I-LAMM FOR SPARSE LEARNING: SIMULTANEOUS CONTROL OF ALGORITHMIC COMPLEXITY AND STATISTICAL ERROR¹

BY JIANQING FAN^{*,†}, HAN LIU[†], QIANG SUN^{†,‡} AND TONG ZHANG^{*,§}

Fudan University^{*}, Princeton University[†], University of Toronto[‡] and Tecent AI Lab[§]

We propose a computational framework named iterative local adaptive majorize-minimization (I-LAMM) to simultaneously control algorithmic complexity and statistical error when fitting high-dimensional models. I-LAMM is a two-stage algorithmic implementation of the local linear approximation to a family of folded concave penalized guasi-likelihood. The first stage solves a convex program with a crude precision tolerance to obtain a coarse initial estimator, which is further refined in the second stage by iteratively solving a sequence of convex programs with smaller precision tolerances. Theoretically, we establish a phase transition: the first stage has a sublinear iteration complexity, while the second stage achieves an improved linear rate of convergence. Though this framework is completely algorithmic, it provides solutions with optimal statistical performances and controlled algorithmic complexity for a large family of nonconvex optimization problems. The iteration effects on statistical errors are clearly demonstrated via a contraction property. Our theory relies on a localized version of the sparse/restricted eigenvalue condition, which allows us to analyze a large family of loss and penalty functions and provide optimality guarantees under very weak assumptions (e.g., I-LAMM requires much weaker minimal signal strength than other procedures). Thorough numerical results are provided to support the obtained theory.

1. Introduction. Modern data acquisitions routinely measure massive amounts of variables, which can be much larger than the sample size, making statistical inference an ill-posed problem. For inferential tractability and interpretability, one common approach is to exploit the penalized M-estimator

(1.1)
$$\widehat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^d}{\operatorname{argmin}} \{ \mathcal{L}(\boldsymbol{\beta}) + \mathcal{R}_{\lambda}(\boldsymbol{\beta}) \},$$

where $\mathcal{L}(\cdot)$ is a smooth loss function, $\mathcal{R}_{\lambda}(\cdot)$ is a sparsity-inducing penalty with a regularization parameter λ . Our framework encompasses the square loss, logistic

Received July 2015; revised March 2017.

¹Supported in part by NIH Grants 5R01-GM072611-11, R01-GM100474-04, R01-MH102339, R01-GM083084 and R01-HG06841, NSF Grants DMS-1206464-04, DMS-1308566, DMS-1454377, DMS-1206464, IIS-1408910, IIS-1332109 and the Science and Technology Commission of Shanghai Municipality 16JC1402600.

MSC2010 subject classifications. Primary 62J07; secondary 62C20, 62H35.

Key words and phrases. Algorithmic statistics, iteration complexity, local adaptive MM, nonconvex statistical optimization, optimal rate of convergence.

loss, Gaussian graphical model negative log-likelihood loss, Huber loss and the family of folded concave penalties [Fan and Li (2001)]. Finding optimal statistical procedures with controlled computational complexity characterizes the efforts of high-dimensional statistical learning in the last two decades. This paper makes an important leap toward this grand challenge by proposing a general algorithmic strategy for solving (1.1) even when $\mathcal{R}_{\lambda}(\boldsymbol{\beta})$ is nonconvex.

A popular choice of $\mathcal{R}_{\lambda}(\boldsymbol{\beta})$ is the Lasso penalty [Tibshirani (1996)], a convex penalty. Though a large literature exists on understanding the theory of penalized M-estimators with convex penalties [Bickel, Ritov and Tsybakov (2009), Bunea, Tsybakov and Wegkamp (2007), van de Geer and Bühlmann (2009), Negahban et al. (2012)], it has been well known [Fan and Li (2001), Zou (2006)] that the convex penalties introduce nonnegligible estimation biases. In addition, the algorithmic issues for finding a global minimizer are rarely addressed. To eliminate the estimation bias, a family of folded-concave penalties was introduced by Fan and Li (2001), which includes the smooth clipped absolute deviation (SCAD) [Fan and Li (2001)], minimax concave penalty (MCP) [Zhang (2010a)], and capped ℓ_1 -penalty [Zhang (2010b)]. Compared to their convex counterparts, these nonconvex penalties eliminate the estimation bias and attain more refined statistical rates of convergence. However, it is more challenging to analyze the theoretical properties of the resulting estimators due to nonconvexity of the penalty functions. Existing work on nonconvex penalized M-estimators treats the statistical properties and practical algorithms separately. On one hand, statistical properties are established for the hypothetical global optimum (or some local minimum), which is usually unobtainable by any practical algorithm in polynomial time. For example, Fan and Li (2001) showed that there exists a local solution that possesses an oracle property; Kim, Choi and Oh (2008) and Fan and Lv (2011) showed that the oracle estimator is a local minimizer with high probability. Later on, Kim and Kwon (2012) and Zhang and Zhang (2012) proved that the global optimum achieves the oracle property under certain conditions. Nevertheless, none of these papers specify an algorithm to find the desired solution. More recently, Agarwal, Negahban and Wainwright (2012), Loh and Wainwright (2015), Negahban et al. (2012) develop a projected gradient algorithm with desired statistical guarantees. However, they need to modify the estimating procedures to include an additional ℓ_1 -ball constraint, $\|\boldsymbol{\beta}\|_1 \leq R$, which depends on the unknown true parameter. On the other hand, practitioners have developed numerous heuristic algorithms for nonconvex optimization problems, but without theoretical guarantees. One such example is the coordinate optimization strategy studied in Breheny and Huang (2011) and Friedman et al. (2007).

So there is a gap between theory and practice: What is actually computed is not the same as what has been proved. To bridge this gap, we propose an iterative local adaptive majorize-minimization (I-LAMM) algorithm for fitting high-dimensional statistical models. Unlike most existing methods, which are mainly motivated from a statistical perspective and ignore the computational consideration, I-LAMM is both algorithmic and statistical: it computes an estimator within polynomial time and achieves optimal statistical accuracy for this estimator. In particular, I-LAMM obtains estimators with the strongest statistical guarantees for a wide family of loss functions under the weakest possible assumptions. Moreover, the statistical properties are established for the estimators computed exactly by our algorithm, which is designed to control the cost of computing resources. Compared to existing works [Agarwal, Negahban and Wainwright (2012), Loh and Wainwright (2015), Negahban et al. (2012)], our method does not impose any constraint that depends on the unknown true parameter.

Inspired by the local linear approximation to the folded concave penalty [Zou and Li (2008)], we use I-LAMM to solve a sequence of convex programs up to a prefixed optimization precision

(1.2)
$$\min_{\boldsymbol{\beta}\in\mathbb{R}^d} \left\{ \mathcal{L}(\boldsymbol{\beta}) + \mathcal{R}(\boldsymbol{\lambda}^{(\ell-1)}\odot\boldsymbol{\beta}) \right\} \quad \text{for } \ell = 1, \dots, T,$$

where $\boldsymbol{\lambda}^{(\ell-1)} = (\lambda w(|\tilde{\boldsymbol{\beta}}_1^{(\ell-1)}|), \dots, \lambda w(|\tilde{\boldsymbol{\beta}}_d^{(\ell-1)}|))^T$, $\tilde{\boldsymbol{\beta}}^{(\ell)}$ is an approximate solution to the ℓ th optimization problem in (1.2), w(·) is a weighting function, $\mathcal{R}(\cdot)$ is a decomposable convex penalty function and " \odot " denotes the Hadamard product. In this paper, we mainly consider $\mathcal{R}(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_1$, though our theory is general. The weighting function corresponds to the derivative of the folded concave penalty in Fan and Li (2001), Zou and Li (2008) and Fan and Lv (2011).

In particular, the I-LAMM algorithm obtains a crude initial estimator $\tilde{\beta}^{(1)}$ and further solves the optimization problem (1.2) for $\ell \geq 2$ with established algorithmic and statistical properties. This provides theoretical insights on how fast the algorithm converges and how much computation is needed, as well as the desired statistical properties of the obtained estimator. The whole procedure consists of T convex programs, each only needs to be solved approximately to control the computational cost. Under mild conditions, we show that only $\log(\lambda \sqrt{n})$ steps are needed to obtain the optimal statistical rate of convergence. Even though I-LAMM solves *approximately* a sequence of convex programs, the solution enjoys the same optimal statistical property of the unobtainable global optimum for the folded-concave penalized regression. The adaptive stopping rule for solving each convex program in (1.2) allows us to control both computational costs and statistical errors. Figure 1 provides a geometric illustration of the I-LAMM procedure. It contains a contraction stage and a tightening stage as described below.

- * Contraction Stage: In this stage ($\ell = 1$), we approximately solve a convex optimization problem (1.2), starting from *any* initial value $\tilde{\beta}^{(0)}$, and terminate the algorithm as long as the approximate solution enters a desired contraction region, which will be characterized in Section 2.3. The obtained estimator is called the contraction estimator, which is very crude and only serves as initialization.
- * Tightening Stage: This stage involves multiple tightening steps ($\ell \ge 2$). Specifically, we iteratively tighten the contraction estimator by solving a sequence



FIG. 1. Geometric illustration of the contraction property. The contraction stage produces an initial estimator, starting from any initial value $\tilde{\beta}^{(0)}$ that falls in the contraction region, which secures the tightening stage to enjoy optimal statistical and computational rates of convergence. The tightening stage adaptively refines the contraction estimator until it enters the optimal region, which is stated in (1.3). Here, λ is a regularization parameter, s the number of nonzero coefficients in β^* and n the sample size.

of convex programs. Each step contracts its initial estimator toward the true parameter until it reaches the optimal region of convergence. At that region, further iteration does not improve statistical performance. See Figure 1. More precisely, we will show the following contraction property:

(1.3)
$$\|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \lesssim \sqrt{\frac{s}{n}} + \delta \cdot \|\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \boldsymbol{\beta}^*\|_2 \quad \text{for } \ell \ge 2.$$

where $\boldsymbol{\beta}^*$ is the true regression coefficient, $\delta \in (0, 1)$ a prefixed contraction parameter and $\sqrt{s/n}$ the order of statistical error. Tightening helps improve the accuracy only when $\|\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \boldsymbol{\beta}^*\|_2$ dominates the statistical error. The iteration effect is clearly demonstrated. Since $\widetilde{\boldsymbol{\beta}}^{(\ell)}$ is only used to create an adaptive weight for $\widetilde{\boldsymbol{\beta}}^{(\ell+1)}$, we can control the iteration complexity by solving each subproblem in (1.2) approximately. What differs from the contraction stage is that the initial estimators in the tightening stage are already in the contraction region, making the optimization algorithm enjoy geometric rate of convergence. This allows us to rapidly solve (1.2) with small optimization error.

* (Phase Transition in Algorithmic Convergence) In the contraction stage $(\ell = 1)$, the optimization problem is not strongly convex and, therefore, our algorithm has only a sublinear convergence rate. Once the solution enters the contraction region, we will show that the feasible solutions are sparse and the objective function is essentially "low" dimensional and becomes (restricted) strongly convex and smooth in that region. Therefore, our algorithm has a linear convergence rate for $\ell > 1$. Indeed, this holds even for $\ell = 1$, which admits a sublinear rate until it enters into the contraction region and enjoys a linear rate of convergence after that; see Figure 2. But this estimator (for $\ell = 1$) is the estimator that corresponds to the LASSO penalty, not the folded concave penalty that we are looking for.



Computational Rate for Constant Correlation Design

FIG. 2. Computational rate of convergence in each stage for the simulation experiment specified in case 2 in Example 6.1. The x-axis is the iteration count k within the ℓ th subproblem. The phase transition from sublinear rate to liner rate of algorithmic convergence is clearly seen once the iterations enter the contraction region. Here, $\hat{\beta}^{(\ell)}$ is the global minimizer of the ℓ th optimization problem in (1.2) and $\beta^{(\ell,k)}$ is its kth iteration (see Figure 3). For $\ell = 1$, the initial estimation sequence has sublinear rate and once the solution sequence enters the contraction region, it becomes linear convergent. For $\ell \geq 2$, the algorithm achieves linear rate, since all estimators $\beta^{(\ell,k-1)}$ are in the contraction region.

This paper makes four major contributions. First, I-LAMM offers an algorithmic approach to obtain the optimal estimator with controlled computing resources. Second, compared to the existing literature, our method requires weaker conditions due to a novel localized analysis of sparse learning problems. Specifically, our method does not need the extra ball constraint as in Loh and Wainwright (2015) and Wang, Liu and Zhang (2014), which is an artifact of their proofs. Third, our computational framework takes the approximate optimization error into analysis and provides theoretical guarantees for the estimator that is computed by the algorithm. Fourth, our method provides new theoretical insights about the adaptive Lasso and folded-concave penalized regression. In particular, we bridge these two methodologies together using a unified framework. See Section 3.2 for more details.

The rest of this paper proceeds as follows. In Section 2, we introduce I-LAMM and its implementation. Section 3 contributes to new insights into existing meth-

ods for high-dimensional regression. In Section 4, we introduce both the localized sparse eigenvalue and localized restricted eigenvalue conditions. Statistical property and computational complexity are then presented. In Section 5, we outline the key proof strategies. Numerical simulations are provided to evaluate the proposed method in Section 6. We conclude by discussions in Section 7. All the proofs are postponed to the Supplementary Material [Fan et al. (2018)].

NOTATION. For $\mathbf{u} = (u_1, u_2, ..., u_d)^T \in \mathbb{R}^d$, we define the ℓ_q -norm of \mathbf{u} by $\|\mathbf{u}\|_q = (\sum_{j=1}^d |u_j|^q)^{1/q}$, where $q \in [1, \infty)$. Let $\|\mathbf{u}\|_{\min} = \min\{u_j : 1 \le j \le d\}$. For a set S, let |S| denote its cardinality. We define the ℓ_0 -pseudo norm of \mathbf{u} as $\|\mathbf{u}\|_0 = |\operatorname{supp}(\mathbf{u})|$, where $\operatorname{supp}(\mathbf{u}) = \{j : u_j \ne 0\}$. For an index set $\mathcal{I} \subseteq \{1, ..., d\}$, $\mathbf{u}_{\mathcal{I}} \in \mathbb{R}^d$ is defined to be the vector whose *i*th entry is equal to u_i if $i \in \mathcal{I}$ and zero otherwise. Let $\mathbf{A} = [a_{i,j}] \in \mathbb{R}^{d \times d}$. For $q \ge 1$, we define $\|\mathbf{A}\|_q$ as the matrix operator q-norm of \mathbf{A} . For index sets $\mathcal{I}, \mathcal{J} \subseteq \{1, ..., d\}$, we define $\mathbf{A}_{\mathcal{I}, \mathcal{J}} \in \mathbb{R}^{d \times d}$ to be the matrix whose (i, j)th entry is equal to $a_{i,j}$ if $i \in \mathcal{I}$ and $j \in \mathcal{J}$, and zero otherwise. We use $\operatorname{sign}(x)$ to denote the sign of x: $\operatorname{sign}(x) = x/|x|$ if $x \ne 0$ and $\operatorname{sign}(x) = 0$ otherwise. For two functionals f(n, d, s) and g(n, d, s), we denote $f(n, d, s) \gtrsim g(n, d, s)$ if $f(n, d, s) \ge Cg(n, d, s)$ for a constant C; $f(n, d, s) \lesssim g(n, d, s)$ otherwise.

2. Methodology. In this paper, we assume that the loss function $\mathcal{L}(\cdot) \in \mathcal{F}_{\mathcal{L}}$, a family of general convex loss functions specified in Appendix A of the Supplementary Material Fan et al. (2018).

2.1. Local adaptive majorize-minimization. Recall that the estimators are obtained by solving a sequence of convex programs in (1.2). We require the function $w(\cdot)$ used therein to be taken from the tightening function class \mathcal{T} , defined as

(2.1)
$$\mathcal{T} = \{ \mathbf{w}(\cdot) \in \mathcal{M} : \mathbf{w}(t_1) \le \mathbf{w}(t_2) \text{ for all } t_1 \ge t_2 \ge 0, \\ 0 \le \mathbf{w}(t) \le 1 \text{ if } t \ge 0, \\ \mathbf{w}(t) = 0 \text{ if } t \le 0 \}.$$

To fix ideas, we take $\mathcal{R}_{\lambda}(\boldsymbol{\beta})$ in (1.1) to be $\sum_{j=1}^{d} p_{\lambda}(|\beta_j|)$, where $p_{\lambda}(\cdot)$ is a folded concave penalty [Fan and Li (2001)] such as the SCAD or MCP. As discussed in Fan and Li (2001), the penalized likelihood function in (1.1) is folded concave with respect to $\boldsymbol{\beta}$, making it difficult to be maximized. We propose to use the adaptive local linear approximation (adaptive LLA) to the penalty function Fan, Xue and Zou (2014), Zou and Li (2008) and approximately solve

(2.2)
$$\underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ \mathcal{L}(\boldsymbol{\beta}) + \sum_{j=1}^{d} p_{\lambda}'(|\widetilde{\boldsymbol{\beta}}_{j}^{(\ell-1)}|)|\boldsymbol{\beta}_{j}| \right\} \quad \text{for } 1 \leq \ell \leq T,$$

where $\widetilde{\beta}_{j}^{(\ell-1)}$ is the *j*th component of $\widetilde{\beta}^{(\ell-1)}$ and $\widetilde{\beta}^{(0)}$ can be an arbitrary bad initial value: $\widetilde{\beta}^{(0)} = \mathbf{0}$, for example. If we assume that $w(\cdot) \equiv \lambda^{-1} p_{\lambda}'(\cdot) \in \mathcal{T}$, such as the

SCAD or MCP, then the adaptive LLA algorithm can be regarded as a special case of our general formulation (1.2). Note that the LLA algorithm with ℓ_q -penalty (q < 1) is not covered by our algorithm since its derivative is unbounded at the origin, and thus $\lambda^{-1} p'_{\lambda}(\cdot) \notin \mathcal{T}$. The latter creates a zero-absorbing state: once a component is shrunk to zero, it will remain zero throughout the remaining iterations, as noted in Fan and Lv (2008). Of course, we can truncate the loss derivative of the loss function to resolve this issue.

We now propose a local adaptive majorize-minimization (LAMM) principal, which will be repeatedly called to practically solve the optimization problem (2.2). We first review the majorize-minimization (MM) algorithm. To minimize a general function $f(\boldsymbol{\beta})$, at a given point $\boldsymbol{\beta}^{(k)}$, MM majorizes it by $g(\boldsymbol{\beta}|\boldsymbol{\beta}^{(k)})$, which satisfies

$$g(\boldsymbol{\beta}|\boldsymbol{\beta}^{(k)}) \ge f(\boldsymbol{\beta}) \text{ and } g(\boldsymbol{\beta}^{(k)}|\boldsymbol{\beta}^{(k)}) = f(\boldsymbol{\beta}^{(k)})$$

and then compute $\boldsymbol{\beta}^{(k+1)} = \operatorname{argmin}_{\boldsymbol{\beta}} \{g(\boldsymbol{\beta} | \boldsymbol{\beta}^{(k)})\}$ [Hunter and Lange (2004), Lange, Hunter and Yang (2000)]. The objective value of such an algorithm is nonincreasing in each step, since

(2.3)
$$f(\boldsymbol{\beta}^{(k+1)}) \stackrel{\text{major.}}{\leq} g(\boldsymbol{\beta}^{(k+1)} | \boldsymbol{\beta}^{(k)}) \stackrel{\text{min.}}{\leq} g(\boldsymbol{\beta}^{(k)} | \boldsymbol{\beta}^{(k)}) \stackrel{\text{init.}}{=} f(\boldsymbol{\beta}^{(k)}).$$

An inspection of the above arguments shows that the majorization requirement is not necessary. It requires only the local property

(2.4)
$$f(\boldsymbol{\beta}^{(k+1)}) \le g(\boldsymbol{\beta}^{(k+1)}|\boldsymbol{\beta}^{(k)}) \text{ and } g(\boldsymbol{\beta}^{(k)}|\boldsymbol{\beta}^{(k)}) = f(\boldsymbol{\beta}^{(k)})$$

for the inequalities in (2.3) to hold.

Inspired by the above observation, we locally majorize (2.2) at the ℓ th step. It is similar to the iteration steps used in the (proximal) gradient method [Boyd and Vandenberghe (2004), Nesterov (2013)]. Instead of computing and storing a large Hessian matrix as in Zou and Li (2008), we majorize $\mathcal{L}(\boldsymbol{\beta})$ in (2.2) at $\boldsymbol{\beta}^{(\ell-1)}$ by an isotropic quadratic function

$$\mathcal{L}(\widetilde{\boldsymbol{\beta}}^{(\ell-1)}) + \langle \nabla \mathcal{L}(\widetilde{\boldsymbol{\beta}}^{(\ell-1)}), \boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}}^{(\ell-1)} \rangle + \frac{\phi}{2} \| \boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}}^{(\ell-1)} \|_{2}^{2},$$

where ∇ is used to denote derivative. By Taylor's expansion, it suffices to take ϕ that is no smaller than the largest eigenvalue of $\nabla^2 \mathcal{L}(\tilde{\boldsymbol{\beta}}^{(\ell-1)})$. More importantly, the isotropic form also allows a simple analytic solution to the subsequent majorized optimization problem:

(2.5)
$$\underset{\boldsymbol{\beta} \in \mathbb{R}^{d}}{\operatorname{argmin}} \left\{ \mathcal{L}(\boldsymbol{\widetilde{\beta}}^{(\ell-1)}) + \langle \nabla \mathcal{L}(\boldsymbol{\widetilde{\beta}}^{(\ell-1)}), \boldsymbol{\beta} - \boldsymbol{\widetilde{\beta}}^{(\ell-1)} \rangle + \frac{\phi}{2} \|\boldsymbol{\beta} - \boldsymbol{\widetilde{\beta}}^{(\ell-1)}\|_{2}^{2} + \sum_{j=1}^{d} p_{\lambda}'(|\boldsymbol{\widetilde{\beta}}_{j}^{(\ell-1)}|)|\boldsymbol{\beta}_{j}| \right\}.$$

820

With $\boldsymbol{\lambda}^{(\ell-1)} = (p'_{\lambda}(|\widetilde{\beta}_1^{(\ell-1)}|), \dots, p'_{\lambda}(|\widetilde{\beta}_d^{(\ell-1)}|))^{\mathrm{T}}$, it is easy to show that (2.5) is minimized at

$$\boldsymbol{\beta}^{(\ell,1)} = T_{\boldsymbol{\lambda}^{(\ell-1)},\phi}(\widetilde{\boldsymbol{\beta}}^{(\ell-1)}) \equiv S(\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \phi^{-1}\nabla \mathcal{L}(\widetilde{\boldsymbol{\beta}}^{(\ell-1)}), \phi^{-1}\boldsymbol{\lambda}^{(\ell-1)}).$$

where $S(\mathbf{x}, \boldsymbol{\lambda})$ is the soft-thresholding operator, defined by $S(\mathbf{x}, \boldsymbol{\lambda}) \equiv (\text{sign}(x_j) \cdot \max\{|x_j| - \lambda_j, 0\})$. The simplicity of this updating rule is due to the fact that (2.5) is an unconstrained optimization problem. This is not the case in Loh and Wainwright (2015) and Wang, Liu and Zhang (2014).

However, finding the value of $\phi \ge \|\nabla^2 \mathcal{L}(\widetilde{\boldsymbol{\beta}}^{(\ell-1)})\|_2$ is not an easy task in computation. To avoid storing and computing the largest eigenvalue of a big matrix, we now state the LAMM algorithm, thanks to the local requirement (2.4). The basic idea of LAMM is to start from a very small isotropic parameter ϕ_0 and then successfully inflate ϕ by a factor $\gamma_u > 1$ (say, 2). If the solution satisfies (2.4), we stop this part of the algorithm, which will make the target value nonincreasing. Since after the *k*th iteration, $\phi = \gamma_u^{k-1} \phi_0$, there always exists a *k* such that it is no larger than $\|\nabla^2 \mathcal{L}(\widetilde{\boldsymbol{\beta}}^{(\ell-1)})\|_2$. In this manner, the LAMM algorithm will find a smallest iteration to make (2.4) hold.

Specifically, our proposed LAMM algorithm to solve (2.5) at $\tilde{\beta}^{(\ell-1)}$ begins with $\phi = \phi_0$, say 10⁻⁶, iteratively increases ϕ by a factor of $\gamma_u > 1$ inside the ℓ th step of optimization, and computes

$$\boldsymbol{\beta}^{(\ell,1)} = T_{\boldsymbol{\lambda}^{(\ell-1)}, \boldsymbol{\phi}^{(\ell,k)}}(\boldsymbol{\beta}^{(\ell,0)}) \qquad \text{with } \boldsymbol{\phi}^{(\ell,k)} = \gamma_u^{k-1} \boldsymbol{\phi}_0, \, \boldsymbol{\beta}^{(\ell,0)} = \widetilde{\boldsymbol{\beta}}^{(\ell-1)}$$

until the local property (2.4) holds. In our context, LAMM stops when

$$\Psi_{\boldsymbol{\lambda}^{(\ell-1)},\boldsymbol{\phi}^{(\ell,k)}}(\boldsymbol{\beta}^{(\ell,1)},\boldsymbol{\beta}^{(\ell,0)}) \geq F(\boldsymbol{\beta}^{(\ell,1)},\boldsymbol{\lambda}^{(\ell-1)}),$$

where $F(\boldsymbol{\beta}, \boldsymbol{\lambda}^{(\ell-1)}) \equiv \mathcal{L}(\boldsymbol{\beta}) + \sum_{j=1}^{d} \lambda_j^{(\ell-1)} |\beta_j|$ and

$$\Psi_{\boldsymbol{\lambda}^{(\ell-1)},\boldsymbol{\phi}^{(\ell,k)}}(\boldsymbol{\beta},\boldsymbol{\beta}^{(\ell,0)}) \equiv \mathcal{L}(\boldsymbol{\beta}^{(\ell,0)}) + \langle \nabla \mathcal{L}(\boldsymbol{\beta}^{(\ell,0)}), \boldsymbol{\beta} - \boldsymbol{\beta}^{(\ell,0)} \rangle \\ + \frac{\boldsymbol{\phi}^{(\ell,k)}}{2} \|\boldsymbol{\beta} - \boldsymbol{\beta}^{(\ell,0)}\|_{2}^{2} + \sum_{j=1}^{d} \lambda_{j}^{(\ell-1)} |\beta_{j}|$$

Inspired by Nesterov (2013), to accelerate LAMM within the next majorizing step, we keep track of the sequence $\{\phi^{(\ell,k)}\}_{\ell,k}$ and set $\phi^{(\ell,k)} = \max\{\phi_0, \gamma_u^{-1}\phi^{(\ell,k-1)}\}$, with the convention that $\phi_{\ell,0} = \tilde{\phi}_{\ell-1}$ and $\tilde{\phi}_0 = \phi_0$, in which $\tilde{\phi}_{\ell-1}$ is the isotropic parameter corresponding to the solution $\tilde{\beta}^{(\ell-1)}$. This is summarized in Algorithm 1 with a generic initial value.

The LAMM algorithm solves only one local majorization step. It corresponds to moving one horizontal step in Figure 3. To solve (2.2), we need to use LAMM iteratively, which we shall call the iterative LAMM (I-LAMM) algorithm, and compute a sequence of solutions $\beta^{(\ell,k)}$ using the initial value $\beta^{(\ell,k-1)}$. Figure 3

Algorithm 1 The LAMM algorithm in the *k*th iteration of the ℓ th tightening subproblem.

1: Algorithm: $\{\beta^{(\ell,k)}, \phi^{(\ell,k)}\} \leftarrow \text{LAMM}(\lambda^{(\ell-1)}, \beta^{(\ell,k-1)}, \phi_0, \phi^{(\ell,k-1)})$ 2: Input: $\lambda^{(\ell-1)}, \beta^{(\ell,k-1)}, \phi_0, \phi^{(\ell,k-1)}$ 3: Initialize: $\phi^{(\ell,k)} \leftarrow \max\{\phi_0, \gamma_u^{-1}\phi^{(\ell,k-1)}\}$ 4: Repeat 5: $\beta^{(\ell,k)} \leftarrow T_{\lambda^{(\ell-1)}, \phi^{(\ell,k)}}(\beta^{(\ell,k-1)})$ 6: If $F(\beta^{(\ell,k)}, \lambda^{(\ell-1)}) > \Psi_{\lambda^{(\ell-1)}, \phi^{(\ell,k)}}(\beta^{(\ell,k)}; \beta^{(\ell,k-1)})$ then $\phi^{(\ell,k)} \leftarrow \gamma_u \phi^{(\ell,k)}$ 7: Until $F(\beta^{(\ell,k)}, \lambda^{(\ell-1)}) \leq \Psi_{\lambda^{(\ell-1)}, \phi^{(\ell,k)}}(\beta^{(\ell,k)}; \beta^{(\ell,k-1)})$ 8: Return $\{\beta^{(\ell,k)}, \phi^{(\ell,k)}\}$

depicts the schematics of our algorithm: the ℓ th row corresponds to solving the ℓ th subproblem in (2.2) approximately, beginning by computing the adaptive weight $\lambda^{(\ell-1)}$. The number of iterations needed within each row will be discussed in the sequel.

2.2. Stopping criterion. I-LAMM recognizes that the exact solutions to (2.2) can never be achieved in practice with algorithmic complexity control. Instead, in the ℓ th optimization subproblem, we compute the approximate solution, $\tilde{\boldsymbol{\beta}}^{(\ell)}$, up to an optimization error ε , the choice of which will be discussed in next subsection. To calculate this approximate solution, starting from the initial value $\boldsymbol{\beta}^{(\ell,0)} = \tilde{\boldsymbol{\beta}}^{(\ell-1)}$, the algorithm constructs a solution sequence $\{\boldsymbol{\beta}^{(\ell,k)}\}_{k=1,2,...}$ using the introduced LAMM algorithm; see Figure 3.

We then introduce a stopping criterion for the I-LAMM algorithm. From optimization theory [Section 5.5 in Boyd and Vandenberghe (2004)], we know that any exact solution $\hat{\beta}^{(\ell)}$ to the ℓ th subproblem in (2.2) satisfies the first-order optimality condition

(2.6)
$$\nabla \mathcal{L}(\widehat{\boldsymbol{\beta}}^{(\ell)}) + \boldsymbol{\lambda}^{(\ell-1)} \odot \boldsymbol{\xi} = \boldsymbol{0}$$
 for some $\boldsymbol{\xi} \in \partial \|\widehat{\boldsymbol{\beta}}^{(\ell)}\|_1 \in [-1, 1]^d$,

where ∂ is used to indicate the subgradient operator. The set of subgradients of a function $f : \mathbb{R}^d \to \mathbb{R}$ at a point x_0 , denoted as $\partial f(x_0)$, is defined as the collection

$$\boldsymbol{\lambda}^{(0)}: \quad \boldsymbol{\beta}^{(1,0)} = \mathbf{0} \stackrel{\text{LAMM}}{\longrightarrow} \boldsymbol{\beta}^{(1,1)} \stackrel{\text{LAMM}}{\longrightarrow} \cdots \stackrel{\text{LAMM}}{\longrightarrow} \boldsymbol{\beta}^{(1,k_1)} = \widetilde{\boldsymbol{\beta}}^{(1)}, \ k_1 \lesssim \varepsilon_c^{-2};$$
$$\boldsymbol{\lambda}^{(1)}: \quad \boldsymbol{\beta}^{(2,0)} = \widetilde{\boldsymbol{\beta}}^{(1)} \stackrel{\text{LAMM}}{\longrightarrow} \boldsymbol{\beta}^{(2,1)} \stackrel{\text{LAMM}}{\longrightarrow} \cdots \stackrel{\text{LAMM}}{\longrightarrow} \boldsymbol{\beta}^{(2,k_2)} = \widetilde{\boldsymbol{\beta}}^{(2)}, \ k_2 \lesssim \log(\varepsilon_t^{-1});$$
$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots \\ \boldsymbol{\lambda}^{(T-1)}: \ \boldsymbol{\beta}^{(T,0)} = \widetilde{\boldsymbol{\beta}}^{(T-1)} \stackrel{\text{LAMM}}{\longrightarrow} \boldsymbol{\beta}^{(T,1)} \stackrel{\text{LAMM}}{\longrightarrow} \cdots \stackrel{\text{LAMM}}{\longrightarrow} \boldsymbol{\beta}^{(T,k_T)} = \widetilde{\boldsymbol{\beta}}^{(T)}, \ k_T \lesssim \log(\varepsilon_t^{-1}).$$

FIG. 3. Paradigm illustration of I-LAMM. k_{ℓ} , $1 \le \ell \le T$, is the iteration index for the ℓ th optimization in (2.2). ε_c and ε_t are the precision parameters for the contraction and tightening stage respectively and will be described in Section 2.3 in detail.

Algorithm 2 I-LAMM algorithm for each subproblem in (2.2)

1: Algorithm: $\{\widetilde{\boldsymbol{\beta}}^{(\ell)}\} \leftarrow \text{I-LAMM}(\boldsymbol{\lambda}^{(\ell-1)}, \boldsymbol{\beta}^{(\ell,0)})$ 2: Input: $\phi_0 > 0$ 3: for k = 0, 1, ... until $\omega_{\boldsymbol{\lambda}^{(\ell-1)}}(\boldsymbol{\beta}^{(\ell,k)}) \leq \varepsilon$ do 4: $\{\boldsymbol{\beta}^{(\ell,k)}, \boldsymbol{\phi}^{(\ell,k)}\} \leftarrow \text{LAMM}(\boldsymbol{\lambda}^{(\ell-1)}, \boldsymbol{\beta}^{(\ell,k-1)}, \phi_0)$ 5: end for 6: Output: $\widetilde{\boldsymbol{\beta}}^{(\ell)} = \boldsymbol{\beta}^{(\ell,k)}$

of vectors, $\boldsymbol{\xi}$, such that $f(x) - f(x_0) \ge \boldsymbol{\xi}^{\mathrm{T}}(x - x_0)$, for any x. Thus, a natural measure for suboptimality of $\boldsymbol{\beta}$ can be defined as

$$\omega_{\boldsymbol{\lambda}^{(\ell-1)}}(\boldsymbol{\beta}) = \min_{\boldsymbol{\xi} \in \partial \|\boldsymbol{\beta}\|_1} \{ \| \nabla \mathcal{L}(\boldsymbol{\beta}) + \boldsymbol{\lambda} \odot \boldsymbol{\xi} \|_{\infty} \}.$$

For a prefixed optimization error ε , we stop the algorithm within the ℓ th subproblem when $\omega_{\lambda^{(\ell-1)}}(\boldsymbol{\beta}^{(\ell,k)}) \leq \varepsilon$. We call $\tilde{\boldsymbol{\beta}}^{(\ell)} \equiv \boldsymbol{\beta}^{(\ell,k)}$ an ε -optimal solution. More details can be found in Algorithm 2.

REMARK 2.1. The I-LAMM algorithm is an early-stop variant of the ISTA algorithm to handle general loss functions and nonconvex penalties [Beck and Teboulle (2009)]. The LAMM principal serves as a novel perspective for the proximal gradient method.

2.3. *Tightening after contraction*. From the computational perspective, optimization in (2.2) can be categorized into two stages: contraction $(\ell = 1)$ and tightening $(2 \le \ell \le T)$. In the contraction stage, we start from an arbitrary initial value, which can be quite remote from the underlying true parameter. We take ε as $\varepsilon_c \simeq \lambda$, reflecting the precision needed to bring the initial solution to a contracting neighborhood of the global minimum. For instance, in linear model with sub-Gaussian errors, ε_c can be taken in the order of $\sqrt{\log d/n}$. This stage aims to find a good initial estimator $\tilde{\beta}^{(1)}$ for the subsequent optimization subproblems in the tightening stage. Recall that $s = \|\boldsymbol{\beta}^*\|_0$ is the sparsity level. We will show in Section 4.3 that with a properly chosen λ , the approximate solution $\tilde{\beta}^{(1)}$, produced by the early stopped I-LAMM algorithm, falls in the region of such good initials estimators:

$$\{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \le C\lambda\sqrt{s} \text{ and } \boldsymbol{\beta} \text{ is sparse}\}.$$

We call this region the contraction region.

However, the estimator $\tilde{\beta}^{(1)}$ suffers from a suboptimal statistical rate of convergence, which is inferior to the refined one obtained by nonconvex regularization. A second stage to tighten this coarse contraction estimator into the optimal region of convergence is needed. This is achieved by the subsequent optimization ($\ell \geq 2$) and referred to as a tightening stage. Because the initial estimators are already good

and sparse at each iteration of the tightening stage, the I-LAMM algorithm at this stage enjoys a geometric rate of convergence, due to the sparse strong convexity. Therefore, the optimization error $\varepsilon = \varepsilon_t$ can be much smaller to simultaneously ensure statistical accuracy and control computational complexity. To achieve the oracle rate $\sqrt{s/n}$: ε_t must be no larger than the order of $\sqrt{1/n}$. A graphical illustration of the full algorithm is presented in Figure 3. Theoretical justifications are provided in Section 4. From this perspective, we shall also call the psuedo-algorithm in (1.2) or (2.2), combined with LAMM, the tightening after contraction (TAC) algorithm.

3. New insights into existing methods.

3.1. Connection to one-step local linear approximation. In the low-dimensional regime, Zou and Li (2008) shows that the one-step LLA algorithm produces an oracle estimator if the maximum likelihood estimator (MLE) is used for initialization. They thus claim that the multi-step LLA is unnecessary. However, this is not the case in high dimensions, under which an unbiased initial estimator, such as the MLE, is not available. In this paper, we show that starting from a possibly arbitrary bad initial value (such as 0), the contraction stage can produce a sparse coarse estimator. Each tightening step then refines the estimator from previous step to the optimal region of convergence by

(3.1)
$$\|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \lesssim \sqrt{\frac{s}{n}} + \delta \cdot \|\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \boldsymbol{\beta}^*\|_2 \quad \text{for } 2 \le \ell \le T,$$

where $\delta \in (0, 1)$ is a prefixed contraction parameter. Unlike the one-step method in Fan, Xue and Zou (2014), the role of iteration is clearly evidenced in (3.1).

An important aspect of our algorithm (2.2) is that we use the solvable approximate solutions, $\tilde{\beta}^{(\ell)}$'s, rather than the exact ones, $\hat{\beta}^{(\ell)}$'s. In order to practically implement (2.2) for a general convex loss function, Zou and Li (2008) propose to locally approximate $\mathcal{L}(\beta)$ by a quadratic function

(3.2)
$$\mathcal{L}(\widehat{\boldsymbol{\beta}}^{(0)}) + \langle \nabla \mathcal{L}(\widehat{\boldsymbol{\beta}}^{(0)}), \boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^{(0)} \rangle + \frac{1}{2} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^{(0)})^{\mathrm{T}} \nabla^{2} \mathcal{L}(\widehat{\boldsymbol{\beta}}^{(0)}) (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^{(0)}),$$

where $\hat{\boldsymbol{\beta}}^{(0)}$ is a "good" initial estimator of $\boldsymbol{\beta}^*$ and $\nabla^2 \mathcal{L}(\hat{\boldsymbol{\beta}}^{(0)})$ is the Hessian evaluated at $\hat{\boldsymbol{\beta}}^{(0)}$. However, in high dimensions, evaluating the $d \times d$ Hessian is not only computationally intensive but also requires a large storage cost. In addition, the optimization problem (2.2) cannot be solved analytically with approximation (3.2). We resolve these issues by proposing the isotropic quadratic approximation; see Section 2. 3.2. New insight into folded-concave regularization and adaptive Lasso. The adaptive local linear approximation (2.2) provides new insight into folded-concave regularization and adaptive Lasso. To correct the Lasso's estimation bias, folded-concave regularization [Fan and Li (2001)] and its one-step implementation, adaptive Lasso [Fan, Xue and Zou (2014), Zou (2006), Zou and Li (2008)] have drawn much research interest due to their attractive statistical properties. For a general loss function $\mathcal{L}(\boldsymbol{\beta})$, the adaptive Lasso solves

$$\widehat{\boldsymbol{\beta}}_{\text{adapt}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ \mathcal{L}(\boldsymbol{\beta}) + \lambda \sum_{j=1}^{d} w(\beta_{\text{init},j}) |\beta_j| \right\},\$$

where $\beta_{\text{init}, j}$ is an initial estimator of β_j . We see that the adaptive Lasso is a special case of (2.2) with $\ell = 2$. Two important open questions for an adaptive Lasso are to obtain a good enough initial estimator in high dimensions and to select a suitable tuning parameter λ , which achieves the optimal statistical performance. Our solution to the first question is to use, the approximate solution to Lasso with controlled computational complexity, which corresponds to $\ell = 1$ in (2.2). For the choice of λ , Bühlmann and van de Geer (2011) suggested sequential tuning: in the first stage, they use cross validation to select the initial tuning parameter, denoted here by $\hat{\lambda}_{\text{init, cv}}$ and the corresponding estimator $\hat{\beta}_{\text{init}}$; in the second stage, they again adopt cross validation to select the adaptive tuning parameter λ in the adaptive Lasso. Despite the popularity of such tuning procedure, there are no theoretical guarantees to support it. As will be shown later in Theorem 4.2 and Corollary 4.3, our framework produces optimal solution by only tuning $\lambda^{(0)} = \lambda \mathbf{1}$ in the contraction stage, indicating that sequential tuning may not be necessary for the adaptive Lasso if $w(\cdot)$ is chosen from the tightening function class \mathcal{T} .

It is worth noting that a classical weight $w(\beta_j) \equiv 1/|\beta_j|$ for the adaptive Lasso does not belong to the tightening function class \mathcal{T} . As pointed out by Fan and Lv (2008), zero is an absorbing state of the adaptive Lasso with this choice of weight function. Hence, when the Lasso estimator in the first stage misses any true positives, it will be missed forever in later stages as well. In contrast, the proposed tightening function class \mathcal{T} overcomes such shortcomings by restricting the weight function $w(\cdot)$ to be bounded. This phenomenon is further elaborated via our numerical experiments in Section 6. The mean square error for the adaptive Lasso can be even worse than the Lasso estimator because the adaptive Lasso may miss true positives in the strongly correlated design case.

Our framework also reveals interesting connections between the adaptive Lasso and folded-concave regularization. Specifically, consider the following foldedconcave penalized regression:

(3.3)
$$\min_{\boldsymbol{\beta} \in \mathbb{R}^d} \{ \mathcal{L}(\boldsymbol{\beta}) + \mathcal{R}_{\lambda}(|\boldsymbol{\beta}|) \} \text{ where } \mathcal{R}_{\lambda}(|\boldsymbol{\beta}|) \text{ is a folded concave penalty.}$$

We assume that $\mathcal{R}_{\lambda}(\cdot)$ is element-wisely decomposable, that is, $\mathcal{R}_{\lambda}(|\boldsymbol{\beta}|) = \sum_{k=1}^{d} p_{\lambda}(|\beta_{k}|)$. Under this assumption, using the concave duality, we can rewrite

 $\mathcal{R}_{\lambda}(|\boldsymbol{\beta}|)$ as

(3.4)
$$\mathcal{R}_{\lambda}(|\boldsymbol{\beta}|) = \inf_{\mathbf{v}} \{|\boldsymbol{\beta}|^{\mathrm{T}} \mathbf{v} - \mathcal{R}_{\lambda}^{\star}(\mathbf{v})\},$$

where $\mathcal{R}^{\star}_{\lambda}(\cdot)$ is the dual of $\mathcal{R}_{\lambda}(\cdot)$. By the duality theory, we know that the minimum of (3.4) is achieved at $\widehat{\mathbf{v}} = \nabla \mathcal{R}_{\lambda}(|\boldsymbol{\mu}|)|_{\boldsymbol{\mu}=\boldsymbol{\beta}}$. We can employ (3.4) to reformulate (3.3) as

$$(\widehat{\boldsymbol{\beta}}, \widehat{\mathbf{v}}) = \underset{\boldsymbol{\beta}, \mathbf{v}}{\operatorname{argmin}} \{ \mathcal{L}(\boldsymbol{\beta}) + \mathbf{v}^{\mathrm{T}} |\boldsymbol{\beta}| - \mathcal{R}_{\lambda}^{\star}(\mathbf{v}) \}.$$

The optimization above can then be solved by exploiting the alternating minimization scheme. In particular, we repeatedly apply the following two steps:

- (1) Optimize over $\boldsymbol{\beta}$ with \mathbf{v} fixed: $\hat{\boldsymbol{\beta}}^{(\ell)} = \operatorname{argmin}_{\boldsymbol{\beta}} \{ \mathcal{L}(\boldsymbol{\beta}) + (\hat{\mathbf{v}}^{(\ell-1)})^{\mathrm{T}} |\boldsymbol{\beta}| \}.$
- (2) Optimize over **v** with $\boldsymbol{\beta}$ fixed. We can obtain closed form solution: $\mathbf{v}^{(\ell)} = \nabla \mathcal{R}_{\lambda}(|\boldsymbol{\mu}|)|_{\boldsymbol{\mu}=\hat{\boldsymbol{\beta}}^{(\ell)}}$.

This is a special case of (1.2) if we take $w(\boldsymbol{\beta}) = \lambda^{-1} \nabla \mathcal{R}_{\lambda}(|\boldsymbol{\mu}|)|_{\boldsymbol{\mu}=\boldsymbol{\beta}}$ and let ℓ grow until convergence. Therefore, with a properly chosen weight function $w(\cdot)$, our proposed algorithm bridges the adaptive Lasso and folded-concave penalized regression together under different choices of ℓ . In Corollary 4.3, we will prove that, when ℓ is in the order of $\log(\lambda\sqrt{n})$, then the proposed estimator enjoys the optimal statistical rate $\|\widehat{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \propto \sqrt{s/n}$, under mild conditions.

4. Theoretical results. We establish the optimal statistical rate of convergence and the computational complexity of the proposed algorithm. To establish these results in a general framework, we first introduce the localized versions of the sparse eigenvalue and restricted eigenvalue conditions.

4.1. Localized eigenvalues and assumptions. The sparse eigenvalue condition [Zhang and Zhang (2012)] is commonly used in the analysis of sparse learning problems. However, it is only valid for the least square loss. For a general loss function, the Hessian matrix depends on the parameter β and can become nearly singular in certain regions. For example, the Hessian matrix of the logistic loss is

$$\nabla^{2} \mathcal{L}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i} \mathbf{x}_{i}^{\mathrm{T}} \cdot \frac{1}{1 + \exp\left(-\mathbf{x}_{i}^{\mathrm{T}} \boldsymbol{\beta}\right)} \cdot \frac{1}{1 + \exp\left(\mathbf{x}_{i}^{\mathrm{T}} \boldsymbol{\beta}\right)},$$

which tends to zero as $\|\boldsymbol{\beta}\|_2 \to \infty$, no matter what the data are. One of our key theoretical observations is that: what we really need are the localized conditions around the true parameters $\boldsymbol{\beta}^*$, which we now introduce.

826

4.1.1. Localized sparse eigenvalue.

DEFINITION 4.1 (Localized sparse eigenvalue, LSE). The localized sparse eigenvalues are defined as

$$\rho_{+}(m,r) = \sup_{\mathbf{u},\boldsymbol{\beta}} \{ \mathbf{u}_{J}^{\mathrm{T}} \nabla^{2} \mathcal{L}(\boldsymbol{\beta}) \mathbf{u}_{J} : \|\mathbf{u}_{J}\|_{2}^{2} = 1, |J| \leq m, \|\boldsymbol{\beta} - \boldsymbol{\beta}^{*}\|_{2} \leq r \};$$

$$\rho_{-}(m,r) = \inf_{\mathbf{u},\boldsymbol{\beta}} \{ \mathbf{u}_{J}^{\mathrm{T}} \nabla^{2} \mathcal{L}(\boldsymbol{\beta}) \mathbf{u}_{J} : \|\mathbf{u}_{J}\|_{2}^{2} = 1, |J| \leq m, \|\boldsymbol{\beta} - \boldsymbol{\beta}^{*}\|_{2} \leq r \}.$$

Both $\rho_+(m, r)$ and $\rho_-(m, r)$ depend on the Hessian matrix $\nabla^2 \mathcal{L}(\boldsymbol{\beta})$, the true coefficient $\boldsymbol{\beta}^*$, the sparsity level *m*, and an extra locality parameter *r*. They reduce to the commonly-used sparse eigenvalues when $\nabla^2 \mathcal{L}(\boldsymbol{\beta})$ does not change with $\boldsymbol{\beta}$ as in the quadratic loss. The following assumption specifies the LSE condition in detail. Recall that $s = \|\boldsymbol{\beta}^*\|_0$.

ASSUMPTION 4.1. We say the LSE condition holds if there exist an integer $\tilde{s} \ge cs$ for some constant *c*, *r* and a constant *C* such that

$$0 < \rho_* \le \rho_-(2s + 2\tilde{s}, r) < \rho_+(2s + 2\tilde{s}, r) \le \rho^* < +\infty \quad \text{and} \\ \rho_+(\tilde{s}, r)/\rho_-(2s + 2\tilde{s}, r) \le 1 + C\tilde{s}/s.$$

Assumption 4.1 is standard for linear regression problems and is commonly referred to as the sparse eigenvalue condition when $r = \infty$. Such conditions have been employed by Bickel, Ritov and Tsybakov (2009), Loh and Wainwright (2015), Negahban et al. (2012), Raskutti, Wainwright and Yu (2010), Wang, Liu and Zhang (2014). The newly proposd LSE condition, to the best of our knowledge, is the weakest one in the literature.

4.1.2. Localized restricted eigenvalue. In this section, we introduce the localized version of the restricted eigenvalue condition [Bickel, Ritov and Tsybakov (2009)]. This is an alternative condition to Assumption 4.1 that allows us to handle general Hessian matrices that depend on β , under which the theoretical properties can be carried out parallelly.

DEFINITION 4.2 (Localized restricted eigenvalue, LRE). The localized restricted eigenvalue is defined as

$$\kappa_{+}(m,\gamma,r) = \sup_{\mathbf{u},\boldsymbol{\beta}} \{ \mathbf{u}^{\mathrm{T}} \nabla^{2} \mathcal{L}(\boldsymbol{\beta}) \mathbf{u} : (\boldsymbol{u},\boldsymbol{\beta}) \in \mathcal{C}(m,\gamma,r) \};$$

$$\kappa_{-}(m,\gamma,r) = \inf_{\mathbf{u},\boldsymbol{\beta}} \{ \mathbf{u}^{\mathrm{T}} \nabla^{2} \mathcal{L}(\boldsymbol{\beta}) \mathbf{u} : (\boldsymbol{u},\boldsymbol{\beta}) \in \mathcal{C}(m,\gamma,r) \},$$

where $C(m, \gamma, r) \equiv {\mathbf{u}, \boldsymbol{\beta} : S \subseteq J, |J| \le m, \|\mathbf{u}_{J^c}\|_1 \le \gamma \|\mathbf{u}_J\|_1, \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \le r}$ is a local ℓ_1 cone.

Similarly, the localized restricted eigenvalue reduces to the restricted eigenvalue when $\nabla^2 \mathcal{L}(\boldsymbol{\beta})$ does not depend on $\boldsymbol{\beta}$. We say the localized restricted eigenvalue condition holds if there exists m, γ, r such that $0 < \kappa_-(m, \gamma, r) \le \kappa_+(m, \gamma, r) < \infty$. In Appendix B, we give a geometric explanation of the local ℓ_1 cone, $\mathcal{C}(m, \gamma, r)$, and the corresponding localized analysis.

4.2. Statistical theory. In this section, we provide theoretical analysis of the proposed estimator under the LSE condition. For completeness, in Appendix B, we also establish similar results under localized restricted eigenvalue condition. We begin with the contraction stage. Recall that the initial value $\tilde{\beta}^{(0)}$ is taken as **0** for simplicity. We need the following assumption on the tightening function.

ASSUMPTION 4.2. Assume that $w(\cdot) \in \mathcal{T}$ and $w(u) \ge 1/2$ for $u = 18\rho_*^{-1}\delta^{-1}\lambda$. Here \mathcal{T} is the tightening function class defined in (2.1).

Our first result characterizes the statistical convergence rate of the estimator in the contraction stage. The key ideas of the proofs are outlined in Section 5. Other technical lemmas and details can be found in the Supplementary Material [Fan et al. (2018)].

PROPOSITION 4.1 (Statistical rate in the contraction stage). Suppose that Assumption 4.1 holds. If λ , ε and r satisfy

(4.1)
$$4(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{\infty} + \varepsilon) \le \lambda \le r\rho_*/(18\sqrt{s}),$$

then any ε_c -optimal solution $\widetilde{\boldsymbol{\beta}}^{(1)}$ satisfies

$$\|\widetilde{\boldsymbol{\beta}}^{(1)} - \boldsymbol{\beta}^*\|_2 \le 18\rho_*^{-1}\lambda\sqrt{s} \lesssim \lambda\sqrt{s}.$$

The result above is a deterministic statement. Its proof is omitted as it directly follows from Lemma 5.1 with $\ell = 1$ and \mathcal{E}_1 there to be *S*, the support of the true parameter β^* . The proof of Lemma 5.1 can be found in Appendix B. In Proposition 4.1, the approximation error ε_c , can be taken to be the order of $\lambda \approx \sqrt{\log d/n}$ in the sub-Gaussian noise case. The contraction stage ensures that the ℓ_2 estimation error is proportional to $\lambda \sqrt{s}$, which is identical to the optimal rate of convergence for the Lasso estimator [Bickel, Ritov and Tsybakov (2009), Zhang (2009)]. Our result can be regarded as a generalization of the usual Lasso analysis to more general losses, which satisfy the localized sparse eigenvalue condition. We are ready to present the main theorem, which demonstrates the effects of optimization error, shrinkage bias and tightening steps on the statistical rate.

THEOREM 4.2 (Optimal statistical rate). Suppose Assumptions 4.1 and 4.2 hold. If $4(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{\infty} + (\varepsilon_t \vee \varepsilon_c)) \leq \lambda \leq r/\sqrt{s}$, then any ε_t -optimal solution

$$\widetilde{\boldsymbol{\beta}}^{(\ell)}, \ell \geq 2, \text{ satisfies the following } \delta \text{-contraction property:} \\ \|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \leq C(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_S\|_2 + \varepsilon_t \sqrt{s} + \lambda \|\mathbf{w}(|\boldsymbol{\beta}_S^*| - u)\|_2) \\ + \delta \|\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \boldsymbol{\beta}^*\|_2,$$

where C is a constant and $u = 18\rho_*^{-1}\delta^{-1}\lambda$. Consequently, there exists a constant C' such that

$$\|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \le C'(\underbrace{\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_S\|_2}_{\text{oracle rate}} + \underbrace{\varepsilon_t \sqrt{s}}_{t \sqrt{s}} + \underbrace{\lambda \|w(|\boldsymbol{\beta}_S^*| - u)\|_2}_{\text{coefficient effect}}) + \underbrace{2C'\delta^{\ell-1}\lambda\sqrt{s}}_{t \sqrt{s}}.$$

The effect of the tightening stage can be clearly seen from the theorem above: each tightening step induces a δ -contraction property, which reduces the influence of the estimation error from the previous step by a δ -fraction. Therefore, in order to achieve the oracle rate $\sqrt{s/n}$, we shall carefully choose the optimization error such that $\varepsilon_t \leq \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_2/\sqrt{s}$ and make the tightening iterations ℓ large enough. As a corollary, we give the explicit statistical rate under the quadratic loss $\mathcal{L}(\boldsymbol{\beta}) = (2n)^{-1} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$. In this case, we take $\lambda \approx \sqrt{n^{-1}\log d}$ so that the scaling condition (4.1) holds with high probability. We use sub-Gaussian(0, σ^2) to denote a sub-Gaussian distribution random variable with mean 0 and variance proxy σ^2 .

COROLLARY 4.3. Let $y_i = \mathbf{x}_i^T \boldsymbol{\beta}^* + \epsilon_i$, $1 \le i \le n$, be independently and identically distributed sub-Gaussian random variables with $\epsilon_i \sim$ sub-Gaussian $(0, \sigma^2)$. The columns of \mathbf{X} are normalized such that $\max_j \|\mathbf{X}_{*j}\|_2 \le \sqrt{n}$. Assume there exists an $\gamma > 0$ such that $\|\boldsymbol{\beta}_S^*\|_{\min} \ge u + \gamma \lambda$ and $w(\gamma \lambda) = 0$. Under Assumptions 4.1 and 4.2, if $\lambda \asymp \sqrt{n^{-1}\log d}$, $\varepsilon_i \le \sqrt{1/n}$ and $T \gtrsim \log \log d$, then with probability at least $1 - 2d^{-\eta_1} - 2\exp\{-\eta_2 s\}$, $\tilde{\boldsymbol{\beta}}^{(T)}$ must satisfy $\|\tilde{\boldsymbol{\beta}}^{(T)} - \boldsymbol{\beta}^*\|_2 \lesssim \sqrt{s/n}$,

where η_1 and η_2 are positive constants.

Corollary 4.3 indicates that I-LAMM can achieve the oracle statistical rate $\sqrt{s/n}$ as if the support for the true coefficients were known in advance. To achieve such rate, we require $\varepsilon_c \leq \sqrt{\log d/n}$ and $\varepsilon_t \leq \sqrt{1/n}$. In other words, we need only a more accurate estimator in the tightening stage rather than in both stages. This will help us to relax the computational burden, which will be discussed in detail in Theorem 4.7. Our last result concerns the oracle property of the obtained estimator $\tilde{\beta}^{(\ell)}$ for ℓ large enough, with the proof postponed to Appendix B in the Supplementary Material [Fan et al. (2018)]. We first define the oracle estimator $\hat{\beta}^{\circ}$ as

$$\widehat{\boldsymbol{\beta}}^{\circ} = \operatorname*{argmin}_{\operatorname{supp}(\boldsymbol{\beta})=S} \mathcal{L}(\boldsymbol{\beta}).$$

THEOREM 4.4 (Strong oracle property). Suppose Assumptions 4.1 and 4.2 hold. Assume $\|\boldsymbol{\beta}_{S}^{*}\|_{\min} \geq u + \gamma \lambda$ and $w(\gamma \lambda) = 0$ for some constant γ . Let $4(\|\nabla \mathcal{L}(\hat{\boldsymbol{\beta}}^{\circ})\|_{\infty} + \varepsilon_{c} \vee \varepsilon_{t}) \leq \lambda \leq r/\sqrt{s}$ and $\varepsilon_{t} \leq \lambda/\sqrt{s}$. If $\|\hat{\boldsymbol{\beta}}^{\circ} - \boldsymbol{\beta}^{*}\|_{\max} \leq \eta_{n} \leq \lambda$, then for ℓ large enough such that $\ell \geq \log\{(1 + \varepsilon_{c}/\lambda)\sqrt{s}\}$, we have

$$\widetilde{\boldsymbol{\beta}}^{(\ell)} = \widehat{\boldsymbol{\beta}}^{\circ}$$

The theorem above is again a deterministic result. Large probability bound can be obtained by bounding the probability of the event $\{4(\|\nabla \mathcal{L}(\hat{\boldsymbol{\beta}}^{\circ})\|_{\infty} + (\varepsilon_c \vee \varepsilon_t)) \leq \lambda\}$. The assumption that $\|\hat{\boldsymbol{\beta}}^{\circ} - \boldsymbol{\beta}^*\|_{\max} \lesssim \lambda$ is very mild, because the oracle estimator only depends on the intrinsic dimension *s* rather than *d*. For instance, under linear model with sub-Gaussian errors, it can be shown that $\|\hat{\boldsymbol{\beta}}^{\circ} - \boldsymbol{\beta}\|_{\max} \leq \sqrt{\log s/n}$ with high probability.

Theorem 4.4 implies that the oracle estimator $\hat{\beta}^{\circ}$ is a fixed point of the I-LAMM algorithm, namely, once the initial estimator is $\hat{\beta}^{\circ}$, the next iteration produces the same estimator. This is in the same spirit as that proved in Fan, Xue and Zou (2014).

4.3. *Computational theory*. In this section, we analyze the computational rate for all of our approximate solutions. We start with the following assumption.

ASSUMPTION 4.3. $\nabla \mathcal{L}(\boldsymbol{\beta})$ is locally ρ_c -Lipschitz continuous, that is, (4.2) $\|\nabla \mathcal{L}(\boldsymbol{\beta}_1) - \nabla \mathcal{L}(\boldsymbol{\beta}_2)\|_2 \le \rho_c \|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2\|_2$ for $\boldsymbol{\beta}_1, \boldsymbol{\beta}_2 \in B_2(R/2, \boldsymbol{\beta}^*)$, where ρ_c is the Lipschitz constant and $R \le \|\boldsymbol{\beta}^*\|_2 + \lambda \sqrt{s}$.

We then give the explicit iteration complexity of the contraction stage in the following proposition. Recall the definition of ϕ_0 and γ_u in Algorithm 2.1, and ρ_* in Assumption 4.1.

PROPOSITION 4.5 (Sublinear rate in the contraction stage). Assume that Assumptions 4.1 and 4.3 hold. Let $4(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{\infty} + \varepsilon_c) \le \lambda \le r/\sqrt{s}$. To achieve an approximate local solution $\tilde{\boldsymbol{\beta}}^{(1)}$ such that $\omega_{\lambda^{(0)}}(\tilde{\boldsymbol{\beta}}^{(1)}) \le \varepsilon_c$ in the contraction stage, we need no more than $((1 + \gamma_u)R\rho_c/\varepsilon_c)^2$ LAMM iterations, where ρ_c is a constant defined in (4.2).

The sublinear rate is due to the lack of strong convexity of the loss function in the contraction stage, because we allow starting with arbitrary bad initial value, say **0**. Once it enters the contracting region (aka, the tightening stage), the problem becomes sparse strongly convex (see Proposition B.3 in Appendix B), which endows the algorithm a linear rate of convergence. This is empirically demonstrated in Figure 2. Our next proposition gives a formal statement on the geometric convergence rate for each subproblem in the tightening stage.

PROPOSITION 4.6 (Geometric rate in the tightening stage). Suppose that the same conditions for Theorem 4.2 hold. To obtain an approximate solution $\tilde{\boldsymbol{\beta}}^{(\ell)}$ satisfying $\omega_{\boldsymbol{\lambda}^{(\ell-1)}}(\tilde{\boldsymbol{\beta}}^{(\ell)}) \leq \varepsilon$ in each step of the ℓ th tightening stage ($\ell \geq 2$), we need at most $C' \log(C'' \lambda \sqrt{s}/\varepsilon)$ LAMM iterations, where C' and C'' are two positive constants.

Proposition 4.6 suggests that we only need to conduct a logarithmic number of LAMM iterations in each tightening step. Simply combining the computational rate in both the contraction and the tightening stages, we manage to obtain the global computational complexity.

THEOREM 4.7. Assume that $\lambda \sqrt{s} = o(1)$. Suppose that the same conditions for Theorem 4.2 hold. To achieve an approximate solution $\tilde{\boldsymbol{\beta}}^{(\ell)}$ such that $\omega_{\boldsymbol{\lambda}^{(0)}}(\tilde{\boldsymbol{\beta}}^{(1)}) \leq \varepsilon_c \lesssim \lambda$ and $\omega_{\boldsymbol{\lambda}^{(k-1)}}(\tilde{\boldsymbol{\beta}}^{(k)}) \leq \varepsilon_t \lesssim \sqrt{1/n}$ for $2 \leq k \leq T$, the total number of LAMM iterations we need is at most

$$C'\frac{1}{\varepsilon_c^2} + C''(T-1)\log\left(\frac{1}{\varepsilon_t}\right),$$

where C' and C'' are two positive constants, and $T \simeq \log(\lambda \sqrt{n})$.

REMARK 4.8. We complete this section with a remark on the sublinear rate in the contraction stage. Without further structures, the sublinear rate in the first stage is the best possible one for the proposed optimization procedure when λ is held fixed. Linear rate can be achieved when we start from a sufficiently good initial value. Another strategy is to use the path-following algorithm which is developed in Wang, Liu and Zhang (2014), where they gradually reduce the size of λ to ensure the solution sequence to be sparse.

5. Proof strategy for main results. In this section, we present the proof strategies for the main statistical and computational theorems, with technical lemmas and other details left in the Supplementary Material [Fan et al. (2018)].

5.1. Proof strategy for statistical recovery result in Section 4.2. Proposition 4.1 indicates that the contraction estimator suffers from a suboptimal rate of convergence $\lambda \sqrt{s}$. The tightening stage helps refine the statistical rate adaptively. To suppress the noise in the ℓ th subproblem, it is necessary to control $\min_j \{ |\tilde{\beta}_j^{(\ell-1)}| : j \in S^c \}$ in high dimensions. For this, we construct an entropy set \mathcal{E}_{ℓ} of *S* in each tightening subproblem to bound the magnitude of $\|\lambda_{\mathcal{E}_{\ell}^c}^{(\ell-1)}\|_{\min}$. The entropy set at the ℓ th step is defined as

(5.1)
$$\mathcal{E}_{\ell} = S \cup \{j : \lambda_j^{(\ell-1)} < \lambda w(u), u = 18\delta^{-1}\rho_*^{-1}\lambda \propto \lambda\}.$$

Under mild conditions, we will show that $|\mathcal{E}_{\ell}| \leq 2s$ and $\|\lambda_{\mathcal{E}_{\ell}^{\ell}}^{(\ell)}\|_{\min} \geq \lambda w(u) \geq \lambda/2$, which is more precisely stated in the following lemma.

LEMMA 5.1. Suppose that Assumptions 4.1 and 4.2 hold. If $4(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{\infty} + \varepsilon_t \vee \varepsilon_c) \leq \lambda \leq r/\sqrt{s}$, we must have $|\mathcal{E}_{\ell}| \leq 2s$, and the ε -optimal solution $\widetilde{\boldsymbol{\beta}}^{(\ell)}$ satisfies

$$\begin{aligned} \|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 &\leq 12\rho_*^{-1}(\|\boldsymbol{\lambda}_S^{(\ell-1)}\|_2 + \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_{\mathcal{E}_\ell}\|_2 + \varepsilon\sqrt{|\mathcal{E}_\ell|}) \\ &\leq 18\rho_*^{-1}\lambda\sqrt{s} \lesssim \lambda\sqrt{s}. \end{aligned}$$

Lemma 5.1 bounds $\|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2$ in terms of $\|\boldsymbol{\lambda}_S^{(\ell-1)}\|_2$, which is further upper bounded by the order of $\lambda\sqrt{s}$. The rate $\lambda\sqrt{s}$ coincides with the convergence rate of the contraction estimator. Later, we will exploit this result in our localized analysis to secure that all the approximate solutions $\{\widetilde{\boldsymbol{\beta}}^{(\ell)}\}_{\ell=1,...,T}$ fall in a local ℓ_2 -ball centered at $\boldsymbol{\beta}^*$ with radius $r \gtrsim \lambda\sqrt{s}$.

The next lemma further bounds $\|\lambda_S^{(\ell-1)}\|_2$ using functionals of $\tilde{\boldsymbol{\beta}}^{(\ell-1)}$, which connects the adaptive regularization parameter to the estimator from previous steps.

LEMMA 5.2. Assume $w \in \mathcal{T}$. Let $\lambda_j^{(\ell-1)} = \lambda w(|\widetilde{\beta}_j^{(\ell-1)}|)$ for $\widetilde{\beta}^{(\ell-1)}$, then for any norm $\|\cdot\|_*$, we have

$$\|\boldsymbol{\lambda}_{S}^{(\ell-1)}\|_{*} \leq \lambda \|\mathbf{w}(|\boldsymbol{\beta}_{S}^{*}|-u)\|_{*} + \lambda u^{-1} \|\boldsymbol{\beta}_{S}^{*} - \widetilde{\boldsymbol{\beta}}_{S}^{(\ell-1)}\|_{*},$$

where $\mathbf{w}(|\boldsymbol{\beta}_{S}^{*}|-u) \equiv (\mathbf{w}(|\boldsymbol{\beta}_{i}^{*}|-u))_{i \in S}.$

Lemma 5.2 bounds the tightening weight $\lambda^{(\ell-1)}$ in the ℓ th subproblem by two terms. The first term describes the coefficient effects: when the coefficients are large enough (in absolute value) such that $\|\boldsymbol{\beta}^*\|_{\min} \ge u + \gamma\lambda$ and $w(\gamma\lambda) = 0$, it becomes 0. The second term concerns the estimation error of the estimator from previous step. Combing the above two lemmas, we prove that $\tilde{\boldsymbol{\beta}}^{(\ell)}$ benefits from the tightening stage and possesses a refined statistical rate of convergence. The proof of Corollary 4.3 is left in Appendix B in the Supplementary Material [Fan et al. (2018)].

PROOF OF THEOREM 4.2. Applying Lemma 5.1, we obtain the size of the entropy set \mathcal{E}_{ℓ} [see the definition in (5.1)] is bounded by 2s and

(5.2)
$$\|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \le C_1(\|\boldsymbol{\lambda}_S^{(\ell-1)}\|_2 + \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_{\mathcal{E}_\ell}\|_2 + \varepsilon_t \sqrt{|\mathcal{E}_\ell|}) \lesssim \lambda \sqrt{s},$$

where $C_1 = 12\rho_*^{-1}$. Using Lemma 5.2 yields that

$$\|\boldsymbol{\lambda}_{S}^{(\ell-1)}\|_{2} \leq \lambda \|\mathbf{w}(|\boldsymbol{\beta}_{S}^{*}|-u)\|_{2} + \lambda u^{-1} \|(\widetilde{\boldsymbol{\beta}}^{(\ell-1)}-\boldsymbol{\beta}^{*})_{S}\|_{2}$$

Plugging the inequality above into (5.2) obtains that

(5.3)
$$\|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \leq C_1(\underbrace{\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_{\mathcal{E}_\ell}\|_2 + \varepsilon_t \sqrt{|\mathcal{E}_\ell|}}_{I} + \lambda \|\mathbf{w}(|\boldsymbol{\beta}^*_S| - u)\|_2) + C_1 \lambda u^{-1} \|(\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \boldsymbol{\beta}^*)_S\|_2.$$

We now simplify the inequality above by providing an upper bound for term I. Decomposing the support set \mathcal{E}_{ℓ} into *S* and $\mathcal{E}_{\ell} \setminus S$ and applying the triangle inequality along with the Hölder inequality, we have

(5.4)
$$\mathbf{I} \leq \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_S\|_2 + \varepsilon_t \sqrt{s} + (\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{\infty} + \varepsilon_t) \sqrt{\mathcal{E}_{\ell}/S}.$$

Following the proof of Lemma 5.1 in Appendix B, $\sqrt{|\mathcal{E}_{\ell} \setminus S|}$ can be bounded by

$$\|\widetilde{\boldsymbol{\beta}}_{\mathcal{E}_{\ell} \setminus S}^{(\ell-1)}\|_{2}/u \leq \|\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \boldsymbol{\beta}^{*}\|_{2}/u \qquad \text{where } u = 18\rho^{*-1}\delta^{-1}\lambda \propto \lambda.$$

Therefore, (5.4) can be simplified to

.....

$$\mathbf{I} \leq \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_S\|_2 + \varepsilon_t \sqrt{s} + \frac{\lambda}{4u} \|\widetilde{\boldsymbol{\beta}}^{(\ell-1)} - \boldsymbol{\beta}^*\|_2$$

which, combining with (5.3), yields the contraction property with δ . Consequently, we obtain

$$\begin{split} \|\widetilde{\boldsymbol{\beta}}^{(\ell)} - \boldsymbol{\beta}^*\|_2 \\ &\leq C(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_{\mathcal{E}_{\ell}}\|_2 + \varepsilon_t \sqrt{s} + \lambda \|\mathbf{w}_S(|\boldsymbol{\beta}_S^*| - u)\|_2) + \delta^{\ell-1} \|\widetilde{\boldsymbol{\beta}}^{(1)} - \boldsymbol{\beta}^*\|_2 \\ &\leq C(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)_{\mathcal{E}_{\ell}}\|_2 + \varepsilon_t \sqrt{s} + \lambda \|\mathbf{w}_S(|\boldsymbol{\beta}_S^*| - u)\|_2) + C\delta^{\ell-1}\lambda \sqrt{s}, \end{split}$$

where $C = C_1/(1 - \delta)$ and the last inequality follows from Proposition 4.1. The proof is complete. \Box

5.2. *Proof strategy for computational result in Section* 4.3. In this section, we present the sketch for the proofs of the results in Section 4.3. We start with the contraction stage. The next lemma shows that the contraction stage enjoys a sublinear rate of convergence. The proof can be found in Appendix C.

LEMMA 5.3. Recall that
$$F(\boldsymbol{\beta}, \boldsymbol{\lambda}) = \mathcal{L}(\boldsymbol{\beta}) + \sum_{j=1}^{d} \lambda_j |\beta_j|$$
. We have
 $F(\boldsymbol{\beta}^{(1,k)}, \boldsymbol{\lambda}^{(0)}) - F(\widehat{\boldsymbol{\beta}}^{(1)}, \boldsymbol{\lambda}^{(0)}) \leq \frac{\phi_c}{2k} \|\boldsymbol{\beta}^{(1,0)} - \widehat{\boldsymbol{\beta}}^{(1)}\|_2^2.$

The result above suggests that the optimization error decreases to zero at the rate of 1/k, while Proposition 4.1 indicates that the best statistical rate for the contraction stage is only in the order of $\lambda \sqrt{s}$. Therefore, one can early stop the LAMM iterations in the contraction stage as soon as it enters the contraction region

 $\{\boldsymbol{\beta} : \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq C\lambda\sqrt{s}, \boldsymbol{\beta} \text{ is sparse}\}$. It is this lemma that helps characterize the iteration complexity in terms of the total number of LAMM updates needed in the contraction stage; see Proposition 4.5.

To utilize the localized sparse eigenvalue condition in the tightening stage, we need the following proposition, which characterizes the sparsity of all the approximate solutions produced by the contraction stage.

LEMMA 5.4. Assume that Assumption 4.1 holds. If $4(\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{\infty} + \varepsilon_c) \leq \lambda \leq r/\sqrt{s}$, then $\widetilde{\boldsymbol{\beta}}^{(1)}$ in the contraction stage is $s + \widetilde{s}$ sparse. In particular, we have $\|(\widetilde{\boldsymbol{\beta}}^{(1)})_{S^c}\|_0 \leq \widetilde{s}$.

Together with Proposition 4.1, it ensures that the contraction estimator $\tilde{\boldsymbol{\beta}}^{(1)}$ falls in the contraction region { $\boldsymbol{\beta} : \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \le C\lambda\sqrt{s}$ and $\boldsymbol{\beta}$ is sparse}. This makes the localized sparse eigenvalue condition useful, and thus makes the geometric rate of convergence possible.

LEMMA 5.5 (Geometric rate in the tightening stage). Under the same conditions for Theorem 4.2, for any $\ell \ge 2$, $\{\beta^{(\ell,k)}\}$ converges geometrically,

$$F(\boldsymbol{\beta}^{(\ell,k)},\boldsymbol{\lambda}^{(\ell-1)}) - F(\widehat{\boldsymbol{\beta}}^{(\ell)},\boldsymbol{\lambda}^{(\ell-1)})$$

$$\leq \left(1 - \frac{1}{4\gamma_{\mu}\kappa}\right)^{k} \{F(\boldsymbol{\beta}^{(\ell,0)},\boldsymbol{\lambda}^{(\ell-1)}) - F(\widehat{\boldsymbol{\beta}}^{(\ell)},\boldsymbol{\lambda}^{(\ell-1)})\}$$

The above result suggests that each subproblem in the tightening stage enjoys a geometric rate of convergence, which is the fastest possible rate among all first-order optimization methods under the blackbox model. Lemma 5.5 can be used to obtain the computational complexity analysis of each single step of the tightening stage, that is, Proposition 4.6.

6. Numerical examples. In this section, we evaluate the statistical performance of the proposed framework through several numerical experiments. We consider the following three examples.

EXAMPLE 6.1 (Linear regression). In the first example, continuous responses were generated according to the model

(6.1)
$$y_i = \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\beta}^* + \epsilon_i$$
, where $\boldsymbol{\beta}^* = (5, 3, 0, 0, -2, \underbrace{0, \dots, 0}_{d-5})^{\mathrm{T}}$,

and n = 100. Moreover, in model (6.1), $\{\mathbf{x}_i\}_{i \in [n]}$ are generated from $N(0, \Sigma)$ distribution with covariance matrix Σ , which is independent of $\epsilon_i \sim N(0, 1)$. We take Σ as a correlation matrix $\Sigma = (\rho_{ij})$ as follows:

- Case 1: independent correlation design with $(\rho_{ij}) = \text{diag}(1, \dots, 1)$.
- Case 2: constant correlation design with $\rho_{ij} = 0.75$ if $i \neq j$; $\rho_{ij} = 1$, otherwise.
- Case 3: autoregressive correlation design with $\rho_{ij} = 0.95^{|i-j|}$.

EXAMPLE 6.2 (Logistic regression). In the second example, independent observations with binary responses are generated according to the model

$$\mathbb{P}(y_i = 1 | \mathbf{x}_i) = \frac{\exp\{\mathbf{x}_i^T \boldsymbol{\beta}^*\}}{1 + \exp\{\mathbf{x}_i^T \boldsymbol{\beta}^*\}}, \qquad i = 1, \dots, n$$

where β^* and $\{\mathbf{x}_i\}_{i \in [n]}$ are generated in the same manner as in the case 1 of Example 6.1.

EXAMPLE 6.3 (Varying dimensions and sample sizes). In this example, we continue Example 6.1 with varying dimensions and sample sizes. Specifically, we consider linear regression under autoregressive correlation design with $\rho_{ij} = 0.90^{|i-j|}$ with *d* varying from 1000 to 3500 and *n* varying from 100 to 500.

In the first two cases, we fix the sample size n at 100 and consider d = 1000. We investigate the sparsity recovery and estimation properties of the I-LAMM (or TAC) estimator via numerical simulations. We compared the I-LAMM estimator with the following methods: the oracle estimator which assumes the availability of the active set S; the refitted Lasso (Refit), which uses a post least square refit on the selected set from Lasso; the adaptive Lasso (ALasso) estimator with weight function w(β_i) = 1/ $|\beta_i|$ proposed by Zou (2006); the smoothly clipped absolute deviation (SCAD) estimator [Fan and Li (2001)] with a = 3.7; and the minimax concave penalty (MCP) estimator with a = 3 [Zhang (2010a)]. For I-LAMM, we used the 3-fold cross-validation to select the constant $c \in 0.5 \times \{1, 2, \dots, 20\}$ in the tuning parameter $\lambda = c \sqrt{\log d/n}$ in the contraction stage, with regularization parameters updated automatically at later steps. We further took $\gamma_u = 2$, $\varepsilon_c = \sqrt{\log d/n}$ and $\varepsilon_t = \sqrt{1/n}$. For the Lasso, we used the I-LAMM algorithm; for the ALasso, sequential tuning in Bühlmann and van de Geer (2011) was used: we employed the 3-fold cross validation in each step with the I-LAMM algorithm used; and the SCAD and MCP estimators were computed using the R package nevreg and the 3-fold cross-validation was used for tuning parameter selection.

For each simulation setting, we generated 100 simulated datasets and applied different estimators to each dataset. We report different statistics for each estimator in Table 1 and Figure 4. To measure the sparsity recovery performance, we calculated the median of the number of zero coefficients incorrectly estimated to be nonzero (i.e., false positive, denoted as FP), the median of the number of nonzero coefficients correctly estimated to be nonzero (i.e., true positive, denoted by TP). To measure the estimation accuracy, we calculated the median of mean squared error (MSE). To evaluate the computational efficiency, we gave the median of time

	MSE	TP	FP	Time	MSE	TP	FP	Time
	Linear\Case 1				Linear\Case 2			
I-LAMM	0.0285	3.00	0.00	0.17	0.0659	3.00	0.00	0.19
Lasso	0.3114	3.00	17.00	0.02	1.3709	3.00	16.00	0.04
Refit	0.5585	3.00	17.00	0.02	2.1573	3.00	16.00	0.04
ALasso	0.4616	3.00	15.00	0.06	1.6077	3.00	13.00	0.08
SCAD	0.0397	3.00	0.00	0.21	0.0695	3.00	0.00	0.23
MCP	0.0344	3.00	0.00	0.17	0.0706	3.00	0.00	0.22
Oracle	0.0258	3.00	0.00	-	0.0565	3.00	0.00	-
	Linear\Case 3				Logistic			
I-LAMM	0.2819	3.00	3.00	0.22	8.94	3.00	0.00	0.20
Lasso	5.8061	2.00	20.00	0.03	26.92	3.00	20.00	0.03
Refit	2.6354	2.00	20.00	0.03	26.85	3.00	20.00	0.03
ALasso	4.4242	2.00	12.00	0.06	8.28	3.00	7.00	0.05
SCAD	14.8680	2.00	5.00	0.25	9.48	3.00	12.00	0.21
MCP	14.9381	1.00	1.00	0.18	11.84	3.00	3.00	0.22
Oracle	0.1661	3.00	0.00	-	3.32	3.00	0.00	-

 TABLE 1

 The median of MSE, TP, FP, Time in seconds under the Case 1, Case 2 and Case 3 for linear regression in Example 6.1 and logistic regression in Example 6.2

(in seconds) used to produce the final estimator for different methods. Note that the computational time provided here is merely for a reference. They depend on optimization errors and implementation.



FIG. 4. The median of MSE with varying dimensions and sample sizes in Example 6.3.

We have several important observations. First, it is not surprising that Lasso tends to overfit. Other procedures improve the performance of Lasso by reducing the estimation bias and the false positive rate. The best overall performance is achieved by the I-LAMM estimator with small MSE and FP in all cases. The MCP and SCAD estimators also have overall good performance in the logistic regression model, and case 1 and case 2 of the linear regression model. However, all of MCP, SCAD and ALasso breaks down by missing true positives in case 3, where the design matrix exhibits a strong correlation between features, while I-LAMM remains the best followed by the Lasso estimator. This suggests the superiority of I-LAMM over other implementation-based nonconvex penalized regression methods under strongly correlated designs. The MSE of the I-LAMM estimator keeps flat when the dimension d varies, which justifies the oracle rate $\sqrt{s/n}$. SCAD and MCP have competitive performance when the dimension is relatively small, but they quickly break down when the dimension gets larger. This is possibly due to the numerical instability for directly solving nonconvex systems. This phenomenon is also observed in Wang, Liu and Zhang (2014). When the sample size is increasing, the performances of I-LAMM, SCAD and MCP are almost identical to each other while other convex methods suffer from slightly worse performance.

In addition, to demonstrate the phase transition phenomenon, in Figure 2, we plot the log estimation error verses the number of iterations for each tightening step for case 2 in Example 6.1. Indeed, the contraction stage suffers a sublinear rate of convergence before getting into the contracting region and enjoys a geometric rate afterwards, while the tightening stage has a geometric rate of convergence. These are in line with our asymptotic theory.

7. Conclusions and discussions. We propose a computational framework, I-LAMM (or TAC), for simultaneous control of algorithmic complexity and statistical error when fitting high-dimensional models. Even though I-LAMM only solves a sequence of convex programs approximately, the solution enjoys the same optimal statistical property of the unobtainable global optimum for the folded-concave penalized regression. Our theoretical treatment relies on a novel localized analysis which avoids the parameter bound contraint, such as $\|\boldsymbol{\beta}\|_1 \leq R$, used in all other recent works. Statistically, a δ -contraction property is established: each convex program contracts the previous estimator by a δ -fraction until the optimal statistical error is reached. Computationally, a phase transition in algorithmic convergence is established. The contraction stage enjoys only a sublinear rate of convergence while the tightening stage converges geometrically fast.

Recently, Negahban et al. (2012) proposed the restricted eigenvalue condition for unified M-estimators. Loh (2017) leveraged this condition, which is more related to our localized conditions. However, there are two major differences. First, their local parameter r is fixed at a constant independent of n, d, s, while we allow it to go to 0 as long as $r \gtrsim \sqrt{s \log d/n}$. Second, their high-dimensional regression problem relies on the ℓ_1 ball constraint $\|\boldsymbol{\beta}\|_1 \leq R$, while our newly developed localized analysis, together with the localized conditions, removes such type of constraint. In Lozano and Meinshausen (2013), the authors only consider the solutions in a local cone, which makes their analysis much simpler than ours. In this paper, we provide a stronger result: with high probability, all local solutions must fall in a local sparse (or ℓ_1) cone, and thus makes the localized eigenvalue conditions applicable.

More recently, Wang, Kim and Li (2013) proposed a two-step approach named calibrated CCCP which achieve strong oracle properties when using the Lasso estimator as initialization. Our work differs from theirs in two aspects. First, their work aims at analyzing the least square loss while our analysis handles much broader families of loss functions. Second, their procedure attains an oracle rate but requires the minimum signal strength to be in the order of $s\sqrt{\log d/n}$. Such a requirement is suboptimal. In contrast, our results requires only $\sqrt{\log d/n}$. This weakened assumption on minimum signal strength also distinguishes I-LAMM from other convex procedures, such as least squares refit after model selection [Belloni and Chernozhukov (2013)]. In Wang, Kim and Li (2013), the authors also proposed a high-dimensional BIC criterion for variable selection and finding the oracle estimator along the solution path. We believe such a criterion can also be applied to our framework under general conditions. In further studies, Loh and Wainwright (2014, 2015) and Wang, Liu and Zhang (2014) study the theoretical properties of nonconvex penalized M-estimators. Specifically, Loh and Wainwright (2015) and Loh and Wainwright (2014) provide conditions under which all the local optima obtained by an ℓ_1 -ball constrained optimization enjoys desired statistical rates. Wang, Liu and Zhang (2014) propose a path-following strategy to obtain optimal computational and statistical rates of convergence, which also relies an extra ball constraint.

Our work differs from the aforementioned literature at least in three aspects:

- (1) Our theory exploits a new notion of localized analysis, which is not available in Loh and Wainwright (2014, 2015) and Wang, Liu and Zhang (2014). Such analysis allows us to eliminate the extra ball constraints in previous work, which introduce more tuning effort and are intuitively redundant given the penalty function.
- (2) Our statistical results tolerate explicit computational precisions and are valid for all obtained approximate solutions, while the analysis in Loh and Wainwright (2015) only targets on the exact local solutions. Moreover, our computational result does not rely on the path-following type strategy as in Wang, Liu and Zhang (2014) and is valid for any algorithm with desired statistical properties as basic building blocks within each of the tightening steps.
- (3) We provide a refined oracle statistical rate $\sqrt{s/n}$ for the obtained approximation solution, while Loh and Wainwright (2015) and Wang, Liu and Zhang (2014) do not provide such a result. Loh and Wainwright (2015) provide a statistical rate which is also achievable using the convex Lasso penalty. Wang, Liu and Zhang (2014) only prove the oracle rate for exact local solutions.

Our work can be applied to many different topics: low-rank matrix completion problems, high-dimensional graphical models, quantile regression and many others. We conjecture that in all of the aforementioned topics, I-LAMM can give a faster rate by approximately solving a sequence of convex programs, with controlled computing resources. It is also interesting to see how our algorithm works in large-scale distributed systems. Is there any fundamental tradeoffs between statistical efficiency, communication and time complexity? We leave these as future research projects.

SUPPLEMENTARY MATERIAL

Supplement to "I-LAMM for Sparse learning: simultaneous control of algorithmic complexity and statistical error" (DOI: 10.1214/17-AOS1568SUPP; .pdf). The Supplementary Material [Fan et al. (2018)] contains proofs for Corollary 4.3, Theorem 4.4, Proposition 4.5, Proposition 4.6 and Theorem 4.7 in Section 4. It collects proofs of the lemmas presented in Section 5. An application to robust linear regression is given in Appendix D. Other technical lemmas are collected in Appendices E and F.

REFERENCES

- AGARWAL, A., NEGAHBAN, S. and WAINWRIGHT, M. J. (2012). Fast global convergence of gradient methods for high-dimensional statistical recovery. *Ann. Statist.* **40** 2452–2482. MR3097609
- BECK, A. and TEBOULLE, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM J. Imaging Sci. 2 183–202. MR2486527
- BELLONI, A. and CHERNOZHUKOV, V. (2013). Least squares after model selection in highdimensional sparse models. *Bernoulli* 19 521–547. MR3037163
- BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. Ann. Statist. 37 1705–1732. MR2533469
- BOYD, S. and VANDENBERGHE, L. (2004). Convex Optimization. Cambridge Univ. Press, Cambridge. MR2061575
- BREHENY, P. and HUANG, J. (2011). Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. Ann. Appl. Stat. 5 232–253. MR2810396
- BÜHLMANN, P. and VAN DE GEER, S. (2011). *Statistics for High-Dimensional Data*. Springer, Heidelberg. MR2807761
- BUNEA, F., TSYBAKOV, A. and WEGKAMP, M. (2007). Sparsity oracle inequalities for the Lasso. *Electron. J. Stat.* **1** 169–194. MR2312149
- FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. J. Amer. Statist. Assoc. 96 1348–1360. MR1946581
- FAN, J. and LV, J. (2008). Sure independence screening for ultrahigh dimensional feature space. J. R. Stat. Soc. Ser. B. Stat. Methodol. 70 849–911. MR2530322
- FAN, J. and LV, J. (2011). Nonconcave penalized likelihood with NP-dimensionality. IEEE Trans. Inform. Theory 57 5467–5484. MR2849368
- FAN, J., XUE, L. and ZOU, H. (2014). Strong oracle optimality of folded concave penalized estimation. Ann. Statist. 42 819–849. MR3210988
- FAN, J., LIU, H., SUN, Q. and ZHANG, T. (2018). Supplement to "I-LAMM for sparse learning: Simultaneous control of algorithmic complexity and statistical error." DOI:10.1214/17-AOS1568SUPP.

- FRIEDMAN, J., HASTIE, T., HÖFLING, H. and TIBSHIRANI, R. (2007). Pathwise coordinate optimization. Ann. Appl. Stat. 1 302–332. MR2415737
- HUNTER, D. R. and LANGE, K. (2004). A tutorial on MM algorithms. Amer. Statist. 58 30–37. MR2055509
- KIM, Y., CHOI, H. and OH, H.-S. (2008). Smoothly clipped absolute deviation on high dimensions. J. Amer. Statist. Assoc. 103 1665–1673. MR2510294
- KIM, Y. and KWON, S. (2012). Global optimality of nonconvex penalized estimators. *Biometrika* 99 315–325. MR2931256
- LANGE, K., HUNTER, D. R. and YANG, I. (2000). Optimization transfer using surrogate objective functions. J. Comput. Graph. Statist. 9 1–59. MR1819865
- LOH, P.-L. (2017). Statistical consistency and asymptotic normality for high-dimensional robust *M*-estimators. *Ann. Statist.* **45** 866–896. MR3650403
- LOH, P.-L. and WAINWRIGHT, M. J. (2014). Support recovery without incoherence: A case for nonconvex regularization. To appear. Available at arXiv:1412.5632.
- LOH, P.-L. and WAINWRIGHT, M. J. (2015). Regularized *M*-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *J. Mach. Learn. Res.* **16** 559–616. MR3335800
- LOZANO, A. C. and MEINSHAUSEN, N. (2013). Minimum distance estimation for robust highdimensional regression. Available at arXiv:1307.3227.
- NEGAHBAN, S. N., RAVIKUMAR, P., WAINWRIGHT, M. J. and YU, B. (2012). A unified framework for high-dimensional analysis of *M*-estimators with decomposable regularizers. *Statist. Sci.* 27 538–557. MR3025133
- NESTEROV, Y. (2013). Gradient methods for minimizing composite functions. *Math. Program.* 140 125–161. MR3071865
- RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2010). Restricted eigenvalue properties for correlated Gaussian designs. J. Mach. Learn. Res. 11 2241–2259. MR2719855
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. MR1379242
- VAN DE GEER, S. A. and BÜHLMANN, P. (2009). On the conditions used to prove oracle results for the Lasso. *Electron. J. Stat.* **3** 1360–1392. MR2576316
- WANG, L., KIM, Y. and LI, R. (2013). Calibrating nonconvex penalized regression in ultra-high dimension. Ann. Statist. 41 2505–2536. MR3127873
- WANG, Z., LIU, H. and ZHANG, T. (2014). Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *Ann. Statist.* 42 2164–2201. MR3269977
- ZHANG, T. (2009). Some sharp performance bounds for least squares regression with L₁ regularization. Ann. Statist. 37 2109–2144. MR2543687
- ZHANG, C.-H. (2010a). Nearly unbiased variable selection under minimax concave penalty. Ann. Statist. 38 894–942. MR2604701
- ZHANG, T. (2010b). Analysis of multi-stage convex relaxation for sparse regularization. J. Mach. Learn. Res. 11 1081–1107. MR2629825
- ZHANG, C.-H. and ZHANG, T. (2012). A general theory of concave regularization for highdimensional sparse estimation problems. *Statist. Sci.* 27 576–593. MR3025135
- ZOU, H. (2006). The adaptive lasso and its oracle properties. J. Amer. Statist. Assoc. 101 1418–1429. MR2279469
- ZOU, H. and LI, R. (2008). One-step sparse estimates in nonconcave penalized likelihood models. Ann. Statist. 36 1509–1533. MR2435443

LOCAL ADAPTIVE MAJORIZE-MINIMIZATION

J. FAN SCHOOL OF DATA SCIENCE FUDAN UNIVERSITY SHANGHAI CHINA AND DEPARTMENT OF OPERATIONS RESEARCH AND FINANCIAL ENGINEERING PRINCETON UNIVERSITY PRINCETON, NEW JERSEY 08544 USA E-MAIL: jqfan@princeton.edu Q. SUN DEPARTMENT OF STATISTICAL SCIENCES

PRINCETON UNIVERSITY PRINCETON, NEW JERSEY 08544 USA E-MAIL: hanliu@princeton.edu

H. LIU

Q. SUN DEPARTMENT OF STATISTICAL SCIENCES DEPARTMENT OF COMPUTER AND MATHEMATICAL SCIENCES UNIVERSITY OF TORONTO TORONTO, ONTARIO M5S 3G3 CANADA E-MAIL: qsun@utstat.toronto.edu T. ZHANG TENCENT AI LAB SHENNAN AVE, NANSHAN DISTRICT SHEN ZHEN GUANGDONG CHINA AND SCHOOL OF DATA SCIENCE FUDAN UNIVERSITY SHANGHAI CHINA E-MAIL: tongzhang@tongzhang-ml.org

DEPARTMENT OF OPERATIONS RESEARCH

AND FINANCIAL ENGINEERING