \alpha^{(0)}) - \Delta_\phi(\beta, \alpha^{(T+1)}). \end{aligned}$$

On the other hand, given $\mu_0 \geq 0$, (A.26) can be written as

$$(A.28) \quad \begin{aligned} & f(\beta^{(t+1)}) - f(\beta) - (1 - \theta_t)[f(\beta^{(t)}) - f(\beta)] \\ & + R_t + \theta_t \Delta_{\bar{\psi}_0}(\beta, \gamma^{(t)}) + \theta_t \Delta_{f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})}(\beta, \alpha^{(t+1)}) \\ & + \theta_t (\mu_0 \Delta_{\Delta_\phi(\cdot, \gamma^{(t)}) - \phi(\cdot)} + \theta_t \rho_t \Delta_{\Delta_\phi(\cdot, \alpha^{(t)}) - \phi(\cdot)})(\beta, \alpha^{(t+1)}) \\ & \leq \theta_t^2 \rho_t \Delta_\phi(\beta, \alpha^{(t)}) - \theta_t^2 (\rho_t + \frac{\mu_0}{\theta_t}) \Delta_\phi(\beta, \alpha^{(t+1)}). \end{aligned}$$

Therefore, we have

$$\begin{aligned} & f(\beta^{(t+1)}) - f(\beta) + \theta_t^2 (\rho_t + \frac{\mu_0}{\theta_t}) \Delta_\phi(\beta, \alpha^{(t+1)}) + \theta_t \mathcal{E}_t(\beta) + R_t \\ & \leq (1 - \theta_t)[f(\beta^{(t)}) - f(\beta)] + \theta_t^2 \rho_t \Delta_\phi(\beta, \alpha^{(t)}), \forall t \geq 0 \end{aligned}$$

and from (69b),

$$\begin{aligned} & f(\beta^{(t+1)}) - f(\beta) + \theta_t^2 (\rho_t + \frac{\mu_0}{\theta_t}) \Delta_\phi(\beta, \alpha^{(t+1)}) + \theta_t \mathcal{E}_t(\beta) + R_t \\ & \leq (1 - \theta_t)[f(\beta^{(t)}) - f(\beta) + \theta_{t-1}^2 (\rho_{t-1} + \frac{\mu_0}{\theta_{t-1}}) \Delta_\phi(\beta, \alpha^{(t)})], \forall t \geq 1 \end{aligned}$$

The second conclusion can be obtained by a recursive argument and $R_T + \theta_T \mathcal{E}_T(\beta) + (1 - \theta_T)(R_{T-1} + \theta_{T-1} \mathcal{E}_{T-1}(\beta)) + \dots + (1 - \theta_T) \dots (1 - \theta_1)(R_0 + \theta_0 \mathcal{E}_0(\beta)) = \sum_{t=0}^T (\prod_{s=t+1}^T (1 - \theta_s))(R_t + \theta_t \mathcal{E}_t(\beta))$.

REMARK A.2. With the ‘=’ in (69b) replaced by ‘ \leq ’, (71) still holds when $\Delta_\phi \geq 0$ (or ϕ is convex), and (70) still holds if we set β to be a minimizer of f . But the equality form of (69b) makes our conclusions applicable to say the noise-free statistical truth $\beta = \beta^*$, which may not be a minimizer of the sample-based objective. The same comment applies to Theorem 7.

Also, it is trivial to see that the conclusions extend to a varying sequence of μ_t . (Concretely, the μ_0 in (67b), (68), and $\mathcal{E}_t(\beta)$ becomes μ_t , and the μ_0 in (69b), (71) becomes μ_{t-1}, μ_T , respectively.) One can add backtracking for μ_t in the algorithm to further reduce its iteration complexity.

Finally, we prove Corollary 2' which implies Corollary 2 and applies to any convex l_0 in (A.29) below. The proof is based on an accumulative R_t bound that can be derived in a more general setup; see (A.32) in Remark A.3.

Here, the optimization problem of interest in “variable screening” is

$$(A.29) \quad \min l(\beta) = l_0(\mathbf{X}\beta) \text{ s.t. } \|\beta\|_0 \leq q,$$

to estimate the target β^* satisfying the strict inequality $\|\beta^*\|_0 < q$. Take $\mu_0 = 0$, $\phi = \|\cdot\|_2^2/2$, $\psi_0(\beta) = l(\beta) - \mathcal{L}\phi(\beta)$ for some $\mathcal{L} \geq 0$.

Given l_0 , \mathbf{X} , and $s \leq p$, we extend the notion of restricted isometry numbers ρ_+, ρ_- [12]:

$$(A.30) \quad \rho_-(s)\mathbf{D}_2(\beta, \gamma) \leq \Delta_{l_0}(\mathbf{X}\beta, \mathbf{X}\gamma) \leq \rho_+(s)\mathbf{D}_2(\beta, \gamma), \forall \beta, \gamma : \|\beta - \gamma\|_0 \leq s.$$

(The dependence on \mathbf{X} and l_0 is dropped for the sake of brevity.)

COROLLARY 2'. *Let l_0 be any convex function and q be any nonnegative integer no more than p . As long as $r = q/\|\beta^*\|_0 > 1$, for any $\mathcal{L} \geq \rho_+(2q)/\sqrt{r}$, there must exist a universal $\rho_t = \rho_0, \forall t$, e.g.,*

$$\rho_0 = (1 - 1/\sqrt{r})\rho_+(2q),$$

and thus $\theta_{t+1} = (\sqrt{\theta_t^4 + 4\theta_t^2} - \theta_t^2)/2$, so that the accelerated iterative quantile-thresholding according to (67a)–(67c) satisfies

$$\begin{aligned} \frac{l(\beta^{(T+1)}) - l(\beta^*)}{\theta_T^2} + T \cdot \text{avg}_{0 \leq t \leq T} \frac{\Delta_{\psi_0}(\beta^*, \gamma^{(t)})}{\theta_t} + \frac{\mathcal{L}(1 - \theta_T)}{\theta_T^3} \mathbf{D}_2(\beta^{(T+1)}, \beta^{(T)}) \\ \leq \rho_0 \mathbf{D}_2(\beta^*, \alpha^{(0)}), \text{ for all } T \geq 0, \end{aligned}$$

or $l(\beta^{(T+1)}) - l(\beta^*) + \min_{0 \leq t \leq T} \Delta_{\psi_0}(\beta^*, \gamma^{(t)}) \leq A/T^2$ with A independent of T .

Seen from the error measure, r should be appropriately large (but cannot be too large from the perspective of statistical accuracy.)

To prove the corollary, we first introduce a useful result [50, Lemma 9].

LEMMA A.6. *Given $q \leq p$, $\hat{\beta} = \Theta^\#(\mathbf{y}; q)$ is a globally optimal solution to $\min_{\beta \in \mathbb{R}^p} l(\beta) = \|\mathbf{y} - \beta\|_2^2/2$ s.t. $\|\beta\|_0 \leq q$. Let $\mathcal{J} = \mathcal{J}(\beta)$, $\hat{\mathcal{J}} = \mathcal{J}(\hat{\beta})$ and assume $J(\hat{\beta}) = q$. Then, for any β with $J(\beta) \leq s = q/r$ and $r \geq 1$,*

$$l(\beta) - l(\hat{\beta}) \geq \{1 - \mathcal{L}(\mathcal{J}, \hat{\mathcal{J}})\} \mathbf{D}_2(\hat{\beta}, \beta),$$

where $\mathcal{L}(\mathcal{J}, \hat{\mathcal{J}}) = (|\mathcal{J} \setminus \hat{\mathcal{J}}|/|\hat{\mathcal{J}} \setminus \mathcal{J}|)^{1/2} \leq (s/q)^{1/2} = r^{-1/2}$.

Set $\beta = \beta^*$ in the previous proof, and apply, instead of Lemma A.2, Lemma A.6 (where $\|\alpha^{(t+1)}\|_0 = q$ due to the no-tie-occurring assumption) to (67b). (A.25) is then replaced by

$$\begin{aligned} l(\alpha^{(t+1)}) - \Delta_{\psi_0}(\alpha^{(t+1)}, \gamma^{(t)}) - l(\beta^*) + \Delta_{\psi_0}(\beta^*, \gamma^{(t)}) + \theta_t \rho_t \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) \\ \leq \theta_t \rho_t \mathbf{D}_2(\beta^*, \alpha^{(t)}) - (\theta_t \rho_t + \mathcal{L})(1 - \frac{1}{\sqrt{r}}) \mathbf{D}_2(\beta^*, \alpha^{(t+1)}). \end{aligned}$$

Accordingly, (A.26) becomes

$$\begin{aligned} l(\beta^{(t+1)}) - l(\beta^*) - (1 - \theta_t)(l(\beta^{(t)}) - l(\beta^*)) + R_t \\ + \theta_t \left\{ \Delta_{\psi_0}(\beta^*, \gamma^{(t)}) + \left[\mathcal{L}(1 - \frac{1}{\sqrt{r}}) - \frac{\theta_t \rho_t}{\sqrt{r}} \right] \mathbf{D}_2(\beta^*, \alpha^{(t+1)}) \right\} \\ \leq \theta_t^2 \rho_t (\mathbf{D}_2(\beta^*, \alpha^{(t)}) - \mathbf{D}_2(\beta^*, \alpha^{(t+1)})), \end{aligned}$$

where R_t is the same as before. Based on Lemma 1, Lemma 2, Lemma 5 (together with some results in its proof), (A.30), and the following facts

$$\begin{aligned}\beta^{(t)} - \gamma^{(t)} &= \theta_t(\beta^{(t)} - \alpha^{(t)}) = \theta_t(1 - \theta_{t-1})(\beta^{(t-1)} - \alpha^{(t)}), \forall t \geq 1 \\ \beta^{(t+1)} - \gamma^{(t)} &= \theta_t(\alpha^{(t+1)} - \alpha^{(t)}), \forall t \geq 0\end{aligned}$$

we obtain for all $t \geq 1$,

$$\begin{aligned}R_t &\geq \theta_t^2 \rho_t \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) - \Delta_{\psi_0}(\beta^{(t+1)}, \gamma^{(t)}) + (1 - \theta_t) \Delta_{\psi_0}(\beta^{(t)}, \gamma^{(t)}) \\ &\quad + \mathcal{L} \theta_t (1 - \theta_t) \mathbf{D}_2(\alpha^{(t+1)}, \beta^{(t)}) \\ &= \theta_t^2 (\rho_t + \mathcal{L}) \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) - \Delta_{l_0}(X\beta^{(t+1)}, X\gamma^{(t)}) + (1 - \theta_t) \Delta_l(\beta^{(t)}, \gamma^{(t)}) \\ &\quad + \mathcal{L} \{ \theta_t (1 - \theta_t) \mathbf{D}_2(\alpha^{(t+1)}, \beta^{(t)}) - (1 - \theta_t) \mathbf{D}_2(\beta^{(t)}, \gamma^{(t)}) \} \\ &= \theta_t^2 (\rho_t + \mathcal{L}) \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) - \Delta_{l_0}(X\beta^{(t+1)}, X\gamma^{(t)}) + (1 - \theta_t) \Delta_l(\beta^{(t)}, \gamma^{(t)}) \\ &\quad + \mathcal{L} \{ \theta_t (1 - \theta_t) \mathbf{D}_2(\alpha^{(t+1)}, \beta^{(t)}) - (1 - \theta_t) \theta_t^2 (1 - \theta_{t-1})^2 \mathbf{D}_2(\alpha^{(t)}, \beta^{(t-1)}) \} \\ &\geq \theta_t^2 (\rho_t + \mathcal{L}) \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) - \theta_t^2 \rho_+ (2q) \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) \\ &\quad + \mathcal{L} \theta_t (1 - \theta_t) \{ \mathbf{D}_2(\alpha^{(t+1)}, \beta^{(t)}) - \theta_t (1 - \theta_{t-1})^2 \mathbf{D}_2(\alpha^{(t)}, \beta^{(t-1)}) \},\end{aligned}$$

and $R_t \geq \theta_t^2 \{ \rho_t + \mathcal{L} - \rho_+ (2q) \} \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)})$ as $t = 0$. It follows that

$$\begin{aligned}\sum_{t=0}^T \frac{R_t}{\theta_t^2 \rho_t} &\geq \sum_{t=0}^T \frac{\rho_t + \mathcal{L} - \rho_+ (2q)}{\rho_t} \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) + \mathcal{L} \frac{1 - \theta_T}{\theta_T \rho_T} \mathbf{D}_2(\alpha^{(T+1)}, \beta^{(T)}) \\ &\quad + \mathcal{L} \sum_{t=0}^{T-1} \left\{ \frac{1 - \theta_t}{\theta_t \rho_t} - \frac{(1 - \theta_t)^2 (1 - \theta_{t+1})}{\rho_{t+1}} \right\} \mathbf{D}_2(\alpha^{(t+1)}, \beta^{(t)}) \\ &= \sum_{t=0}^T \frac{\rho_t + \mathcal{L} - \rho_+ (2q)}{\rho_t} \mathbf{D}_2(\alpha^{(t+1)}, \alpha^{(t)}) + \frac{\mathcal{L} (1 - \theta_T)}{\theta_T^3 \rho_T} \mathbf{D}_2(\beta^{(T+1)}, \beta^{(T)}) \\ &\quad + \mathcal{L} \sum_{t=0}^{T-1} \frac{1 - \theta_t}{\theta_t \rho_t} \left\{ 1 - (1 - \theta_t) \theta_{t+1} \frac{\theta_{t+1}}{\theta_t} \right\} \mathbf{D}_2(\alpha^{(t+1)}, \beta^{(t)}).\end{aligned}$$

Therefore, choosing a universal $\rho_t = \rho_0 \geq \rho_+ (2q) - \mathcal{L}$ (which implies $\theta_t \downarrow$) ensures $\sum_{t=0}^T R_t / (\theta_t^2 \rho_t) \geq 0$.

Moreover, $\mathcal{E}_t(\beta^*) \geq \Delta_{\psi_0}(\beta^*, \gamma^{(t)})$ holds under $\{ \mathcal{L} (1 - \frac{1}{\sqrt{r}}) - \frac{\theta_t \rho_t}{\sqrt{r}} \} \mathbf{D}_2(\beta^*, \alpha^{(t+1)}) \geq 0$ or $L(1 - 1/\sqrt{r}) \geq \rho_t / \sqrt{r}$. It is easy to see that as long as $r > 1$, there exist positive \mathcal{L}, ρ_0 satisfying

$$(A.31) \quad \begin{cases} \mathcal{L}(\sqrt{r} - 1) \geq \rho_0 \\ \rho_0 \geq \rho_+ (2q) - \mathcal{L}. \end{cases}$$

Furthermore, for any $\mathcal{L} \geq \rho_+ (2q) / \sqrt{r}$, we can always choose $\rho_0 = (1 - 1/\sqrt{r}) \rho_+ (2q)$. The rest of the proof proceeds as before.

REMARK A.3. The idea of controlling the overall $\sum_{t \leq T} \frac{R_t}{\theta_t^2 \rho_t}$ can be extended with a proper choice of ψ_0 to a general problem $\min f(\beta)$ that may be nonconvex. In fact, if $f(\beta)$ can be decomposed as $l(\beta) + P(\beta)$ with $0 \leq \Delta_l \leq L \mathbf{D}_2$ and $\Delta_P + \mathcal{L}_0 \mathbf{D}_2 \geq 0$ for some

finite $\mathcal{L}_0 \geq 0$, then setting $\mu_0 = 0$, $\psi_0 = l - \mathcal{L} \|\cdot\|_2^2/2$ with $\mathcal{L} \geq \mathcal{L}_0$ and repeating the previous arguments, we obtain

$$\begin{aligned}
 \sum_{t=0}^T \frac{R_t}{\theta_t^2 \rho_t} &\geq \sum_{t=0}^T \frac{1}{\rho_t} (\rho_t \mathbf{D}_\phi + \mathcal{L} \mathbf{D}_2 - L \mathbf{D}_2) (\boldsymbol{\alpha}^{(t+1)}, \boldsymbol{\alpha}^{(t)}) \\
 (A.32) \quad &+ \sum_{t=0}^{T-1} \frac{1}{\theta_t^2 \rho_t} \Delta_l(\boldsymbol{\beta}^{(t+1)}, \boldsymbol{\gamma}^{(t+1)}) + \frac{(\mathcal{L} - \mathcal{L}_0)(1 - \theta_T)}{\theta_T^3 \rho_T} \mathbf{D}_2(\boldsymbol{\beta}^{(T+1)}, \boldsymbol{\beta}^{(T)}) \\
 &+ \sum_{t=0}^{T-1} \mathcal{L} \frac{1 - \theta_t}{\theta_t \rho_t} \left\{ \frac{\mathcal{L} - \mathcal{L}_0}{\mathcal{L}} - (1 - \theta_t) \theta_{t+1} \frac{\theta_{t+1}}{\theta_t} \right\} \mathbf{D}_2(\boldsymbol{\alpha}^{(t+1)}, \boldsymbol{\beta}^{(t)}).
 \end{aligned}$$

With $\rho_t = \rho_0$, it can be shown that $(1 - \theta_t) \theta_{t+1} \frac{\theta_{t+1}}{\theta_t}$ achieves the maximum value 0.1612 at $t = 3$ and so $\mathcal{L} \geq 1.2\mathcal{L}_0$ makes the last term nonnegative. (A varying ρ_t may reduce \mathcal{L} further.) Therefore, under $\mathbf{D}_\phi \geq \sigma \mathbf{D}_2$, we can choose any $\rho_0 \geq (L - \mathcal{L})/\sigma$ and $\mathcal{L} \geq 1.2\mathcal{L}_0$ so that for any $\boldsymbol{\beta}$,

$$\begin{aligned}
 \frac{f(\boldsymbol{\beta}^{(T+1)}) - f(\boldsymbol{\beta})}{\theta_T^2} + T \cdot \text{avg}_{0 \leq t \leq T} \frac{\mathcal{E}_t(\boldsymbol{\beta})}{\theta_t} + \frac{(\mathcal{L} - \mathcal{L}_0)(1 - \theta_T)}{\theta_T^3} \mathbf{D}_2(\boldsymbol{\beta}^{(T+1)}, \boldsymbol{\beta}^{(T)}) \\
 \leq \rho_0 \mathbf{D}_2(\boldsymbol{\beta}, \boldsymbol{\alpha}^{(0)}), \forall T \geq 0.
 \end{aligned}$$

Hence for some A independent of T , $f(\boldsymbol{\beta}^{(T+1)}) - f(\boldsymbol{\beta}) + \min_{0 \leq t \leq T} \mathcal{E}_t(\boldsymbol{\beta}) \leq A/T^2$, or $f(\boldsymbol{\beta}^{(T+1)}) - f(\boldsymbol{\beta}) + \min_{0 \leq t \leq T} \{\Delta_{\psi_0}(\boldsymbol{\beta}, \boldsymbol{\gamma}^{(t)}) + \Delta_{f(\cdot) - \Delta_{\psi_0}(\cdot, \boldsymbol{\gamma}^{(t)})}(\boldsymbol{\beta}, \boldsymbol{\alpha}^{(t+1)})\} \leq A/T^2$ when ϕ is differentiable.

A.15. Proof of Theorem 7. The construction of the new acceleration scheme and the proof are motivated by Proposition 2 of [57], with the use of GBF calculus. First, from Lemma A.2, given any $\boldsymbol{\beta}'_t$,

$$\begin{aligned}
 (A.33) \quad &f(\boldsymbol{\beta}^{(t+1)}) - f(\boldsymbol{\beta}'_t) + (\rho_t \mathbf{D}_2 - \Delta_{\bar{\psi}_0})(\boldsymbol{\beta}^{(t+1)}, \boldsymbol{\gamma}^{(t)}) + \Delta_f(\boldsymbol{\beta}'_t, \boldsymbol{\beta}^{(t+1)}) \\
 &\leq (\rho_t \mathbf{D}_2 - \Delta_{\bar{\psi}_0})(\boldsymbol{\beta}'_t, \boldsymbol{\gamma}^{(t)}) - (\rho_t \mathbf{D}_2 - \Delta_{\Delta_{\bar{\psi}_0}(\cdot, \boldsymbol{\gamma}^{(t)})})(\boldsymbol{\beta}'_t, \boldsymbol{\beta}^{(t+1)})
 \end{aligned}$$

Let $\boldsymbol{\beta}'_t = \theta_t \boldsymbol{\beta} + (1 - \theta_t) \boldsymbol{\beta}^{(t)}$ with θ_t to be determined. Define $h_t(\cdot) = f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \boldsymbol{\gamma}^{(t)})$. By the definition of \mathbf{C} ,

$$\begin{aligned}
 -f(\boldsymbol{\beta}'_t) &= \theta_t \Delta_{\bar{\psi}_0}(\boldsymbol{\beta}, \boldsymbol{\gamma}^{(t)}) + (1 - \theta_t) \Delta_{\bar{\psi}_0}(\boldsymbol{\beta}^{(t)}, \boldsymbol{\gamma}^{(t)}) - \Delta_{\bar{\psi}_0}(\boldsymbol{\beta}'_t, \boldsymbol{\gamma}^{(t)}) \\
 &\quad - \theta_t f(\boldsymbol{\beta}) - (1 - \theta_t) f(\boldsymbol{\beta}^{(t)}) + \mathbf{C}_{h_t}(\boldsymbol{\beta}, \boldsymbol{\beta}^{(t)}, \theta_t).
 \end{aligned}$$

Plugging the last equality into (A.33) yields

$$\begin{aligned}
 &f(\boldsymbol{\beta}^{(t+1)}) - f(\boldsymbol{\beta}) - (1 - \theta_t)(f(\boldsymbol{\beta}^{(t)}) - f(\boldsymbol{\beta})) \\
 &\quad + (\rho_t \mathbf{D}_2 - \Delta_{\bar{\psi}_0})(\boldsymbol{\beta}^{(t+1)}, \boldsymbol{\gamma}^{(t)}) + \mathbf{C}_{h_t}(\boldsymbol{\beta}, \boldsymbol{\beta}^{(t)}, \theta_t) \\
 &\quad + \theta_t \Delta_{\bar{\psi}_0}(\boldsymbol{\beta}, \boldsymbol{\gamma}^{(t)}) + (1 - \theta_t) \Delta_{\bar{\psi}_0}(\boldsymbol{\beta}^{(t)}, \boldsymbol{\gamma}^{(t)}) \\
 &\leq (\rho_t \mathbf{D}_2 - \Delta_{\bar{\psi}_0})(\boldsymbol{\beta}'_t, \boldsymbol{\gamma}^{(t)}) - (\rho_t \mathbf{D}_2 - \Delta_{\Delta_{\bar{\psi}_0}(\cdot, \boldsymbol{\gamma}^{(t)})})(\boldsymbol{\beta}'_t, \boldsymbol{\beta}^{(t+1)}) \\
 &\quad + \Delta_{\bar{\psi}_0}(\boldsymbol{\beta}'_t, \boldsymbol{\gamma}^{(t)}) - \Delta_f(\boldsymbol{\beta}'_t, \boldsymbol{\beta}^{(t+1)}) \\
 &= \rho_t [\mathbf{D}_2(\boldsymbol{\beta}'_t, \boldsymbol{\gamma}^{(t)}) - \mathbf{D}_2(\boldsymbol{\beta}'_t, \boldsymbol{\beta}^{(t+1)})] - \Delta_{h_t}(\boldsymbol{\beta}'_t, \boldsymbol{\beta}^{(t+1)})
 \end{aligned}$$

and based on the definition of $\bar{\psi}_0$ and R_t ,

$$\begin{aligned}
& f(\beta^{(t+1)}) - f(\beta) + R_t \\
& + \rho_t \mathbf{D}_2(\beta'_t, \beta^{(t+1)}) + \mathbf{C}_{\mu_0 \mathbf{D}_2(\cdot, \gamma^{(t)})}(\beta, \beta^{(t)}, \theta_t) + \mu_0 \mathbf{D}_2(\beta'_t, \beta^{(t+1)}) \\
& + \theta_t \Delta_{\bar{\psi}_0}(\beta, \gamma^{(t)}) + \mathbf{C}_{f(\cdot) - \Delta_{\psi_0}(\cdot, \gamma^{(t)})}(\beta, \beta^{(t)}, \theta_t) + \Delta_{f(\cdot) - \Delta_{\psi_0}(\cdot, \gamma^{(t)})}(\beta'_t, \beta^{(t+1)}) \\
& \leq (1 - \theta_t)(f(\beta^{(t)}) - f(\beta)) + \rho_t \mathbf{D}_2(\beta'_t, \gamma^{(t)}).
\end{aligned}$$

From Section A.5, $\mathbf{C}_{\mu_0 \mathbf{D}_2(\cdot, \gamma^{(t)})}(\beta, \beta^{(t)}, \theta_t) = \mu_0 \mathbf{C}_2(\beta, \beta^{(t)}, \theta_t) = \mu_0 \theta_t (1 - \theta_t) \mathbf{D}_2(\beta, \beta^{(t)})$ and so

$$\begin{aligned}
& f(\beta^{(t+1)}) - f(\beta) + R_t + \theta_t \mathcal{E}_t(\beta) \\
& + (\rho_t + \mu_0) \mathbf{D}_2(\beta'_t, \beta^{(t+1)}) + \mu_0 \theta_t (1 - \theta_t) \mathbf{D}_2(\beta, \beta^{(t)}) \\
& \leq (1 - \theta_t)(f(\beta^{(t)}) - f(\beta)) + \rho_t \mathbf{D}_2(\beta'_t, \gamma^{(t)}).
\end{aligned} \tag{A.34}$$

We would like to write $(\rho_t + \mu_0) \mathbf{D}_2(\theta_t \beta + (1 - \theta_t) \beta^{(t)}, \beta^{(t+1)}) + \mu_0 \theta_t (1 - \theta_t) \mathbf{D}_2(\beta, \beta^{(t)})$ into the form of a multiple of $\mathbf{D}_2(\beta, \nu^{(t+1)})$ for some $\nu^{(t+1)}$. This can be done by solving the gradient equation with respect to β :

$$\nu^{(t+1)} = \frac{(\rho_t + \mu_0) \theta_t \beta^{(t+1)} - \rho_t \theta_t (1 - \theta_t) \beta^{(t)}}{\rho_t \theta_t^2 + \mu_0 \theta_t}. \tag{A.35}$$

On the other hand, $\nabla \rho_t \mathbf{D}_2(\theta_t \beta + (1 - \theta_t) \beta^{(t)}, \gamma^{(t)}) = \mathbf{0}$ gives

$$\nu^{(t)} = \frac{\gamma^{(t)}}{\theta_t} - \frac{1 - \theta_t}{\theta_t} \beta^{(t)}. \tag{A.36}$$

Combining (A.35) and (A.36) results in

$$\gamma^{(t)} = \beta^{(t)} + \frac{\rho_{t-1} \theta_t (1 - \theta_{t-1})}{\rho_{t-1} \theta_{t-1} + \mu_0} (\beta^{(t)} - \beta^{(t-1)}), \tag{A.37}$$

as in (76a). Therefore, (A.34) becomes

$$\begin{aligned}
& f(\beta^{(t+1)}) - f(\beta) + (\theta_t^2 \rho_t + \mu_0 \theta_t) \mathbf{D}_2(\beta, \nu^{(t+1)}) + R_t + \theta_t \mathcal{E}_t(\beta) \\
& \leq (1 - \theta_t)(f(\beta^{(t)}) - f(\beta)) + \theta_t^2 \rho_t \mathbf{D}_2(\beta, \nu^{(t)}).
\end{aligned} \tag{A.38}$$

Let $\mu_0 = 0$. It follows from (A.38) that

$$\begin{aligned}
& \frac{1}{\theta_t^2 \rho_t} [f(\beta^{(t+1)}) - f(\beta)] + \mathbf{D}_2(\beta, \nu^{(t+1)}) + \frac{\mathcal{E}_t(\beta)}{\theta_t \rho_t} + \frac{R_t}{\theta_t^2 \rho_t} \\
& \leq \frac{(1 - \theta_t)}{\theta_t^2 \rho_t} [f(\beta^{(t)}) - f(\beta)] + \mathbf{D}_2(\beta, \nu^{(t)}), \forall t \geq 0.
\end{aligned} \tag{A.39}$$

Under (77b), we have

$$\begin{aligned}
& \frac{1}{\theta_t^2 \rho_t} [f(\beta^{(t+1)}) - f(\beta)] + \mathbf{D}_2(\beta, \nu^{(t+1)}) + \frac{\mathcal{E}_t(\beta)}{\theta_t \rho_t} + \frac{R_t}{\theta_t^2 \rho_t} \\
& \leq \frac{1}{\theta_{t-1}^2 \rho_{t-1}} [f(\beta^{(t)}) - f(\beta)] + \mathbf{D}_2(\beta, \nu^{(t)}), \quad \forall t \geq 1.
\end{aligned} \tag{A.40}$$

Summing (A.40) for $t = T, \dots, 1$ and (A.39) for $t = 0$ gives

$$\begin{aligned} & \frac{1}{\theta_T^2 \rho_T} [f(\beta^{(T+1)}) - f(\beta)] + \sum_{t=0}^T \left(\frac{\mathcal{E}_t(\beta)}{\theta_t \rho_t} + \frac{R_t}{\theta_t^2 \rho_t} \right) \\ & \leq \frac{1 - \theta_0}{\theta_0^2 \rho_0} [f(\beta^{(0)}) - f(\beta)] + \mathbf{D}_2(\beta, \nu^{(0)}) - \mathbf{D}_2(\beta, \nu^{(T+1)}), \end{aligned}$$

and so the first bound noticing that $\nu^{(0)} = \gamma^{(0)} = \beta^{(0)}$.

Moreover, given any $\mu_0 \geq 0$, from (77b), (A.38) implies for any $t \geq 1$,

$$\begin{aligned} & f(\beta^{(t+1)}) - f(\beta) + (\theta_t^2 \rho_t + \mu_0 \theta_t) \mathbf{D}_2(\beta, \nu^{(t+1)}) + R_t + \theta_t \mathcal{E}_t(\beta) \\ (A.41) \quad & \leq (1 - \theta_t) [f(\beta^{(t)}) - f(\beta) + (\theta_{t-1}^2 \rho_{t-1} + \mu_0 \theta_{t-1}) \mathbf{D}_2(\beta, \nu^{(t)})]. \end{aligned}$$

Similar to the proof of Theorem 6, a recursive argument using (A.41) and (A.38) gives the second bound.

REMARK A.4. Compared with the proof of Theorem 6, the proof here needs to perform a finer analysis of \mathbf{C}_{h_t} (the proof of Corollary 2' uses a similar treatment). Otherwise one would get $\gamma^{(t)} = \beta^{(t)} + \theta_t(\theta_{t-1}^{-1} - 1)(\beta^{(t)} - \beta^{(t-1)})$ and $\rho_t \theta_t^2 / (1 - \theta_t) = \theta_{t-1}^2 (\rho_{t-1} + \mu_0)$ in place of (76a), (77b), respectively. Following the same proof, we can show that the resultant algorithm does result in a linear rate when $\mu_0 = \mu > 0$, but offers no acceleration ($\theta_0 = 1/\kappa$) in strongly smooth and convex optimization.

Finally, Remark A.2 still applies. For example, the second conclusion holds when the '=' in (77b) is replaced by ' \leq ', and it is straightforward to see that μ_0 can be similarly replaced by a sequence of varying μ_t to speed the convergence.

A.16. Statistical accuracy of LLA iterates. In this subsection, assume $f(\beta) = l(\beta) + P(\varrho\beta)$, $l(\beta) = l_0(\mathbf{X}\beta)$, $P(\varrho\beta) = \sum_j P(\varrho\beta_j)$ (by a slight abuse of notation), $P(0) = 0$, $P'_+(0) < +\infty$, $P(t) = P(-t) \geq 0$, $P(t)$ is differentiable for any $t > 0$, and P is concave on $(0, +\infty)$. Recall $\Delta_{\text{LLA}}^{(t)} = \Delta_{\|\alpha^{(t)} \circ (\cdot)\|_1 - P(\cdot)}$ which does *not* satisfy the strong idempotence.

ASSUMPTION $\mathcal{A}(\varepsilon, \delta, \vartheta, K, \alpha, \beta)$ Given $\mathbf{X}, \alpha, \beta$, there exist $\varepsilon > 0, \delta > 0, \vartheta > 0, K \geq 0$ such that the following inequality holds

$$\begin{aligned} & (1 + \varepsilon) \Delta_{\|\alpha \circ (\cdot)\|_1 - P(\cdot)}(\varrho\beta^*, \varrho\beta) + \delta \mathbf{D}_2(\mathbf{X}\beta^*, \mathbf{X}\beta) + \vartheta P_H(\varrho(\beta - \beta^*); \lambda) \\ & \leq 2\bar{\Delta}_l(\beta^*, \beta) + P(\varrho\beta; \lambda) - P(\varrho\beta^*; \lambda) + K\lambda^2 J^*. \end{aligned}$$

PROPOSITION 5. Assume that for any given $T \geq 1$, $\mathcal{A}(\varepsilon, \delta, \vartheta, K, \alpha^{(t)}, \beta^{(t)})$ ($1 \leq t \leq T$) is satisfied for some $\varepsilon > 0, \delta > 0, \vartheta > 0, K \geq 0$. Let $\lambda = A\sigma\sqrt{\log(ep)}/\sqrt{(\delta \wedge \vartheta)\vartheta}$. Then the following inequality holds with probability at least $1 - Cp^{-cA^2}$

$$\Delta_{\text{LLA}}^{(T)}(\varrho\beta^*, \varrho\beta^{(T)}) \leq \kappa^T \Delta_{\text{LLA}}^{(0)}(\varrho\beta^*, \varrho\beta^{(0)}) + \frac{\kappa}{1 - \kappa} K\lambda^2 J^*,$$

where $\kappa = 1/(1 + \varepsilon)$ and C, c are universal positive constants.

PROOF. From the proof of Proposition 4, for any β ,

$$\begin{aligned} & \Delta_f(\beta, \beta^{(t+1)}) + \Delta_{\Delta_{\text{LLA}}^{(t)}(\cdot, \varrho\beta^{(t)})}(\varrho\beta, \varrho\beta^{(t+1)}) \\ & \leq f(\beta) - f(\beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta, \varrho\beta^{(t)}) - \Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t+1)}, \varrho\beta^{(t)}). \end{aligned}$$

Using the definition of $\Delta_{\text{LLA}}^{(t)}$, we have

$$\begin{aligned}
 & \Delta_P(\varrho\beta, \varrho\beta^{(t+1)}) - \Delta_{\Delta_{P(\cdot, \varrho\beta^{(t)})}}(\varrho\beta, \varrho\beta^{(t+1)}) \\
 & + \Delta_l(\beta, \beta^{(t+1)}) + \sum_j \alpha_j^{(t)} \Delta_{\Delta_1(\cdot, \varrho\beta_j^{(t)})}(\varrho\beta_j, \varrho\beta_j^{(t+1)}) \\
 & \leq f(\beta) - f(\beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta, \varrho\beta^{(t)}),
 \end{aligned}
 \tag{A.42}$$

where we used $\Delta_{\text{LLA}}^{(t)}(\varrho\beta^{(t+1)}, \varrho\beta^{(t)}) \geq 0$ since $P(\cdot)$ is concave on $(0, +\infty)$.

LEMMA A.7. *For any $P(\cdot)$ which is differentiable on $(0, +\infty)$ and satisfies $P(t) = P(-t) \geq 0$, $P(0) = 0$ and $P'_+(0) < +\infty$, we have $\Delta_{\Delta_P(\cdot, \alpha)}(\beta, \gamma) = \Delta_P(\beta, \gamma) - P'_+(0)\Delta_1(\beta, \gamma)1_{\alpha=0}$ for any $\alpha, \beta, \gamma \in \mathbb{R}$. In particular, $\Delta_{\Delta_1(\cdot, \alpha)}(\beta, \gamma) = \Delta_1(\beta, \gamma)1_{\alpha \neq 0}$.*

The result can be shown from the proof of Lemma 2. Indeed, from (29),

$$\Delta_P(\cdot, \alpha) = \begin{cases} P(\cdot) - P(\alpha) - P'(\alpha)(\cdot - \alpha), & \alpha \neq 0 \\ P(\cdot) - P'_+(0)|\cdot|, & \alpha = 0. \end{cases}$$

When $\alpha \neq 0$, by Lemma 1 and Lemma 2, $\Delta_{\Delta_P(\cdot, \alpha)}(\beta, \gamma) = \Delta_P(\beta, \gamma) - \Delta_{P(\alpha) + P'(\alpha)(\cdot - \alpha)}(\beta, \gamma) = \Delta_P(\beta, \gamma)$. When $\alpha = 0$, $\Delta_{\Delta_P(\cdot, \alpha)}(\beta, \gamma) = \Delta_P(\beta, \gamma) - P'_+(0)\Delta_1(\beta, \gamma)$. Combining the two cases gives

$$\Delta_{\Delta_P(\cdot, \alpha)}(\beta, \gamma) = \Delta_P(\beta, \gamma) - P'_+(0)\Delta_1(\beta, \gamma)1_{\alpha=0}.$$

When $P(\beta) = \|\beta\|_1$, $\Delta_{\Delta_1(\cdot, \alpha)}(\beta, \gamma) = \Delta_1(\beta, \gamma) - \Delta_1(\beta, \gamma)1_{\alpha=0} = \Delta_1(\beta, \gamma)1_{\alpha \neq 0}$.

From Lemma A.7,

$$\Delta_P(\varrho\beta, \varrho\beta^{(t+1)}) - \Delta_{\Delta_{P(\cdot, \varrho\beta^{(t)})}}(\varrho\beta, \varrho\beta^{(t+1)}) = \sum_{j: \beta_j^{(t)} = 0} P'_+(0)\Delta_1(\varrho\beta_j, \varrho\beta_j^{(t+1)})$$

and

$$\sum_j \alpha_j^{(t)} \Delta_{\Delta_1(\cdot, \varrho\beta_j^{(t)})}(\varrho\beta_j, \varrho\beta_j^{(t+1)}) = \sum_{j: \beta_j^{(t)} \neq 0} \alpha_j^{(t)} \Delta_1(\varrho\beta_j, \varrho\beta_j^{(t+1)}).$$

Plugging these into (A.42) gives $\Delta_l(\beta, \beta^{(t+1)}) + \sum_{j: \beta_j^{(t)} = 0} P'_+(0)\Delta_1(\varrho\beta_j, \varrho\beta_j^{(t+1)}) + \sum_{j: \beta_j^{(t)} \neq 0} \alpha_j^{(t)} \Delta_1(\varrho\beta_j, \varrho\beta_j^{(t+1)}) \leq f(\beta) - f(\beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta, \varrho\beta^{(t)})$.

Together with $\alpha_j^{(t)} = |P'_+(\beta_j^{(t)})| \leq P'_+(0)$, we have

$$\Delta_l(\beta, \beta^{(t+1)}) + \sum_j \alpha_j^{(t)} \Delta_1(\varrho\beta_j, \varrho\beta_j^{(t+1)}) \leq f(\beta) - f(\beta^{(t+1)}) + \Delta_{\text{LLA}}^{(t)}(\varrho\beta, \varrho\beta^{(t)}).
 \tag{A.43}$$

Letting $\beta = \beta^*$ and using the definition of ϵ , we obtain

$$\begin{aligned}
 & 2\bar{\Delta}_l(\beta^*, \beta^{(t+1)}) + \Delta_{\|\alpha^{(t)} \circ (\cdot)\|_1}(\varrho\beta^*, \varrho\beta^{(t+1)}) + P(\varrho\beta^{(t+1)}; \lambda) \\
 & \leq \Delta_{\text{LLA}}^{(t)}(\varrho\beta^*, \varrho\beta^{(t)}) + \langle \epsilon, X\beta^{(t+1)} - X\beta^* \rangle + P(\varrho\beta^*; \lambda).
 \end{aligned}$$

From the regularity condition,

$$\begin{aligned}
 & (1 + \varepsilon)\Delta_{\text{LLA}}^{(t+1)}(\varrho\beta^*, \varrho\beta^{(t+1)}) + \delta D_2(X\beta^*, X\beta^{(t+1)}) + \vartheta P_H(\varrho(\beta^{(t+1)} - \beta^*); \lambda) \\
 & \leq 2\bar{\Delta}_l(\beta^*, \beta^{(t+1)}) + \Delta_{\|\alpha^{(t)} \circ (\cdot)\|_1}(\varrho\beta^*, \varrho\beta^{(t+1)}) + P(\varrho\beta^{(t+1)}; \lambda) - P(\varrho\beta^*; \lambda) + K\lambda^2 J^*
 \end{aligned}$$

for $1 \leq t \leq T$. The final conclusion can be proved by combining the last two inequalities and then applying a similar probabilistic argument as in Theorem 5. \square

A.17. A-estimators as F-estimators. In this part, we show that an important class of A-estimators that has *alternative* optimality, typically arising from block coordinate descent (BCD) algorithms, can often be converted to F-estimators, and analyzed in a similar way. Let $\beta = [\beta_{[1]}^T, \dots, \beta_{[K]}^T]^T$ where $\beta_{[k]}$ is the k th block, $1 \leq k \leq K$, and we use $\beta_{[-k]}$ to denote the subvector after removing the k th block. Assume

$$f = l + P$$

where l is differentiable, and P is separable: $P(\beta) = \sum P_k(\beta_{[k]})$. When viewed as a function of $\beta_{[k]}$ only, f is denoted by $f(\beta_{[k]}; \beta_{[-k]})$. We say $\hat{\beta}$ has alternative optimality or is an A-estimator if

$$(A.44) \quad \hat{\beta}_{[k]} \in \arg \min_{\beta_{[k]}} f(\beta_{[k]}; \hat{\beta}_{[-k]}), 1 \leq k \leq K.$$

LEMMA A.8. *Let $\hat{\beta}$ be an A-estimator of $\min f(\beta)$. Construct a surrogate function:*

$$(A.45) \quad g_{\rho}(\beta; \beta^-) = f(\beta) - \Delta_l(\beta, \beta^-) + \sum \rho_k \mathbf{D}_2(\beta_{[k]}, \beta_{[k]}^-)$$

where $\rho = (\rho_1, \dots, \rho_K)$ with $\rho_k \geq 0$.

(i) If P_k are directionally differentiable and

$$(A.46) \quad \Delta_{P_k} + \mathcal{L}_k \mathbf{D}_2 \geq 0$$

for some $\mathcal{L}_k \geq 0$, then for any $\rho_k \geq \mathcal{L}_k$, $\hat{\beta}$ must satisfy

$$(A.47) \quad \hat{\beta} \in \arg \min_{\beta} g_{\rho}(\beta; \beta^-)|_{\beta^- = \hat{\beta}}.$$

(ii) If l as a function of β_k satisfies $\Delta_{l(\cdot; \beta_{[-k]})} \leq L_k \mathbf{D}_2, \forall \beta_{[-k]}, 1 \leq k \leq K$, or less restrictively,

$$(A.48) \quad \Delta_{\hat{l}_k(\cdot)}(\beta_{[k]}, \hat{\beta}_{[k]}) \leq L_k \mathbf{D}_2(\beta_{[k]}, \hat{\beta}_{[k]}), \forall \beta_{[k]}$$

where $\hat{l}_k(\beta_{[k]})$ denotes $l_k(\beta_{[k]}; \hat{\beta}_{[-k]})$, then for any $\rho_k > L_k$, (A.47) still holds. In addition, if $\hat{\beta}_{[k]}$ is the unique solution to (A.44), then $\rho_k \geq L_k$ suffices.

Overall, (A.47) provides a useful *joint* optimization form that can be used as the so-called “basic inequality” in empirical process theory, and so with the lemma, A-estimators can be analyzed like F-estimators. Moreover, the quality of the initial point can be incorporated in the analysis; see [54].

PROOF. (i) The condition (A.46) means that g_{ρ} is convex in $\beta_{[k]}$. By Lemma 4 and Lemma 1, and the fact that g_{ρ} is separable in $\beta_{[k]}, 1 \leq k \leq K$, we immediately know that $\hat{\beta}$ is necessarily a solution to $\min_{\beta} g_{\rho}(\beta; \hat{\beta})$.

(ii) We use a shorthand notation $\hat{g}_k(\beta_{[k]})$ to denote $g_{\rho}(\beta; \hat{\beta})$ as a function of $\beta_{[k]}$ when $\beta_{[-k]} = \hat{\beta}_{[-k]}$. Let $\tilde{\beta}_{[k]} \in \arg \min_{\beta_{[k]}} \hat{g}_k(\beta_{[k]})$. It suffices to show $\tilde{\beta}_{[k]} = \hat{\beta}_{[k]}$. Because of the separability of g_{ρ} ,

$$f(\tilde{\beta}_{[k]}; \hat{\beta}_{[-k]}) + (\rho_k - L_k) \mathbf{D}_2(\hat{\beta}_{[k]}, \tilde{\beta}_{[k]}) \leq f(\hat{\beta}_{[k]}; \hat{\beta}_{[-k]})$$

and so $(\rho_k - L_k) \mathbf{D}_2(\hat{\beta}_{[k]}, \tilde{\beta}_{[k]}) = 0$. The conclusion follows. \square

Some conclusions like (i) can be extended to functions defined on Riemannian manifolds. It is also worth mentioning that in the regression setup, which is of primary interest in many statistical applications, we can use some surrogates with $\rho_k = 1$, regardless of the design or penalty, to convert alternative optimality to joint optimality. The following lemma exemplifies the point in matrix regression, and is condition free.

LEMMA A.9. *Let $l_0(\mathbf{A}; \mathbf{Y}) = \|\mathbf{Y} - \mathbf{A}\|_F^2/2$, and \mathbf{A} be defined differently as follows.*

(i) *Let $\mathbf{A} = \sum \mathbf{X}_k \mathbf{B}_k$ with $\mathbf{B} = (\mathbf{B}_1, \dots, \mathbf{B}_K)$, where the dependence of \mathbf{B} (and \mathbf{X}_k) is dropped for simplicity. Consider the problem*

$$(A.49) \quad \min_{\mathbf{B}_1, \dots, \mathbf{B}_K} l_0(\mathbf{A}; \mathbf{Y}) + \sum P_k(\mathbf{B}_k) \text{ s.t. } \mathbf{A} = \sum \mathbf{X}_k \mathbf{B}_k.$$

Then the set of A-estimators of (A.49) is exactly the set of F-estimators associated with the following surrogate

$$(A.50) \quad g(\mathbf{B}, \mathbf{B}^-) = l_0(\mathbf{A}; \mathbf{Y}) - \mathbf{D}_{l_0}(\mathbf{A}, \mathbf{A}^-) + \sum P_k(\mathbf{B}_k) + \sum \mathbf{D}_2(\mathbf{X}_k \mathbf{B}_k, \mathbf{X}_k \mathbf{B}_k^-).$$

(ii) *Let $\mathbf{A} = \mathbf{X} \mathbf{B}_1 \cdots \mathbf{B}_K$ with $\mathbf{B} = (\mathbf{B}_1, \dots, \mathbf{B}_K)$, where the dependence of \mathbf{B} (and \mathbf{X}_k) is dropped for simplicity. Consider*

$$(A.51) \quad \min_{\mathbf{B}_1, \dots, \mathbf{B}_K} l_0(\mathbf{A}; \mathbf{Y}) + \sum P_k(\mathbf{B}_k) \text{ s.t. } \mathbf{A} = \mathbf{X} \mathbf{B}_1 \cdots \mathbf{B}_K.$$

Redefine l_0 as a function \bar{l}_0 of \mathbf{B} and introduce a discrepancy measure d^2 as follows

$$l_0(\mathbf{A}; \mathbf{Y}) = \bar{l}_0(\mathbf{B}; \mathbf{X}, \mathbf{Y})$$

$$d^2(\mathbf{B}, \mathbf{B}^-) = \frac{1}{2} \sum_{k=1}^K \|\mathbf{X} \mathbf{B}_1^- \cdots \mathbf{B}_{k-1}^- (\mathbf{B}_k - \mathbf{B}_k^-) \mathbf{B}_{k+1}^- \cdots \mathbf{B}_K^-\|_F^2.$$

Then the set of A-estimators of (A.51) is exactly the set of F-estimators associated with the surrogate

$$(A.52) \quad \begin{aligned} g(\mathbf{B}, \mathbf{B}^-) &= \bar{l}_0(\mathbf{B}) - \Delta_{\bar{l}_0}(\mathbf{B}, \mathbf{B}^-) + \sum P_k(\mathbf{B}_k) + d^2(\mathbf{B}, \mathbf{B}^-) \\ &= \frac{1}{2} \|\mathbf{A}^- - \mathbf{Y}\|_F^2 + \langle \mathbf{X} \mathbf{B}_1^- \cdots \mathbf{B}_K^- - \mathbf{Y}, \\ &\quad \sum_{k=1}^K \mathbf{X} \mathbf{B}_1^- \cdots \mathbf{B}_{k-1}^- (\mathbf{B}_k - \mathbf{B}_k^-) \mathbf{B}_{k+1}^- \cdots \mathbf{B}_K^- \rangle + \sum P_k(\mathbf{B}_k) \\ &\quad \frac{1}{2} \sum_{k=1}^K \|\mathbf{X} \mathbf{B}_1^- \cdots \mathbf{B}_{k-1}^- (\mathbf{B}_k - \mathbf{B}_k^-) \mathbf{B}_{k+1}^- \cdots \mathbf{B}_K^-\|_F^2. \end{aligned}$$

The lemma can be directly proved by the definition of GBF and matrix differentiation and its proof is omitted. For the application of the first result (i), see [50] for example. The second result can be used to study bilinear problems or NMF like matrix decomposition problems. One could show a statistical accuracy result in terms of d^1 (which satisfies $d^1 \leq K d^2$),

$$d^1(\mathbf{B}, \mathbf{B}^-) = \frac{1}{2} \left\| \sum_{k=1}^K \mathbf{X} \mathbf{B}_1^- \cdots \mathbf{B}_{k-1}^- (\mathbf{B}_k - \mathbf{B}_k^-) \mathbf{B}_{k+1}^- \cdots \mathbf{B}_K^- \right\|_F^2,$$

under a proper regularity condition involving d^2 ; see, for example, [53].

A.18. Statistical error analysis of a general optimal solution. This part demonstrates that using the statistical notions and Bregman calculus developed earlier can perform statistical analysis of a general optimization problem that may not be in the MLE setup:

$$(A.53) \quad \min_{\beta} f(\beta) \text{ s.t. } \beta \in \mathcal{S}$$

where f is directionally differentiable and $\mathcal{S} \subset \mathbb{R}^p$ can be formulated by linear equality constraints $A\beta = \alpha$, sparsity constraints $\|\beta\|_0 \leq s$, nonnegativity constraints $\beta \geq 0$, and so on.

Statistically, we would like to study how a target parameter can be recovered from solving (A.53) in the presence of data noise. Following (53), let β^* be a statistical truth and define the associated effective noise by $\epsilon = -\nabla f(\beta^*)$, assuming f is differentiable at β^* .

Although β^* in the above definition can be any point, a meaningful recovery must be under some conditions satisfied the associated ϵ . Consider the following three scenarios:

(a) Statistical estimation often assumes a zero mean noise:

$$(A.54) \quad \mathbb{E}\epsilon = 0,$$

which essentially means that the statistical truth makes the gradient of the *expectation* of f (so as to remove data randomness) vanish—see Section 3.2. Yet (A.54) alone does not always guarantee a unique β^* .

(b) Stronger conclusions can be obtained for the β^* that satisfies the no-model-ambiguity assumption: f is differentiable at $\beta^* \in D = \text{dom}(f)$ with the gradient $\nabla f(\beta^*) = -\epsilon$, β^* is a finite optimal solution to the Fenchel conjugate as $\zeta = -\epsilon$:

$$(A.55) \quad f^*(\zeta) = \sup_{\beta} \langle \zeta, \beta \rangle - f(\beta),$$

and the extended real-valued convex function f^* is differentiable at $-\epsilon$. This assumption simply means that $(\beta^*, -\epsilon)$ makes a so-called “conjugate pair”. Note that f need not be overall strictly convex, especially when D is compact according to Danskin’s min-max theorem [6].

(c) Another popular assumption in statistical learning is strong convexity in a restricted sense (especially when $\beta = A\alpha$ with $\|\alpha\|_0 \leq s$):

$$(A.56) \quad (\Delta_f - \mu D_2)(\beta_1, \beta_2) \geq 0, \quad \forall \beta_1, \beta_2 \in \mathcal{S}$$

for some $\mu > 0$. The condition may hold even when the number of unknowns is much larger than the sample size [12, 36].

The following theorem uses the GBF calculus to argue how the statistical accuracy of the obtained solutions is determined by the (tail decay of) effective noise. Probabilistic arguments can follow to bound the stochastic terms more explicitly.

THEOREM A.1. *Let $\hat{\beta}$ be an optimal solution to (A.53).*

(i) *Under the zero mean assumption (A.54) and $\beta^* \in \mathcal{S}$, the risk of $\hat{\beta}$ in terms of Δ_f satisfies a Fenchel-Young form bound*

$$(A.57) \quad \mathbb{E}\Delta_f(\hat{\beta}, \beta^*) \leq \mathbb{E}[f^*(\epsilon) + f(\beta^*)].$$

(ii) *Under the no-model-ambiguity assumption in (b) and $\beta^* \in \mathcal{S}$, we have*

$$(A.58) \quad \Delta_f(\hat{\beta}, \beta^*) \leq \Delta_{f^*}(\epsilon, -\epsilon).$$

(iii) *An oracle inequality holds for any $\delta > 0$ and any reference $\beta \in \mathcal{S}$:*

$$(A.59) \quad (\Delta_f - \delta D_2)(\hat{\beta}, \beta^*) \leq \Delta_f(\beta, \beta^*) + \frac{1}{2\delta} \left[\sup_{\theta \in \Gamma(\beta)} \langle \epsilon, \theta \rangle \right]^2,$$

where $\Gamma(\beta) = \{\theta : \|\theta\|_2 \leq 1, \theta = \bar{\beta} - \beta \text{ for some } \bar{\beta} \in \mathcal{S}\}$. In particular, under (A.56),

$$(A.60) \quad D_2(\hat{\beta}, \beta^*) \leq \frac{1}{\mu} \Delta_f(\beta, \beta^*) + \frac{1}{2\mu^2} \left[\sup_{\theta \in \Gamma(\beta)} \langle \epsilon, \theta \rangle \right]^2.$$

The first two bounds reveal the important role of the Fenchel conjugate of the loss, and can be made more explicit under proper Orlicz norm conditions of ϵ ; the third conclusion, on the basis of the supremum of an empirical process [60], demonstrates how modern probabilistic tools can be used to derive finite-sample error bounds of $\hat{\beta}$ in a general noisy setup.

PROOF. First, by definition, $f(\hat{\beta}) \leq f(\beta^*)$, from which it follows that $\Delta_f(\hat{\beta}, \beta^*) \leq \langle \epsilon, \hat{\beta} - \beta^* \rangle$. Define

$$h(\delta) = \Delta_f(\delta + \beta^*, \beta^*).$$

By assumption, f is a proper function and applying Fenchel-Young's inequality gives

$$\Delta_f(\hat{\beta}, \beta^*) \leq \langle \epsilon, \delta \rangle|_{\delta=\hat{\beta}-\beta^*} \leq \frac{1}{c} \Delta_f(\delta + \beta^*, \beta^*)|_{\delta=\hat{\beta}-\beta^*} + \frac{1}{c} h^*(c\epsilon)$$

or $(1 - 1/c) \Delta_f(\hat{\beta}, \beta^*) \leq h^*(c\epsilon)/c$ for any $c > 0$.

On the other hand,

$$\begin{aligned} h^*(\zeta) &= \sup_{\delta} \langle \zeta, \delta \rangle - f(\beta^* + \delta) + f(\beta^*) + \langle \nabla f(\beta^*), \delta \rangle \\ &= \sup_{\delta} \langle \zeta + \nabla f(\beta^*), \delta \rangle - f(\beta^* + \delta) + f(\beta^*) \\ &= \sup_{\delta} \langle \zeta + \nabla f(\beta^*), \beta^* + \delta \rangle - f(\beta^* + \delta) + f(\beta^*) - \langle \zeta + \nabla f(\beta^*), \eta^* \rangle \\ &= f^*(\zeta + \nabla f(\beta^*)) + f(\beta^*) - \langle \zeta + \nabla f(\beta^*), \beta^* \rangle, \end{aligned}$$

where $f(\beta^*), \nabla f(\beta^*)$ are known to be finite. Therefore we obtain

$$(A.61) \quad \Delta_f(\hat{\beta}, \beta^*) \leq \frac{1}{c-1} [f^*((c-1)\epsilon) + f(\beta^*) - (c-1)\langle \epsilon, \beta^* \rangle], \quad \forall c > 0.$$

Taking $c = 2$ and using $\mathbb{E}\epsilon = \mathbf{0}$ gives the Δ_f risk bound (A.57).

Next, we prove the second bound under the no-model-ambiguity assumption. Using the optimality of β^* , we have further

$$h^*(\zeta) = f^*(\zeta + \nabla f(\beta^*)) - f^*(-\epsilon) - \langle \zeta, \beta^* \rangle.$$

Moreover, from the assumption and definition (A.55), it is easy to show that $\beta^* \in \partial f^*(-\epsilon)$, and so $\nabla f^*(-\epsilon) = \beta^*$, from which it follows that

$$(A.62) \quad h^*(\zeta) = \Delta_{f^*}(\zeta - \epsilon, -\epsilon).$$

Taking $c = 2$ gives (A.58) (even though $\mathbb{E}\epsilon$ may not be $\mathbf{0}$).

Finally, for any $\beta \in \mathcal{S}$, $f(\hat{\beta}) \leq f(\beta)$ and so

$$\Delta_f(\hat{\beta}, \beta^*) \leq \Delta_f(\beta, \beta^*) + \langle \epsilon, (\hat{\beta} - \beta) / \|\hat{\beta} - \beta\|_2 \rangle \|\hat{\beta} - \beta\|_2.$$

We obtain a general result

$$(A.63) \quad (\Delta_f - \delta \mathbf{D}_2)(\hat{\beta}, \beta^*) \leq \Delta_f(\beta, \beta^*) + \frac{1}{2\delta} \left[\sup_{\theta \in \Gamma(\beta)} \langle \epsilon, \theta \rangle \right]^2,$$

for any $\delta > 0$.

Based on the regularity condition,

$$\frac{\delta}{2} \mathbf{D}_2(\hat{\beta}, \beta^*) \leq \Delta_f(\hat{\beta}, \beta^*) - \frac{\delta}{2} \mathbf{D}_2(\hat{\beta}, \beta^*) \leq \Delta_f(\beta, \beta^*) + \frac{1}{\delta} \left(\sup_{\theta \in \Gamma(\beta)} \langle \epsilon, \theta \rangle \right)^2$$

for any $\delta \leq \mu$. Taking $\delta = \mu$ gives the desired result. \square

Algorithm B.1 Accelerated Bregman of the second kind

Input $\beta^{(0)}$: initial value; $\rho_{\min} > 0$, $\alpha > 0$, $M \in \mathbb{N}$, $\mu_0 \geq 0$ (e.g., $\rho_{\min} = 1$, $\alpha = 2$, $M = 3$)

- 1: $\theta_0 \in (0, 1]$, $t \leftarrow 0$, $\alpha^{(0)} \leftarrow \beta^{(0)}$;
- 2: **while** not converged **do**
- 3: $\rho_t \leftarrow \rho_{\min}/\alpha$, $s \leftarrow 0$
- 4: **repeat**
- 5: $s \leftarrow s + 1$
- 6: $\rho_t \leftarrow \alpha \rho_t$
- 7: if $t \geq 1$, then $\theta_t \leftarrow (\sqrt{r^2 + 4r} - r)/2$ with $r = (\rho_{t-1}\theta_{t-1} + \mu_0)\theta_{t-1}/\rho_t$
- 8: $\gamma^{(t)} \leftarrow (1 - \theta_t)\beta^{(t)} + \theta_t\alpha^{(t)}$
- 9: $\alpha^{(t+1)} \leftarrow \arg \min_{\beta} \{f(\beta) - \Delta_{\psi_0}(\beta, \gamma^{(t)}) + \mu_0\Delta_{\phi}(\beta, \gamma^{(t)}) + \theta_t\rho_t\Delta_{\phi}(\beta, \alpha^{(t)})\}$
- 10: $\beta^{(t+1)} \leftarrow (1 - \theta_t)\beta^{(t)} + \theta_t\alpha^{(t+1)}$
- 11: $R_t \leftarrow \theta_t^2\rho_t\Delta_{\phi}(\alpha^{(t+1)}, \alpha^{(t)}) - \Delta_{\bar{\psi}_0}(\beta^{(t+1)}, \gamma^{(t)})$
 $\quad + (1 - \theta_t)\Delta_{\bar{\psi}_0}(\beta^{(t)}, \gamma^{(t)}) + \mathbf{C}_{f(\cdot) - \Delta_{\bar{\psi}_0}(\cdot, \gamma^{(t)})}(\alpha^{(t+1)}, \beta^{(t)}, \theta_t)$
- 12: **until** $R_t \geq 0$ or $s > M$
- 13: if $s > M$, pick $(\alpha^{(t+1)}, \beta^{(t+1)}, \gamma^{(t)}, \rho_t, \theta_t)$ with the largest $R_t/(\theta_t^2\rho_t)$
- 14: $t \leftarrow t + 1$
- 15: **end while**
- 16: **return** $\beta^{(t+1)}$.

APPENDIX B: ALGORITHMS FOR ACCELERATIONS

For clarity, we give an outline of the algorithms for acceleration.

Algorithm B.2 Accelerated Bregman of the first kind

Input $\beta^{(0)}$: initial value; $\rho_{\min} > 0$, $\alpha > 0$, $M \in \mathbb{N}$, $\mu_0 \geq 0$ ($\rho_{\min} = 1$, $\alpha = 2$, $M = 3$)

- 1: $\theta_0 \in (0, 1]$, $t \leftarrow 0$;
- 2: **while** not converged **do**
- 3: $\rho_t \leftarrow \rho_{\min}/\alpha$, $s \leftarrow 0$
- 4: **repeat**
- 5: $s \leftarrow s + 1$
- 6: $\rho_t \leftarrow \alpha \rho_t$
- 7: if $t \geq 1$, then $\theta_t \leftarrow (\sqrt{r^2 + 4r} - r)/2$ with $r = (\rho_{t-1}\theta_{t-1} + \mu_0)\theta_{t-1}/\rho_t$
- 8: $\gamma^{(t)} \leftarrow \beta^{(t)} + \{\rho_{t-1}\theta_t(1 - \theta_{t-1})/(\rho_{t-1}\theta_{t-1} + \mu_0)\}(\beta^{(t)} - \beta^{(t-1)})$
 \quad if $t \geq 1$ and $\beta^{(t)}$ if $t = 0$
- 9: $\beta^{(t+1)} \leftarrow \arg \min_{\beta} \{f(\beta) - \Delta_{\psi_0}(\beta, \gamma^{(t)}) + \mu_0\mathbf{D}_2(\beta, \gamma^{(t)}) + \rho_t\mathbf{D}_2(\beta, \gamma^{(t)})\}$
- 10: $R_t \leftarrow (\rho_t\mathbf{D}_2 - \Delta_{\bar{\psi}_0})(\beta^{(t+1)}, \gamma^{(t)}) + (1 - \theta_t)\Delta_{\bar{\psi}_0}(\beta^{(t)}, \gamma^{(t)})$
- 11: **until** $R_t \geq 0$ or $s > M$
- 12: if $s > M$, pick $(\beta^{(t+1)}, \gamma^{(t)}, \rho_t, \theta_t)$ with the largest associated $R_t/(\theta_t^2\rho_t)$
- 13: $t \leftarrow t + 1$
- 14: **end while**
- 15: **return** $\beta^{(t+1)}$.

APPENDIX C: EXPERIMENTS

This section performs some simulation studies to support the theoretical results.

C.1. Computational error. In this part, we use mirror descent and DC programming to solve two nonconvex problems.

• *Nonconvex mirror descent for IS divergence minimization.* In infrared astronomical satellite (IRAS) image reconstruction [13] and audio signal processing [21, 35], the Itakura-Saito (IS) divergence (or the negative cross Burg entropy), $\text{IS}(\mathbf{a}, \mathbf{b}) = \sum_i (a_i/b_i - \log(a_i/b_i) - 1)$, is popularly used to measure the discrepancy between the observed data and the reconstructed data. Given $\mathbf{X} \in \mathbb{R}_+^{n \times p}$, and $\mathbf{y} \in \mathbb{R}_+^n$, the problem can be defined by $\min f(\boldsymbol{\beta}) := \text{IS}(\mathbf{y}, \mathbf{X}\boldsymbol{\beta})$ s.t. $\boldsymbol{\beta} \in \mathbb{R}_+^p$, which is nonconvex in $\boldsymbol{\beta}$. To maintain the nonnegativity constraint in updating $\boldsymbol{\beta}$ automatically, we develop a mirror descent algorithm. Concretely, define $g(\boldsymbol{\beta}; \boldsymbol{\beta}^-) = f(\boldsymbol{\beta}) + (\rho \mathbf{D}_\varphi - \boldsymbol{\Delta}_f)(\boldsymbol{\beta}, \boldsymbol{\beta}^-)$, where $\varphi(\boldsymbol{\beta}) = \sum_j \beta_j \log \beta_j - \beta_j$. Then, minimizing $g(\boldsymbol{\beta}; \boldsymbol{\beta}^{(t)})$ with respect to $\boldsymbol{\beta}$ gives rise to a *multiplicative* rule

$$(C.1) \quad \beta_j^{(t+1)} = \beta_j^{(t)} \exp \left[-\frac{1}{\rho} \sum_i \frac{(\mathbf{X}\boldsymbol{\beta}^{(t)})_i - y_i}{(\mathbf{X}\boldsymbol{\beta}^{(t)})_i^2} X_{ij} \right].$$

From Theorem 2, the $\mathcal{O}(1/T)$ rate of convergence holds for the optimization error $\text{avg}_{0 \leq t \leq T} (2\rho \mathbf{D}_\varphi - \boldsymbol{\Delta}_f)(\boldsymbol{\beta}^{(t+1)}, \boldsymbol{\beta}^{(t)})$. To verify this, we generated a design matrix of size 1000×1000 with all elements drawn from $U(0, 1)$, and set $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \mathbf{e}$ with β_j^* chosen uniformly from the interval $(0, 5)$ and $e_i \sim \mathcal{N}(0, \sigma^2)$ with $\sigma^2 = 10$. We fixed $1/\rho = 0.01$. Figure C.1 shows how the logarithm of optimization error converges to $-\infty$ (since $\log 0 = -\infty$) for 50 different $\boldsymbol{\beta}^{(0)}$ with $\beta_j^{(0)}$ randomly chosen from $U(0, 1)$. Observe that all the error curves in the log-log plot are bounded by a line with slope -1 .

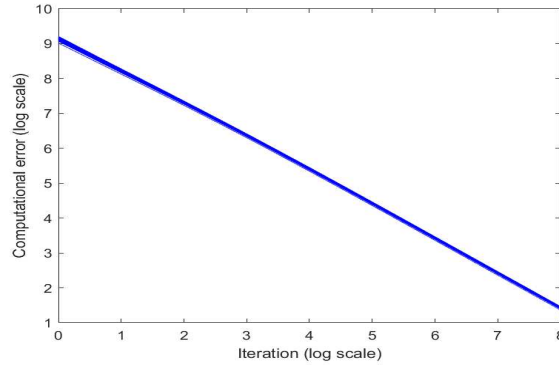


FIG C.1. Log-log plot of optimization error v.s. number of iterations: mirror descent for IS divergence minimization with 50 random starting points. All error curves are bounded above by the dashed line which has slope -1 .

• *DC programming for capped- ℓ_1 SVM.* High-dimensional classification with concurrent feature selection can be achieved by minimizing a composite objective function. Given $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \{-1, 1\}^n$, let $l(\boldsymbol{\beta}) = \sum_{i=1}^n (1 - y_i \mathbf{x}_i^\top \boldsymbol{\beta})_+$ be the hinge loss [61] that is nondifferentiable, and $P(\boldsymbol{\beta}; \lambda) = \sum_{j=1}^p \min(\lambda |\beta_j|, \lambda^2/2)$ be the capped- ℓ_1 penalty [66] which is nonsmooth and nonconvex. [42] proposed an effective DC algorithm for solving $\min_{\boldsymbol{\beta}} f(\boldsymbol{\beta}) := l(\boldsymbol{\beta}) + P(\boldsymbol{\beta}; \lambda)$ based on the decomposition $P(\boldsymbol{\beta}; \lambda) = d_1(\boldsymbol{\beta}; \lambda) - d_2(\boldsymbol{\beta}; \lambda)$ with

$$(C.2) \quad d_1(\boldsymbol{\beta}; \lambda) = \lambda \|\boldsymbol{\beta}\|_1, \quad d_2(\boldsymbol{\beta}; \lambda) = \sum_{j=1}^p \max(\lambda |\beta_j| - \lambda^2/2, 0).$$

As stated in Example 4, we can recharacterize DC as a Bregman-surrogate algorithm

$$(C.3a) \quad \beta^{(t+1)} \in \arg \min f(\beta) + \Delta_{d_2}(\beta, \beta^{(t)})$$

$$(C.3b) \quad \in \arg \min_{\beta} \sum_{i=1}^n \max(0, 1 - y_i x_i^\top \beta) + \lambda \sum_{j=1}^p (|\beta_j| - \beta_j 1_{|\beta_j^{(t)}| \geq \lambda/2}).$$

(C.3b) is equivalent to a linear program: $\min_{\xi, \zeta, \beta} \sum_{i=1}^n \xi_i + \lambda \sum_{j=1}^p \zeta_j - \lambda \sum_{j=1}^p \beta_j 1_{|\beta_j^{(t)}| \geq \lambda/2}$ s.t. $\xi_i \geq 1 - y_i x_i^\top \beta$ for $i \in [n]$, $-\zeta_j \leq \beta_j \leq \zeta_j$ for $j \in [p]$ and $\xi_i, \zeta_j \geq 0$ for $i \in [n], j \in [p]$, which can be efficiently solved by standard linear programming (LP) solvers. For the convergence of the DC algorithm, a similar result can be shown for the optimization error $\text{avg}_{0 \leq t \leq T} (\Delta_t + \tilde{\Delta}_{d_2} + \Delta_{d_1})(\beta^{(t)}, \beta^{(t+1)})$, following the lines of the proof of Proposition 4.

We generated $\mathbf{X} \in \mathbb{R}^{400 \times 800}$ with each row following $\mathcal{N}(\mathbf{0}, \Sigma)$ and $\Sigma_{ij} = 0.5^{|i-j|}$, $\beta^* = [15, 10, 0, \dots, 0]^\top$, and $\mathbf{y} = \text{sgn}(\mathbf{X}\beta^* + \mathbf{e})$ where $e_i \sim \mathcal{N}(0, \sigma^2)$ with $\sigma^2 = 10$. We fixed $\lambda = 1$ and ran the DC algorithm for 50 different starting points with each component randomly drawn from $U(0, 1)$. The corresponding optimization error curves are plotted in Figure C.2, where the $\mathcal{O}(1/T)$ rate of convergence is impressive.

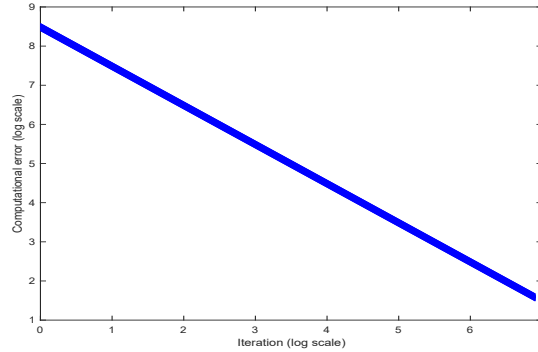


FIG C.2. Log-log plot of optimization error v.s. number of iterations: DC programming for capped- ℓ_1 SVM (50 different initial points).

C.2. Statistical error. In this part, we consider two algorithms for sparse regression: LLA and iterative thresholding. The nonconvex “hard” penalty defined by (55) is applied, which is constructed from the hard-thresholding rule via (48). Given $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^n$, we study the following regularized problem $\min_{\beta} f(\beta) := l(\beta) + \sum_{j=1}^p P_H(\varrho \beta_j; \lambda)$, where l is the loss function and $\varrho = \|\mathbf{X}\|_2$. The loss functions under consideration are the ordinary quadratic loss and a nonconvex loss which is resistant to gross outliers:

- (i) ℓ_2 loss: $l(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2/2$;
- (ii) Tukey’s biweight loss: $l(\beta) = \sum_{i=1}^n \int_0^{|x_i^\top \beta - y_i|} \psi(t) dt$ and $\psi(t) = t[1 - (t/c)^2]^2$ if $|t| \leq c$ and 0 otherwise, where $c = 4.685\sigma$ with σ a robust estimate of the standard deviation of errors [26].

In either case, we have a nonconvex optimization problem. In simulations, the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ has i.i.d. rows drawn from $\mathcal{N}(\mathbf{0}, \Sigma)$ with $\Sigma_{ij} = 0.15^{|i-j|}$, the response is given by $\mathbf{y} = \mathbf{X}\beta^* + \mathbf{e}$ with $e_i \sim \mathcal{N}(0, \sigma^2)$ and $\sigma^2 = 10$, and the regularization parameter λ is set to $A\sigma\sqrt{\log(ep)}$. We set $n = 800, p = 1000, \beta^* = [12, 8, 0, \dots, 0]^\top$ and $A = 2$.

First, we tested the statistical accuracy of LLA (cf. Theorem 3 and Proposition 5). We generated 15 initial points $\beta^{(0)}$ with each element following $U(-a, a)$, with $a \in \{0.5, 1, 1.5\}$ and 5 for each. Figure C.3 shows how the statistical error varies as the cycles progress, with each curve representing an average over 20 implementations of the same setting. Here, the errors are plotted on a log scale for a better view of the convergence rate. Unlike Figure C.1 and Figure C.2, the statistical errors can not reach 0 (or $-\infty$ in the log plot) due to the existence of noise. But they all achieved essentially the same order of statistical precision, which verifies Theorem 3, and the statistical convergence of LLA was really fast.

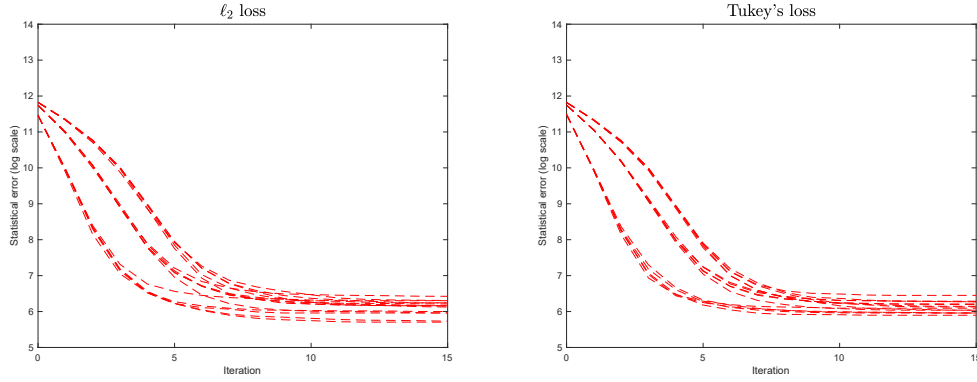


FIG C.3. Log plot of the statistical accuracy of LLA iterates in P_H -penalized sparse regression (left) and robust regression (right).

On the other hand, the computational burden of LLA turned out to be pretty high, mainly due to the cost of solving a weighted lasso problem at each iteration. We thus turned to iterative thresholding because of its low per-iteration complexity. Figure C.4 shows some analogous results. According to Figure C.4, all final statistical errors were controlled within the same order of precision. The convergence process seems to conform to the bound in Theorem 5: when t is small, $\log \Delta_\psi(\beta^*, \beta^{(t)}) \lesssim -\log(1/\kappa) t + \log(\Delta_\psi(\beta^*, \beta^{(0)}))$, and when t is large, $\log \Delta_\psi(\beta^*, \beta^{(t)}) \lesssim \kappa^t \Delta_\psi(\beta^*, \beta^{(0)}) + \log(\kappa K \lambda^2 J^* / (1 - \kappa))$, demonstrating an exponential decay.

C.3. Accelerations. We test the acceleration schemes in IS divergence minimization and robust sparse regression in this subsection.

Figure C.5 shows the power of applying the (second) acceleration in IS divergence minimization problem in Section C.1, where we used 50 starting points with $\beta_j^{(0)} \sim U(0, 1)$, $1 \leq j \leq p$. With the acceleration, the number of iterations was brought down from 1000 to less than 50 to reach the same value of the objective function, and the overall computational time was saved by nearly 90%.

Figure C.6 shows the convergence of statistical error when applying the (first) acceleration scheme in iterative thresholding for the P_H -penalized Tukey's loss minimization problem as mentioned in Section C.2. The simulation setting remains the same as before and we sampled 20 initial points with $\beta_j^{(0)} \sim U(-1, 1)$, $1 \leq j \leq p$. A substantial reduction in the number of iterations was achieved. Of course, the line search causes some overhead in computation. But the accelerated iterative thresholding still reduced the overall running time by more than 30%, and obtained slightly better statistical accuracy.

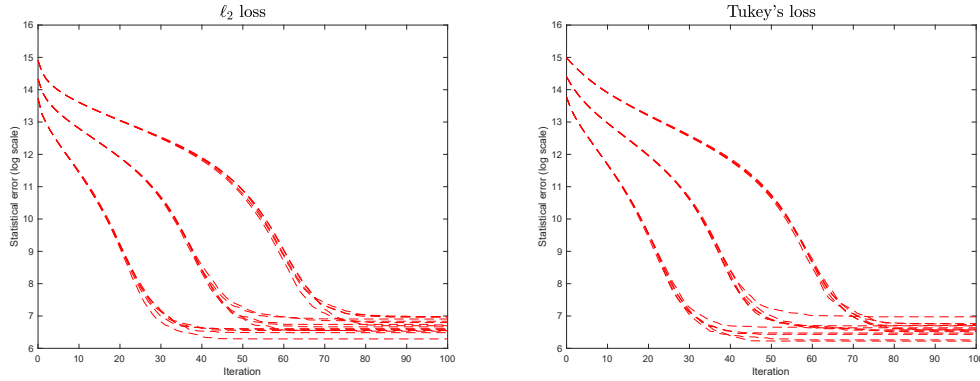


FIG C.4. Log plot of the statistical accuracy of iterative thresholding iterates in P_H -penalized sparse regression (left) and robust regression (right).

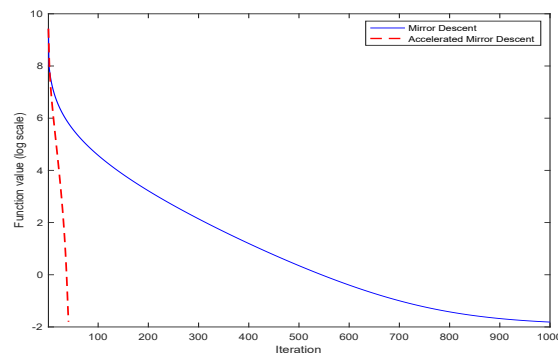


FIG C.5. Objective function value (shown on log scale) v.s. number of iterations for the plain and accelerated exponentiated gradient descent algorithms in nonconvex Burg entropy optimization.

REFERENCES

- [1] AGARWAL, A., NEGAHBAN, S. and WAINWRIGHT, M. J. (2012). Fast global convergence of gradient methods for high-dimensional statistical recovery. *The Annals of Statistics* **40** 2452–2482.
- [2] AN, L. T. H. and TAO, P. D. (2005). The DC (difference of convex functions) programming and DCA revisited with DC models of real world nonconvex optimization problems. *Annals of Operations Research* **133** 23–46.
- [3] BANERJEE, A., MERUGU, S., DHILLON, I. S. and GHOSH, J. (2005). Clustering with Bregman Divergences. *Journal of Machine Learning Research (JMLR)* **6**.
- [4] BECK, A. and TEBoulLE, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences* **2** 183–202.
- [5] BEN-TAL, A. and NEMIROVSKI, A. (2013). Optimization III: Convex Analysis, Nonlinear Programming Theory, Standard Nonlinear Programming Algorithms. Lecture Notes.
- [6] BERTSEKAS, D. P. (1999). *Nonlinear Programming*, 2nd ed. Athena Scientific.
- [7] BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *The Annals of Statistics* 1705–1732.
- [8] BLUMENSATH, T. and DAVIES, M. E. (2009). Iterative hard thresholding for compressed sensing. *Applied and Computational Harmonic Analysis* **27** 265–274.
- [9] BREGMAN, L. M. (1967). The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics* **7** 200–217.

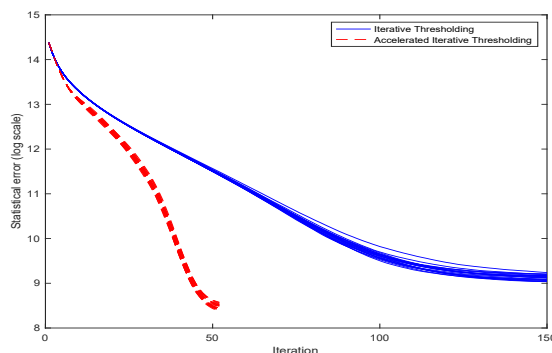


FIG C.6. Statistical error (shown on log scale) v.s. iteration number for iterative thresholding and accelerated iterative thresholding in robust sparse regression.

- [10] BUNEA, F., TSYBAKOV, A. and WEGKAMP, M. (2007). Sparsity oracle inequalities for the LASSO. *Electronic Journal of Statistics* **1** 169–194.
- [11] CANDÈS, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when p is much larger than n . *The Annals of Statistics* 2313–2351.
- [12] CANDÈS, E. J. and TAO, T. (2005). Decoding by Linear Programming. *IEEE Trans. Inf. Theor.* **51** 4203–4215.
- [13] CAO, EGGERMONT, P. P. B. and TEREBEY, S. (1999). Cross Burg entropy maximization and its application to ringing suppression in image reconstruction. *IEEE Transactions on Image Processing* **8** 286–292.
- [14] CHEN, G. and TEOULLE, M. (1993). Convergence Analysis of a Proximal-Like Minimization Algorithm Using Bregman Functions. *SIAM Journal on Optimization* **3** 538–543.
- [15] CICHOCKI, A., ICHI AMARI, S., ZDUNEK, R., KOMPASS, R., HORI, G. and HE, Z. (2006). Extended SMART algorithms for non-negative matrix factorization. In *ICAISC* (L. RUTKOWSKI, R. TADEUSIEWICZ, L. A. ZADEH and J. M. ZURADA, eds.). *Lecture Notes in Computer Science* **4029** 548–562. Springer.
- [16] CORREA, R., JOFRE, A. and THIBAUT, T. (1994). Subdifferential Monotonicity as Characterization of Convex Functions. *Numerical Functional Analysis and Optimization* **15** 531–535.
- [17] DONOHO, D. and JOHNSTONE, I. (1994). Ideal Spatial Adaptation via Wavelet Shrinkages. *Biometrika* **81** 425–455.
- [18] DUCHI, J. C., HAZAN, E. and SINGER, Y. (2011). Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *J. Mach. Learn. Res.* **12** 2121–2159.
- [19] FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* **96** 1348–1360.
- [20] FAN, J., XUE, L. and ZOU, H. (2014). Strong oracle optimality of folded concave penalized estimation. *The Annals of Statistics* **42** 819–849.
- [21] FÉVOTTE, C., BERTIN, N. and DURRIEU, J.-L. (2009). Nonnegative Matrix Factorization with the Itakura-saito Divergence: With Application to Music Analysis. *Neural Computation* **21** 793–830.
- [22] FRANK, I. E. and FRIEDMAN, J. H. (1993). A Statistical View of Some Chemometrics Regression Tools. *Technometrics* **35** 109–135.
- [23] GASSO, G., RAKOTOMAMONJY, A. and CANU, S. (2009). Recovering sparse signals with a certain family of nonconvex penalties and DC programming. *IEEE Transactions on Signal Processing* **57** 4686–4698.
- [24] GHADIMI, S. and LAN, G. (2016). Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming* **156** 59–99.
- [25] HAGOOD, J. W. and THOMSON, B. S. (2006). Recovering a function from a Dini derivative. *The American Mathematical Monthly* **113** 34–46.
- [26] HAMPEL, F. R., RONCHETTI, E. M., ROUSSEEUW, P. J. and STAHEL, W. A. (2005). *Robust Statistics*. John Wiley & Sons, New York.
- [27] HUBER, P. J. (1981). *Robust Statistics*. John Wiley and Sons, New York.
- [28] HUNTER, D. R. and LANGE, K. (2004). A tutorial on MM algorithms. *The American Statistician* 30–37.
- [29] HUNTER, D. R. and LI, R. (2005). Variable selection using MM algorithms. *The Annals of Statistics* **33** 1617–1642.

- [30] JØRGENSEN, B. (1987). Exponential Dispersion Models. *Journal of the Royal Statistical Society. Series B* **49** 127–145.
- [31] KOLTCHINSKII, V. (2011). *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems*. Springer.
- [32] KRICHENE, W., BAYEN, A. and BARTLETT, P. L. (2015). Accelerated Mirror Descent in Continuous and Discrete Time. In *Advances in Neural Information Processing Systems* (C. CORTES, N. LAWRENCE, D. LEE, M. SUGIYAMA and R. GARNETT, eds.) **28**. Curran Associates, Inc.
- [33] LANGE, K. and ZHOU, H. (2014). MM algorithms for geometric and signomial programming. *Mathematical Programming* **143** 339–356.
- [34] LEE, D. D. and SEUNG, H. S. (1999). Learning the parts of objects by nonnegative matrix factorization. *Nature* **401** 788–791.
- [35] LEFÈVRE, A., BACH, F. R. and FÉVOTTE, C. (2011). Online algorithms for nonnegative matrix factorization with the Itakura-Saito divergence. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, WASPAA 2011* 313–316.
- [36] LOH, P.-L. and WAINWRIGHT, M. J. (2015). Regularized M-estimators with nonconvexity: statistical and algorithmic theory for local optima. *Journal Machine Learning Research* **16** 559–616.
- [37] LOUNICI, K., PONTIL, M., TSYBAKOV, A. B. and VAN DE GEER, S. (2011). Oracle Inequalities and Optimal Inference under Group Sparsity. *Annals of Statistics* **39** 2164–2204.
- [38] NEMIROVSKI, A. and YUDIN, D. (1983). *Problem complexity and method efficiency in optimization*. Wiley-Interscience series in discrete mathematics. Wiley, Chichester, New York.
- [39] NESTEROV, Y. (1983). A method of solving a convex programming problem with convergence rate $O(1/k^2)$. *Soviet Mathematics Doklady* **27** 372–376.
- [40] NESTEROV, Y. (1988). On an approach to the construction of optimal methods of minimization of smooth convex functions. *Ekonom. i. Mat. Metody (In Russian)* **24** 509–517.
- [41] NESTEROV, Y. (2004). *Introductory Lectures on Convex Optimization: A Basic Course*. Applied optimization. Kluwer Academic Publ., Boston, Dordrecht, London.
- [42] ONG, C. S. and AN, L. T. H. (2013). Learning sparse classifiers with difference of convex functions algorithms. *Optimization Methods and Software* **28** 830–854.
- [43] PAN, W., SHEN, X. and LIU, B. (2013). Cluster analysis: unsupervised learning via supervised learning with a non-convex penalty. *Journal of Machine Learning Research* **14** 1865–1889.
- [44] RIGOLLET, P. and TSYBAKOV, A. (2011). Exponential Screening and optimal rates of sparse estimation. *Annals of Statistics* **39** 731–771.
- [45] ROCKAFELLAR, R. T. (1970). *Convex Analysis*. Princeton University Press, Princeton, NJ.
- [46] SCHMIDT, M. (2010). Graphical Model Structure Learning with ℓ_1 -Regularization, PhD thesis, University of British Columbia.
- [47] SHE, Y. (2009). Thresholding-based Iterative Selection Procedures for Model Selection and Shrinkage. *Electronic Journal of Statistics* **3** 384–415.
- [48] SHE, Y. (2012). An iterative algorithm for fitting nonconvex penalized generalized linear models with grouped predictors. *Computational Statistics and Data Analysis* **9** 2976–2990.
- [49] SHE, Y. (2016). On the finite-sample analysis of Θ -estimators. *Electronic Journal of Statistics* **10** 1874–1895.
- [50] SHE, Y. and CHEN, K. (2017). Robust reduced-rank regression. *Biometrika* **104** 633–647.
- [51] SHE, Y., HE, Y. and WU, D. (2014). Learning Topology and Dynamics of Large Recurrent Neural Networks. *IEEE Transactions on Signal Processing* **62** 5881–5891.
- [52] SHE, Y., LI, H., WANG, J. and WU, D. (2013). Grouped Iterative Spectrum Thresholding for Super-Resolution Sparse Spectrum Selection. *IEEE Transactions on Signal Processing* **61** 6371–6386.
- [53] SHE, Y., SHEN, J. and ZHANG, C. Supervised Multivariate Learning with Simultaneous Feature Auto-grouping and Dimension Reduction. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. To appear.
- [54] SHE, Y., WANG, Z. and SHEN, J. Gaining Outlier Resistance with Progressive Quantiles: Fast Algorithms and Theoretical Studies. *Journal of the American Statistical Association*. To appear.
- [55] TAO, P. D. and SQUAD, E. B. (1986). Algorithms for solving a class of nonconvex optimization problems. *Methods of subgradients. North-Holland Mathematics Studies* **129** 249–271.
- [56] TIBSHIRANI, R. (1996). Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society (Series B)* **58** 267–288.
- [57] TSENG, P. (2008). On accelerated proximal gradient methods for convex-concave optimization. Technical report, Department of Mathematics, University of Washington.
- [58] TSYBAKOV, A. B. (2008). *Introduction to Nonparametric Estimation*. Springer, New York, NY.
- [59] VAN DE GEER, S. A. and BÜHLMANN, P. (2009). On the conditions used to prove oracle results for the Lasso. *Electronic Journal of Statistics* **3** 1360–1392.

- [60] VAN DER VAART, A. W. and WELLNER, J. (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer.
- [61] VAPNIK, V. N. (1995). *The Nature of Statistical Learning Theory*. Springer-Verlag, New York, NY, USA.
- [62] WAINWRIGHT, M. J. and JORDAN, M. I. (2008). Graphical Models, Exponential Families, and Variational Inference. *Foundations and Trends in Machine Learning* **1** 1-305.
- [63] WANG, Z., LIU, H. and ZHANG, T. (2014). Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *The Annals of Statistics* **42** 2164–2201.
- [64] ZHANG, C., JIANG, Y. and CHAI, Y. (2010). Penalized Bregman divergence for large-dimensional regression and classification. *Biometrika* **97** 551–566.
- [65] ZHANG, C. H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics* **38** 894-942.
- [66] ZHANG, T. (2010). Analysis of Multi-stage Convex Relaxation for Sparse Regularization. *Journal of Machine Learning Research* **11** 1081–1107.
- [67] ZOU, H. and LI, R. (2008). One-step Sparse Estimates in Nonconcave Penalized Likelihood Models. *Annals of Statistics* **36** 1509–1533.