
High Dimensional EM Algorithm: Statistical Optimization and Asymptotic Normality*

Zhaoran Wang
Princeton University

Quanquan Gu
University of Virginia

Yang Ning
Princeton University

Han Liu
Princeton University

Abstract

We provide a general theory of the expectation-maximization (EM) algorithm for inferring high dimensional latent variable models. In particular, we make two contributions: (i) For parameter estimation, we propose a novel high dimensional EM algorithm which naturally incorporates sparsity structure into parameter estimation. With an appropriate initialization, this algorithm converges at a geometric rate and attains an estimator with the (near-)optimal statistical rate of convergence. (ii) Based on the obtained estimator, we propose a new inferential procedure for testing hypotheses for low dimensional components of high dimensional parameters. For a broad family of statistical models, our framework establishes the first computationally feasible approach for optimal estimation and asymptotic inference in high dimensions.

1 Introduction

The expectation-maximization (EM) algorithm [12] is the most popular approach for calculating the maximum likelihood estimator of latent variable models. Nevertheless, due to the nonconcavity of the likelihood function of latent variable models, the EM algorithm generally only converges to a local maximum rather than the global one [30]. On the other hand, existing statistical guarantees for latent variable models are only established for global optima [3]. Therefore, there exists a gap between computation and statistics.

Significant progress has been made toward closing the gap between the local maximum attained by the EM algorithm and the maximum likelihood estimator [2, 18, 25, 30]. In particular, [30] first establish general sufficient conditions for the convergence of the EM algorithm. [25] further improve this result by viewing the EM algorithm as a proximal point method applied to the Kullback-Leibler divergence. See [18] for a detailed survey. More recently, [2] establish the first result that characterizes explicit statistical and computational rates of convergence for the EM algorithm. They prove that, given a suitable initialization, the EM algorithm converges at a geometric rate to a local maximum close to the maximum likelihood estimator. All these results are established in the low dimensional regime where the dimension d is much smaller than the sample size n .

In high dimensional regimes where the dimension d is much larger than the sample size n , there exists no theoretical guarantee for the EM algorithm. In fact, when $d \gg n$, the maximum likelihood estimator is in general not well defined, unless the models are carefully regularized by sparsity-type assumptions. Furthermore, even if a regularized maximum likelihood estimator can be obtained in a computationally tractable manner, establishing the corresponding statistical properties, especially asymptotic normality, can still be challenging because of the existence of high dimensional nuisance parameters. To address such a challenge, we develop a general inferential theory of the EM algorithm for parameter estimation and uncertainty assessment of high dimensional latent variable models. In particular, we make two contributions in this paper:

- For high dimensional parameter estimation, we propose a novel high dimensional EM algorithm by attaching a truncation step to the expectation step (E-step) and maximization step (M-step). Such a

*Research supported by NSF IIS1116730, NSF IIS1332109, NSF IIS1408910, NSF IIS1546482-BIGDATA, NSF DMS1454377-CAREER, NIH R01GM083084, NIH R01HG06841, NIH R01MH102339, and FDA HHSF223201000072C.

truncation step effectively enforces the sparsity of the attained estimator and allows us to establish significantly improved statistical rate of convergence.

- Based upon the estimator attained by the high dimensional EM algorithm, we propose a decorrelated score statistic for testing hypotheses related to low dimensional components of the high dimensional parameter.

Under a unified analytic framework, we establish simultaneous statistical and computational guarantees for the proposed high dimensional EM algorithm and the respective uncertainty assessment procedure. Let $\beta^* \in \mathbb{R}^d$ be the true parameter, s^* be its sparsity level and $\{\beta^{(t)}\}_{t=0}^T$ be the iterative solution sequence of the high dimensional EM algorithm with T being the total number of iterations. In particular, we prove that:

- Given an appropriate initialization β^{init} with relative error upper bounded by a constant $\kappa \in (0, 1)$, i.e., $\|\beta^{\text{init}} - \beta^*\|_2 / \|\beta^*\|_2 \leq \kappa$, the iterative solution sequence $\{\beta^{(t)}\}_{t=0}^T$ satisfies

$$\|\beta^{(t)} - \beta^*\|_2 \leq \underbrace{\Delta_1 \cdot \rho^{t/2}}_{\text{Optimization Error}} + \underbrace{\Delta_2 \cdot \sqrt{s^* \cdot \log d/n}}_{\text{Statistical Error: Optimal Rate}} \quad (1.1)$$

with high probability. Here $\rho \in (0, 1)$, and Δ_1, Δ_2 are quantities that possibly depend on ρ, κ and β^* . As the optimization error term in (1.1) decreases to zero at a geometric rate with respect to t , the overall estimation error achieves the $\sqrt{s^* \cdot \log d/n}$ statistical rate of convergence (up to an extra factor of $\log n$), which is (near-)minimax-optimal. See Theorem 3.4 for details.

- The proposed decorrelated score statistic is asymptotically normal. Moreover, its limiting variance is optimal in the sense that it attains the semiparametric information bound for the low dimensional components of interest in the presence of high dimensional nuisance parameters. See Theorem 4.6 for details.

Our framework allows two implementations of the M-step: the exact maximization versus approximate maximization. The former one calculates the maximizer exactly, while the latter one conducts an approximate maximization through a gradient ascent step. Our framework is quite general. We illustrate its effectiveness by applying it to two high dimensional latent variable models, that is, Gaussian mixture model and mixture of regression model.

Comparison with Related Work: A closely related work is by [2], which considers the low dimensional regime where d is much smaller than n . Under certain initialization conditions, they prove that the EM algorithm converges at a geometric rate to some local optimum that attains the $\sqrt{d/n}$ statistical rate of convergence. They cover both maximization and gradient ascent implementations of the M-step, and establish the consequences for the two latent variable models considered in our paper under low dimensional settings. Our framework adopts their view of treating the EM algorithm as a perturbed version of gradient methods. However, to handle the challenge of high dimensionality, the key ingredient of our framework is the truncation step that enforces the sparsity structure along the solution path. Such a truncation operation poses significant challenges for both computational and statistical analysis. In detail, for computational analysis we need to carefully characterize the evolution of each intermediate solution’s support and its effects on the evolution of the entire iterative solution sequence. For statistical analysis, we need to establish a fine-grained characterization of the entrywise statistical error, which is technically more challenging than just establishing the ℓ_2 -norm error employed by [2]. In high dimensional regimes, we need to establish the $\sqrt{s^* \cdot \log d/n}$ statistical rate of convergence, which is much sharper than their $\sqrt{d/n}$ rate when $d \gg n$. In addition to point estimation, we further construct hypothesis tests for latent variable models in the high dimensional regime, which have not been established before.

High dimensionality poses significant challenges for assessing the uncertainty (e.g., testing hypotheses) of the constructed estimators. For example, [15] show that the limiting distribution of the Lasso estimator is not Gaussian even in the low dimensional regime. A variety of approaches have been proposed to correct the Lasso estimator to attain asymptotic normality, including the debiasing method [13], the desparsification methods [26, 32] as well as instrumental variable-based methods [4]. Meanwhile, [16, 17, 24] propose the post-selection procedures for exact inference. In addition, several authors propose methods based on data splitting [20, 29], stability selection [19] and ℓ_2 -confidence sets [22]. However, these approaches mainly focus on generalized linear models rather than latent variable models. In addition, their results heavily rely on the fact that the estimator is a global optimum of a convex program. In comparison, our approach applies to a much broader family of statistical models with latent structures. For these latent variable models, it is computationally infeasible to

obtain the global maximum of the penalized likelihood due to the nonconcavity of the likelihood function. Unlike existing approaches, our inferential theory is developed for the estimator attained by the proposed high dimensional EM algorithm, which is not necessarily a global optimum to any optimization formulation.

Another line of research for the estimation of latent variable models is the tensor method, which exploits the structures of third or higher order moments. See [1] and the references therein. However, existing tensor methods primarily focus on the low dimensional regime where $d \ll n$. In addition, since the high order sample moments generally have a slow statistical rate of convergence, the estimators obtained by the tensor methods usually have a suboptimal statistical rate even for $d \ll n$. For example, [9] establish the $\sqrt{d^6/n}$ statistical rate of convergence for mixture of regression model, which is suboptimal compared with the $\sqrt{d/n}$ minimax lower bound. Similarly, in high dimensional settings, the statistical rates of convergence attained by tensor methods are significantly slower than the statistical rate obtained in this paper.

The latent variable models considered in this paper have been well studied. Nevertheless, only a few works establish theoretical guarantees for the EM algorithm. In particular, for Gaussian mixture model, [10, 11] establish parameter estimation guarantees for the EM algorithm and its extensions. For mixture of regression model, [31] establish exact parameter recovery guarantees for the EM algorithm under a noiseless setting. For high dimensional mixture of regression model, [23] analyze the gradient EM algorithm for the ℓ_1 -penalized log-likelihood. They establish support recovery guarantees for the attained local optimum but have no parameter estimation guarantees. In comparison with existing works, this paper establishes a general inferential framework for simultaneous parameter estimation and uncertainty assessment based on a novel high dimensional EM algorithm. Our analysis provides the first theoretical guarantee of parameter estimation and asymptotic inference in high dimensional regimes for the EM algorithm and its applications to a broad family of latent variable models.

Notation: The matrix (p, q) -norm, i.e., $\|\cdot\|_{p,q}$, is obtained by taking the ℓ_p -norm of each row and then taking the ℓ_q -norm of the obtained row norms. We use C, C', \dots to denote generic constants. Their values may vary from line to line. We will introduce more notations in §2.2.

2 Methodology

We first introduce the high dimensional EM Algorithm and then the respective inferential procedure. As examples, we consider their applications to Gaussian mixture model and mixture of regression model. For compactness, we defer the details to §A of the appendix. More models are included in the longer version of this paper.

Algorithm 1 High Dimensional EM Algorithm

- 1: **Parameter:** Sparsity Parameter \hat{s} , Maximum Number of Iterations T
 - 2: **Initialization:** $\hat{\mathcal{S}}^{\text{init}} \leftarrow \text{supp}(\beta^{\text{init}}, \hat{s})$, $\beta^{(0)} \leftarrow \text{trunc}(\beta^{\text{init}}, \hat{\mathcal{S}}^{\text{init}})$
 $\{\text{supp}(\cdot, \cdot) \text{ and } \text{trunc}(\cdot, \cdot) \text{ are defined in (2.2) and (2.3)}\}$
 - 3: **For** $t = 0$ to $T - 1$
 - 4: **E-step:** Evaluate $Q_n(\beta; \beta^{(t)})$
 - 5: **M-step:** $\beta^{(t+0.5)} \leftarrow M_n(\beta^{(t)})$ $\{M_n(\cdot) \text{ is implemented as in Algorithm 2 or 3}\}$
 - 6: **T-step:** $\hat{\mathcal{S}}^{(t+0.5)} \leftarrow \text{supp}(\beta^{(t+0.5)}, \hat{s})$, $\beta^{(t+1)} \leftarrow \text{trunc}(\beta^{(t+0.5)}, \hat{\mathcal{S}}^{(t+0.5)})$
 - 7: **End For**
 - 8: **Output:** $\hat{\beta} \leftarrow \beta^{(T)}$
-

Algorithm 2 Maximization Implementation of the M-step

- 1: **Input:** $\beta^{(t)}$, $Q_n(\beta; \beta^{(t)})$ **Output:** $M_n(\beta^{(t)}) \leftarrow \text{argmax}_{\beta} Q_n(\beta; \beta^{(t)})$
-

Algorithm 3 Gradient Ascent Implementation of the M-step

- 1: **Input:** $\beta^{(t)}$, $Q_n(\beta; \beta^{(t)})$ **Parameter:** Stepsize $\eta > 0$
 - 2: **Output:** $M_n(\beta^{(t)}) \leftarrow \beta^{(t)} + \eta \cdot \nabla Q_n(\beta^{(t)}; \beta^{(t)})$
-

2.1 High Dimensional EM Algorithm

Before we introduce the proposed high dimensional EM Algorithm (Algorithm 1), we briefly review the classical EM algorithm. Let $h_{\beta}(\mathbf{y})$ be the probability density function of $\mathbf{Y} \in \mathcal{Y}$, where $\beta \in \mathbb{R}^d$ is the model parameter. For latent variable models, we assume that $h_{\beta}(\mathbf{y})$ is obtained by marginalizing over an unobserved latent variable $\mathbf{Z} \in \mathcal{Z}$, i.e., $h_{\beta}(\mathbf{y}) = \int_{\mathcal{Z}} f_{\beta}(\mathbf{y}, \mathbf{z}) d\mathbf{z}$. Let $k_{\beta}(\mathbf{z} | \mathbf{y})$ be the density

of \mathbf{Z} conditioning on the observed variable $\mathbf{Y} = \mathbf{y}$, i.e., $k_{\beta}(\mathbf{z} | \mathbf{y}) = f_{\beta}(\mathbf{y}, \mathbf{z})/h_{\beta}(\mathbf{y})$. We define

$$Q_n(\beta; \beta') = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} k_{\beta'}(\mathbf{z} | \mathbf{y}_i) \cdot \log f_{\beta}(\mathbf{y}_i, \mathbf{z}) \, d\mathbf{z}. \quad (2.1)$$

See §B of the appendix for a detailed derivation. At the t -th iteration of the classical EM algorithm, we evaluate $Q_n(\beta; \beta^{(t)})$ at the E-step and then perform $\max_{\beta} Q_n(\beta; \beta^{(t)})$ at the M-step. The proposed high dimensional EM algorithm (Algorithm 1) is built upon the E-step and M-step (lines 4 and 5) of the classical EM algorithm. In addition to the exact maximization implementation of the M-step (Algorithm 2), we allow the gradient ascent implementation of the M-step (Algorithm 3), which performs an approximate maximization via a gradient ascent step. To handle the challenge of high dimensionality, in line 6 of Algorithm 1 we perform a truncation step (T-step) to enforce the sparsity structure. In detail, we define

$$\text{supp}(\beta, s): \text{The set of index } j\text{'s corresponding to the top } s \text{ largest } |\beta_j|\text{'s.} \quad (2.2)$$

Also, for an index set $\mathcal{S} \subseteq \{1, \dots, d\}$, we define the $\text{trunc}(\cdot, \cdot)$ function in line 6 as

$$[\text{trunc}(\beta, \mathcal{S})]_j = \beta_j \cdot \mathbf{1}\{j \in \mathcal{S}\}. \quad (2.3)$$

Note that $\beta^{(t+0.5)}$ is the output of the M-step (line 5) at the t -th iteration of the high dimensional EM algorithm. To obtain $\beta^{(t+1)}$, the T-step (line 6) preserves the entries of $\beta^{(t+0.5)}$ with the top \hat{s} large magnitudes and sets the rest to zero. Here \hat{s} is a tuning parameter that controls the sparsity level (line 1). By iteratively performing the E-step, M-step and T-step, the high dimensional EM algorithm attains an \hat{s} -sparse estimator $\hat{\beta} = \beta^{(T)}$ (line 8). Here T is the total number of iterations.

2.2 Asymptotic Inference

Notation: Let $\nabla_1 Q(\beta; \beta')$ be the gradient with respect to β and $\nabla_2 Q(\beta; \beta')$ be the gradient with respect to β' . If there is no confusion, we simply denote $\nabla Q(\beta; \beta') = \nabla_1 Q(\beta; \beta')$ as in the previous sections. We define the higher order derivatives in the same manner, e.g., $\nabla_{1,2}^2 Q(\beta; \beta')$ is calculated by first taking derivative with respect to β and then with respect to β' . For $\beta = (\beta_1^\top, \beta_2^\top)^\top \in \mathbb{R}^d$ with $\beta_1 \in \mathbb{R}^{d_1}$, $\beta_2 \in \mathbb{R}^{d_2}$ and $d_1 + d_2 = d$, we use notations such as $\mathbf{v}_{\beta_1} \in \mathbb{R}^{d_1}$ and $\mathbf{A}_{\beta_1, \beta_2} \in \mathbb{R}^{d_1 \times d_2}$ to denote the corresponding subvector of $\mathbf{v} \in \mathbb{R}^d$ and the submatrix of $\mathbf{A} \in \mathbb{R}^{d \times d}$.

We aim to conduct asymptotic inference for low dimensional components of the high dimensional parameter β^* . Without loss of generality, we consider a single entry of β^* . In particular, we assume $\beta^* = [\alpha^*, (\gamma^*)^\top]^\top$, where $\alpha^* \in \mathbb{R}$ is the entry of interest, while $\gamma^* \in \mathbb{R}^{d-1}$ is treated as the nuisance parameter. In the following, we construct a high dimensional score test named decorrelated score test. It is worth noting that, our method and theory can be easily generalized to perform statistical inference for an arbitrary low dimensional subvector of β^* .

Decorrelated Score Test: For score test, we are primarily interested in testing $H_0 : \alpha^* = 0$, since this null hypothesis characterizes the uncertainty in variable selection. Our method easily generalizes to $H_0 : \alpha^* = \alpha_0$ with $\alpha_0 \neq 0$. For notational simplicity, we define the following key quantity

$$T_n(\beta) = \nabla_{1,1}^2 Q_n(\beta; \beta) + \nabla_{1,2}^2 Q_n(\beta; \beta) \in \mathbb{R}^{d \times d}. \quad (2.4)$$

Let $\beta = (\alpha, \gamma^\top)^\top$. We define the decorrelated score function $S_n(\cdot, \cdot) \in \mathbb{R}$ as

$$S_n(\beta, \lambda) = [\nabla_1 Q_n(\beta; \beta)]_{\alpha} - w(\beta, \lambda)^\top \cdot [\nabla_1 Q_n(\beta; \beta)]_{\gamma}. \quad (2.5)$$

Here $w(\beta, \lambda) \in \mathbb{R}^{d-1}$ is obtained using the following Dantzig selector [8]

$$w(\beta, \lambda) = \underset{\mathbf{w} \in \mathbb{R}^{d-1}}{\text{argmin}} \|\mathbf{w}\|_1, \quad \text{subject to } \|[T_n(\beta)]_{\gamma, \alpha} - [T_n(\beta)]_{\gamma, \gamma} \cdot \mathbf{w}\|_{\infty} \leq \lambda, \quad (2.6)$$

where $\lambda > 0$ is a tuning parameter. Let $\hat{\beta} = (\hat{\alpha}, \hat{\gamma}^\top)^\top$, where $\hat{\beta}$ is the estimator attained by the high dimensional EM algorithm (Algorithm 1). We define the decorrelated score statistic as

$$\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda) / \{-[T_n(\hat{\beta}_0)]_{\alpha|\gamma}\}^{1/2}, \quad (2.7)$$

where $\hat{\beta}_0 = (0, \hat{\gamma}^\top)^\top$, and $[T_n(\hat{\beta}_0)]_{\alpha|\gamma} = [1, -w(\hat{\beta}_0, \lambda)^\top] \cdot T_n(\hat{\beta}_0) \cdot [1, -w(\hat{\beta}_0, \lambda)^\top]^\top$.

Here we use $\hat{\beta}_0$ instead of $\hat{\beta}$ since we are interested in the null hypothesis $H_0 : \alpha^* = 0$. We can also replace $\hat{\beta}_0$ with $\hat{\beta}$ and the theoretical results will remain the same. In §4 we will prove the proposed decorrelated score statistic in (2.7) is asymptotically $N(0, 1)$. Consequently, the decorrelated score

test with significance level $\delta \in (0, 1)$ takes the form

$$\psi_S(\delta) = \mathbf{1}\{\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda) / \{-[T_n(\hat{\beta}_0)]_{\alpha|\gamma}\}^{1/2} \notin [-\Phi^{-1}(1 - \delta/2), \Phi^{-1}(1 - \delta/2)]\},$$

where $\Phi^{-1}(\cdot)$ is the inverse function of the Gaussian cumulative distribution function. If $\psi_S(\delta) = 1$, we reject the null hypothesis $H_0 : \alpha^* = 0$. The intuition of this decorrelated score test is explained in §D of the appendix. The key theoretical observation is Theorem 2.1, which connects $\nabla_1 Q_n(\cdot; \cdot)$ in (2.5) and $T_n(\cdot)$ in (2.7) with the score function and Fisher information in the presence of latent structures. Let $\ell_n(\beta)$ be the log-likelihood. Its score function is $\nabla \ell_n(\beta)$ and the Fisher information is $I(\beta^*) = -\mathbb{E}_{\beta^*} [\nabla^2 \ell_n(\beta^*)] / n$, where $\mathbb{E}_{\beta^*}(\cdot)$ is the expectation under the model with parameter β^* .

Theorem 2.1. For the true parameter β^* and any $\beta \in \mathbb{R}^d$, it holds that

$$\nabla_1 Q_n(\beta; \beta) = \nabla \ell_n(\beta) / n, \quad \text{and} \quad \mathbb{E}_{\beta^*} [T_n(\beta^*)] = -I(\beta^*) = \mathbb{E}_{\beta^*} [\nabla^2 \ell_n(\beta^*)] / n. \quad (2.8)$$

Proof. See §I.1 of the appendix for a detailed proof. \square

Based on the decorrelated score test, it is easy to establish the decorrelated Wald test, which allows us to construct confidence intervals. For compactness we defer it to the longer version of this paper.

3 Theory of Computation and Estimation

Before we present the main results, we introduce three technical conditions, which will significantly ease our presentation. They will be verified for specific latent variable models in §E of the appendix. The first two conditions, proposed by [2], characterize the properties of the population version lower bound function $Q(\cdot; \cdot)$, i.e., the expectation of $Q_n(\cdot; \cdot)$ defined in (2.1). We define the respective population version M-step as follows. For the M-step in Algorithm 2, we define

$$M(\beta) = \operatorname{argmax}_{\beta'} Q(\beta'; \beta). \quad (3.1)$$

For the M-step in Algorithm 3, we define

$$M(\beta) = \beta + \eta \cdot \nabla_1 Q(\beta; \beta), \quad (3.2)$$

where $\eta > 0$ is the stepsize in Algorithm 3. We use \mathcal{B} to denote the basin of attraction, i.e., the local region where the high dimensional EM algorithm enjoys desired guarantees.

Condition 3.1. We define two versions of this condition.

- *Lipschitz-Gradient-1*(γ_1, \mathcal{B}). For the true parameter β^* and any $\beta \in \mathcal{B}$, we have

$$\|\nabla_1 Q[M(\beta); \beta^*] - \nabla_1 Q[M(\beta); \beta]\|_2 \leq \gamma_1 \cdot \|\beta - \beta^*\|_2, \quad (3.3)$$

where $M(\cdot)$ is the population version M-step (maximization implementation) defined in (3.1).

- *Lipschitz-Gradient-2*(γ_2, \mathcal{B}). For the true parameter β^* and any $\beta \in \mathcal{B}$, we have

$$\|\nabla_1 Q(\beta; \beta^*) - \nabla_1 Q(\beta; \beta)\|_2 \leq \gamma_2 \cdot \|\beta - \beta^*\|_2. \quad (3.4)$$

Condition 3.1 defines a variant of Lipschitz continuity for $\nabla_1 Q(\cdot; \cdot)$. In the sequel, we will use (3.3) and (3.4) in the analysis of the two implementations of the M-step respectively.

Condition 3.2 *Concavity-Smoothness*(μ, ν, \mathcal{B}). For any $\beta_1, \beta_2 \in \mathcal{B}$, $Q(\cdot; \beta^*)$ is μ -smooth, i.e.,

$$Q(\beta_1; \beta^*) \geq Q(\beta_2; \beta^*) + (\beta_1 - \beta_2)^\top \cdot \nabla_1 Q(\beta_2; \beta^*) - \mu/2 \cdot \|\beta_2 - \beta_1\|_2^2, \quad (3.5)$$

and ν -strongly concave, i.e.,

$$Q(\beta_1; \beta^*) \leq Q(\beta_2; \beta^*) + (\beta_1 - \beta_2)^\top \cdot \nabla_1 Q(\beta_2; \beta^*) - \nu/2 \cdot \|\beta_2 - \beta_1\|_2^2. \quad (3.6)$$

This condition indicates that, when the second variable of $Q(\cdot; \cdot)$ is fixed to be β^* , the function is ‘sandwiched’ between two quadratic functions. The third condition characterizes the statistical error between the sample version and population version M-steps, i.e., $M_n(\cdot)$ defined in Algorithms 2 and 3, and $M(\cdot)$ in (3.1) and (3.2). Recall $\|\cdot\|_0$ denotes the total number of nonzero entries in a vector.

Condition 3.3 *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$). For any fixed $\beta \in \mathcal{B}$ with $\|\beta\|_0 \leq s$, we have that

$$\|M(\beta) - M_n(\beta)\|_\infty \leq \epsilon \quad (3.7)$$

holds with probability at least $1 - \delta$. Here $\epsilon > 0$ possibly depends on δ , sparsity level s , sample size n , dimension d , as well as the basin of attraction \mathcal{B} .

In (3.7) the statistical error ϵ quantifies the ℓ_∞ -norm of the difference between the population version and sample version M-steps. Particularly, we constrain the input β of $M(\cdot)$ and $M_n(\cdot)$ to be s -sparse. Such a condition is different from the one used by [2]. In detail, they quantify the statistical error

with the ℓ_2 -norm and do not constrain the input of $M(\cdot)$ and $M_n(\cdot)$ to be sparse. Consequently, our subsequent statistical analysis is different from theirs. The reason we use the ℓ_∞ -norm is that, it characterizes the more refined entrywise statistical error, which converges at a fast rate of $\sqrt{\log d/n}$ (possibly with extra factors depending on specific models). In comparison, the ℓ_2 -norm statistical error converges at a slow rate of $\sqrt{d/n}$, which does not decrease to zero as n increases with $d \gg n$. Furthermore, the fine-grained entrywise statistical error is crucial to our key proof for quantifying the effects of the truncation step (line 6 of Algorithm 1) on the iterative solution sequence.

3.1 Main Results

To simplify the technical analysis of the high dimensional EM algorithm, we focus on its resampling version, which is illustrated in Algorithm 4 in §C of the appendix.

Theorem 3.4. We define $\mathcal{B} = \{\beta : \|\beta - \beta^*\|_2 \leq R\}$, where $R = \kappa \cdot \|\beta^*\|_2$ for some $\kappa \in (0, 1)$. We assume Condition *Concavity-Smoothness*(μ, ν, \mathcal{B}) holds and $\|\beta^{\text{init}} - \beta^*\|_2 \leq R/2$.

- For the maximization implementation of the M-step (Algorithm 2), we suppose that Condition *Lipschitz-Gradient-1*(γ_1, \mathcal{B}) holds with $\rho_1 := \gamma_1/\nu \in (0, 1)$ and

$$\hat{s} = \lceil C \cdot \max \{ 16/(1/\rho_1 - 1)^2, 4 \cdot (1 + \kappa)^2/(1 - \kappa)^2 \} \cdot s^* \rceil, \quad (3.8)$$

$$(\sqrt{\hat{s}} + C'/\sqrt{1 - \kappa} \cdot \sqrt{s^*}) \cdot \epsilon \leq \min \{ (1 - \sqrt{\rho_1})^2 \cdot R, (1 - \kappa)^2/[2 \cdot (1 + \kappa)] \cdot \|\beta^*\|_2 \}. \quad (3.9)$$

Here $C \geq 1$ and $C' > 0$ are constants. Under Condition *Statistical-Error*($\epsilon, \delta/T, \hat{s}, n/T, \mathcal{B}$) we have that, for $t = 1, \dots, T$,

$$\|\beta^{(t)} - \beta^*\|_2 \leq \underbrace{\rho_1^{t/2} \cdot R}_{\text{Optimization Error}} + \underbrace{(\sqrt{\hat{s}} + C'/\sqrt{1 - \kappa} \cdot \sqrt{s^*})/(1 - \sqrt{\rho_1}) \cdot \epsilon}_{\text{Statistical Error}} \quad (3.10)$$

holds with probability at least $1 - \delta$, where C' is the same constant as in (3.9).

- For the gradient ascent implementation of the M-step (Algorithm 3), we suppose that Condition *Lipschitz-Gradient-2*(γ_2, \mathcal{B}) holds with $\rho_2 := 1 - 2 \cdot (\nu - \gamma_2)/(\nu + \mu) \in (0, 1)$ and the stepsize in Algorithm 3 is set to $\eta = 2/(\nu + \mu)$. Meanwhile, we assume (3.8) and (3.9) hold with ρ_1 replaced by ρ_2 . Under Condition *Statistical-Error*($\epsilon, \delta/T, \hat{s}, n/T, \mathcal{B}$) we have that, for $t = 1, \dots, T$, (3.10) holds with probability at least $1 - \delta$, in which ρ_1 is replaced with ρ_2 .

Proof. See §G.1 of the appendix for a detailed proof. \square

The assumption in (3.8) states that the sparsity parameter \hat{s} is chosen to be sufficiently large and also of the same order as the true sparsity level s^* . This assumption ensures that the error incurred by the truncation step can be upper bounded. In addition, as is shown for specific latent variable models in §E of the appendix, the error term ϵ in Condition *Statistical-Error*($\epsilon, \delta/T, \hat{s}, n/T, \mathcal{B}$) decreases as sample size n increases. By the assumption in (3.8), $(\sqrt{\hat{s}} + C'/\sqrt{1 - \kappa} \cdot \sqrt{s^*})$ is of the same order as $\sqrt{s^*}$. Therefore, the assumption in (3.9) suggests the sample size n is sufficiently large such that $\sqrt{s^*} \cdot \epsilon$ is sufficiently small. These assumptions guarantee that the entire iterative solution sequence remains within the basin of attraction \mathcal{B} in the presence of statistical error.

Theorem 3.4 illustrates that, the upper bound of the overall estimation error can be decomposed into two terms. The first term is the upper bound of optimization error, which decreases to zero at a geometric rate of convergence, because we have $\rho_1, \rho_2 < 1$. Meanwhile, the second term is the upper bound of statistical error, which does not depend on t . Since $(\sqrt{\hat{s}} + C'/\sqrt{1 - \kappa} \cdot \sqrt{s^*})$ is of the same order as $\sqrt{s^*}$, this term is proportional to $\sqrt{s^*} \cdot \epsilon$, where ϵ is the entrywise statistical error between $M(\cdot)$ and $M_n(\cdot)$. In §E of the appendix we prove that, for each specific latent variable model, ϵ is roughly of the order $\sqrt{\log d/n}$. (There may be extra factors attached to ϵ depending on each specific model.) Therefore, the statistical error term is roughly of the order $\sqrt{s^*} \cdot \log d/n$. Consequently, for a sufficiently large $t = T$ such that the optimization and statistical error terms in (3.10) are of the same order, the final estimator $\hat{\beta} = \beta^{(T)}$ attains a (near-)optimal $\sqrt{s^*} \cdot \log d/n$ (possibly with extra factors) statistical rate. For compactness, we give the following example and defer the details to §E.

Implications for Gaussian Mixture Model: We assume $\mathbf{y}_1, \dots, \mathbf{y}_n$ are the n i.i.d. realizations of $\mathbf{Y} = Z \cdot \beta^* + \mathbf{V}$. Here Z is a Rademacher random variable, i.e., $\mathbb{P}(Z = +1) = \mathbb{P}(Z = -1) = 1/2$, and $\mathbf{V} \sim N(\mathbf{0}, \sigma^2 \cdot \mathbf{I}_d)$ is independent of Z , where σ is the standard deviation. Suppose that we have $\|\beta^*\|_2/\sigma \geq r$, where $r > 0$ is a sufficiently large constant that denotes the minimum signal-to-noise ratio. In §E of the appendix we prove that there exists some constant $C > 0$ such that Conditions

Lipschitz-Gradient-1(γ_1, \mathcal{B}) and *Concavity-Smoothness*(μ, ν, \mathcal{B}) hold with $\gamma_1 = \exp(-C \cdot r^2)$, $\mu = \nu = 1$, $\mathcal{B} = \{\beta : \|\beta - \beta^*\|_2 \leq R\}$ with $R = \kappa \cdot \|\beta^*\|_2$, $\kappa = 1/4$. For a sufficiently large n , we have that Condition *Statistical-Error*($\epsilon, \delta, s, n, \mathcal{B}$) holds with

$$\epsilon = C \cdot (\|\beta^*\|_\infty + \sigma) \cdot \sqrt{[\log d + \log(2/\delta)]/n}.$$

Then the first part of Theorem 3.4 implies $\|\widehat{\beta} - \beta^*\|_2 \leq C \cdot \sqrt{s^* \cdot \log d \cdot \log n/n}$ for a sufficiently large T , which is near-optimal with respect to the minimax lower bound $\sqrt{s^* \log d/n}$.

4 Theory of Inference

To simplify the presentation of the unified framework, we lay out several technical conditions, which will be verified for each model. Let $\zeta^{\text{EM}}, \zeta^{\text{G}}, \zeta^{\text{T}}$ and ζ^{L} be four quantities that scale with s^*, d and n . These conditions will be verified for specific latent variable models in §F of the appendix.

Condition 4.1 *Parameter-Estimation*(ζ^{EM}). We have $\|\widehat{\beta} - \beta^*\|_1 = O_{\mathbb{P}}(\zeta^{\text{EM}})$.

Condition 4.2 *Gradient-Statistical-Error*(ζ^{G}). We have

$$\|\nabla_1 Q_n(\beta^*; \beta^*) - \nabla_1 Q(\beta^*; \beta^*)\|_\infty = O_{\mathbb{P}}(\zeta^{\text{G}}).$$

Condition 4.3 *$T_n(\cdot)$ -Concentration*(ζ^{T}). We have $\|T_n(\beta^*) - \mathbb{E}_{\beta^*}[T_n(\beta^*)]\|_{\infty, \infty} = O_{\mathbb{P}}(\zeta^{\text{T}})$.

Condition 4.4 *$T_n(\cdot)$ -Lipschitz*(ζ^{L}). For any β , we have

$$\|T_n(\beta) - T_n(\beta^*)\|_{\infty, \infty} = O_{\mathbb{P}}(\zeta^{\text{L}}) \cdot \|\beta - \beta^*\|_1.$$

In the sequel, we lay out an assumption on several population quantities and the sample size n . Recall that $\beta^* = [\alpha^*, (\gamma^*)^\top]^\top$, where $\alpha^* \in \mathbb{R}$ is the entry of interest, while $\gamma^* \in \mathbb{R}^{d-1}$ is the nuisance parameter. By the notations in §2.2, $[I(\beta^*)]_{\gamma, \gamma} \in \mathbb{R}^{(d-1) \times (d-1)}$ and $[I(\beta^*)]_{\gamma, \alpha} \in \mathbb{R}^{(d-1) \times 1}$ denote the submatrices of the Fisher information matrix $I(\beta^*) \in \mathbb{R}^{d \times d}$. We define $\mathbf{w}^*, s_{\mathbf{w}}^*$ and $\mathcal{S}_{\mathbf{w}}^*$ as

$$\mathbf{w}^* = [I(\beta^*)]_{\gamma, \gamma}^{-1} \cdot [I(\beta^*)]_{\gamma, \alpha} \in \mathbb{R}^{d-1}, \quad s_{\mathbf{w}}^* = \|\mathbf{w}^*\|_0, \quad \text{and} \quad \mathcal{S}_{\mathbf{w}}^* = \text{supp}(\mathbf{w}^*). \quad (4.1)$$

We define $\lambda_1[I(\beta^*)]$ and $\lambda_d[I(\beta^*)]$ as the largest and smallest eigenvalues of $I(\beta^*)$, and

$$[I(\beta^*)]_{\alpha|\gamma} = [I(\beta^*)]_{\alpha, \alpha} - [I(\beta^*)]_{\gamma, \alpha}^\top \cdot [I(\beta^*)]_{\gamma, \gamma}^{-1} \cdot [I(\beta^*)]_{\gamma, \alpha} \in \mathbb{R}. \quad (4.2)$$

According to (4.1) and (4.2), we can easily verify that

$$[I(\beta^*)]_{\alpha|\gamma} = [1, -(\mathbf{w}^*)^\top] \cdot I(\beta^*) \cdot [1, -(\mathbf{w}^*)^\top]^\top. \quad (4.3)$$

The following assumption ensures that $\lambda_d[I(\beta^*)] > 0$. Hence, $[I(\beta^*)]_{\gamma, \gamma}$ in (4.1) is invertible. Also, according to (4.3) and the fact that $\lambda_d[I(\beta^*)] > 0$, we have $[I(\beta^*)]_{\alpha|\gamma} > 0$.

Assumption 4.5. We impose the following assumptions.

- For positive constants ρ_{\max} and ρ_{\min} , we assume

$$\rho_{\max} \geq \lambda_1[I(\beta^*)] \geq \lambda_d[I(\beta^*)] \geq \rho_{\min}, \quad [I(\beta^*)]_{\alpha|\gamma} = O(1), \quad [I(\beta^*)]_{\alpha|\gamma}^{-1} = O(1). \quad (4.4)$$

- The tuning parameter λ of the Dantzig selector in (2.6) is set to

$$\lambda = C \cdot (\zeta^{\text{T}} + \zeta^{\text{L}} \cdot \zeta^{\text{EM}}) \cdot (1 + \|\mathbf{w}^*\|_1), \quad (4.5)$$

where $C \geq 1$ is a sufficiently large constant. The sample size n is sufficiently large such that

$$\begin{aligned} \max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1), \quad \zeta^{\text{EM}} = o(1), \quad s_{\mathbf{w}}^* \cdot \lambda \cdot \zeta^{\text{G}} = o(1/\sqrt{n}), \\ \lambda \cdot \zeta^{\text{EM}} = o(1/\sqrt{n}), \quad \max\{1, \|\mathbf{w}^*\|_1\} \cdot \zeta^{\text{L}} \cdot (\zeta^{\text{EM}})^2 = o(1/\sqrt{n}). \end{aligned} \quad (4.6)$$

The assumption on $\lambda_d[I(\beta^*)]$ guarantees that the Fisher information matrix is positive definite. The other assumptions in (4.4) guarantee the existence of the asymptotic variance of $\sqrt{n} \cdot S_n(\widehat{\beta}_0, \lambda)$ in the score statistic defined in (2.7). Similar assumptions are standard in existing asymptotic inference results. For example, for mixture of regression model, [14] impose variants of these assumptions. For specific models, we will show that $\zeta^{\text{EM}}, \zeta^{\text{G}}, \zeta^{\text{T}}$ and λ all decrease with n , while ζ^{L} increases with n at a slow rate. Therefore, the assumptions in (4.6) ensure that the sample size n is sufficiently large. We will make these assumptions more explicit after we specify $\zeta^{\text{EM}}, \zeta^{\text{G}}, \zeta^{\text{T}}$ and ζ^{L} for each

model. Note the assumptions in (4.6) imply that $s_{\mathbf{w}}^* = \|\mathbf{w}^*\|_0$ needs to be small. For instance, for λ specified in (4.5), $\max\{\|\mathbf{w}^*\|_1, 1\} \cdot s_{\mathbf{w}}^* \cdot \lambda = o(1)$ in (4.6) implies $s_{\mathbf{w}}^* \cdot \zeta^T = o(1)$. In the following, we will prove that ζ^T is of the order $\sqrt{\log d/n}$. Hence, we require that $s_{\mathbf{w}}^* = o(\sqrt{n/\log d}) \ll d-1$, i.e., $\mathbf{w}^* \in \mathbb{R}^{d-1}$ is sparse. Such a sparsity assumption can be understood as follows. According to the definition of \mathbf{w}^* in (4.1), we have $[I(\beta^*)]_{\gamma, \gamma} \cdot \mathbf{w}^* = [I(\beta^*)]_{\gamma, \alpha}$. Therefore, such a sparsity assumption suggests $[I(\beta^*)]_{\gamma, \alpha}$ lies within the span of a few columns of $[I(\beta^*)]_{\gamma, \gamma}$. Such a sparsity assumption on \mathbf{w}^* is necessary, because otherwise it is difficult to accurately estimate \mathbf{w}^* in high dimensional regimes. In the context of high dimensional generalized linear models, [26, 32] impose similar sparsity assumptions.

4.1 Main Results

Decorrelated Score Test: The next theorem establishes the asymptotic normality of the decorrelated score statistic defined in (2.7).

Theorem 4.6. We consider $\beta^* = [\alpha^*, (\gamma^*)^T]^T$ with $\alpha^* = 0$. Under Assumption 4.5 and Conditions 4.1-4.4, we have that for $n \rightarrow \infty$,

$$\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda) / \{-[T_n(\hat{\beta}_0)]_{\alpha|\gamma}\}^{1/2} \xrightarrow{D} N(0, 1), \quad (4.7)$$

where $\hat{\beta}_0$ and $[T_n(\hat{\beta}_0)]_{\alpha|\gamma} \in \mathbb{R}$ are defined in (2.7). The limiting variance of the decorrelated score function $\sqrt{n} \cdot S_n(\hat{\beta}_0, \lambda)$ is $[I(\beta^*)]_{\alpha|\gamma}$, which is defined in (4.2).

Proof. See §G.2 of the appendix for a detailed proof. \square

Optimality: [27] prove that for inferring α^* in the presence of nuisance parameter γ^* , $[I(\beta^*)]_{\alpha|\gamma}$ is the semiparametric efficient information, i.e., the minimum limiting variance of the (rescaled) score function. Our proposed decorrelated score function achieves such a semiparametric information lower bound and is therefore in this sense optimal.

In the following, we use Gaussian mixture model to illustrate the effectiveness of Theorem 4.6. We defer the details and the implications for mixture of regression to §F of the appendix.

Implications for Gaussian Mixture Model: Under the same model considered in §3.1, if we assume all quantities except $s_{\mathbf{w}}^*$, s^* , d and n are constant, then we have that Conditions 4.1-4.4 hold with $\zeta^{\text{EM}} = s^* \sqrt{\log d \cdot \log n/n}$, $\zeta^{\text{G}} = \sqrt{\log d/n}$, $\zeta^{\text{T}} = \sqrt{\log d/n}$ and $\zeta^{\text{L}} = (\log d + \log n)^{3/2}$. Thus, under Assumption 4.5, (4.7) holds when $n \rightarrow \infty$. Also, we can verify that (4.6) in Assumption 4.5 holds if $\max\{s_{\mathbf{w}}^*, s^*\}^2 \cdot (s^*)^2 \cdot (\log d)^5 = o[n/(\log n)^2]$.

5 Conclusion

We propose a novel high dimensional EM algorithm which naturally incorporates sparsity structure. Our theory shows that, with a suitable initialization, the proposed algorithm converges at a geometric rate and achieves an estimator with the (near-)optimal statistical rate of convergence. Beyond point estimation, we further propose the decorrelated score and Wald statistics for testing hypotheses and constructing confidence intervals for low dimensional components of high dimensional parameters. We apply the proposed algorithmic framework to a broad family of high dimensional latent variable models. For these models, our framework establishes the first computationally feasible approach for optimal parameter estimation and asymptotic inference under high dimensional settings.

References

- [1] ANANDKUMAR, A., GE, R., HSU, D., KAKADE, S. M. and TELGARSKY, M. (2014). Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research* **15** 2773–2832.
- [2] BALAKRISHNAN, S., WAINWRIGHT, M. J. and YU, B. (2014). Statistical guarantees for the EM algorithm: From population to sample-based analysis. *arXiv preprint arXiv:1408.2156*.
- [3] BARTHOLOMEW, D. J., KNOTT, M. and MOUSTAKI, I. (2011). *Latent variable models and factor analysis: A unified approach*, vol. 899. Wiley.
- [4] BELLONI, A., CHEN, D., CHERNOZHUKOV, V. and HANSEN, C. (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* **80** 2369–2429.
- [5] BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics* **37** 1705–1732.

- [6] BOUCHERON, S., LUGOSI, G. and MASSART, P. (2013). *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press.
- [7] CAI, T., LIU, W. and LUO, X. (2011). A constrained ℓ_1 minimization approach to sparse precision matrix estimation. *Journal of the American Statistical Association* **106** 594–607.
- [8] CANDÈS, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when p is much larger than n . *Annals of Statistics* **35** 2313–2351.
- [9] CHAGANTY, A. T. and LIANG, P. (2013). Spectral experts for estimating mixtures of linear regressions. *arXiv preprint arXiv:1306.3729*.
- [10] CHAUDHURI, K., DASGUPTA, S. and VATTANI, A. (2009). Learning mixtures of Gaussians using the k -means algorithm. *arXiv preprint arXiv:0912.0086*.
- [11] DASGUPTA, S. and SCHULMAN, L. (2007). A probabilistic analysis of EM for mixtures of separated, spherical Gaussians. *Journal of Machine Learning Research* **8** 203–226.
- [12] DEMPSTER, A. P., LAIRD, N. M. and RUBIN, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* **39** 1–38.
- [13] JAVANMARD, A. and MONTANARI, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research* **15** 2869–2909.
- [14] KHALILI, A. and CHEN, J. (2007). Variables selection in finite mixture of regression models. *Journal of the American Statistical Association* **102** 1025–1038.
- [15] KNIGHT, K. and FU, W. (2000). Asymptotics for Lasso-type estimators. *Annals of Statistics* **28** 1356–1378.
- [16] LEE, J. D., SUN, D. L., SUN, Y. and TAYLOR, J. E. (2013). Exact inference after model selection via the Lasso. *arXiv preprint arXiv:1311.6238*.
- [17] LOCKHART, R., TAYLOR, J., TIBSHIRANI, R. J. and TIBSHIRANI, R. (2014). A significance test for the Lasso. *Annals of Statistics* **42** 413–468.
- [18] MCLACHLAN, G. and KRISHNAN, T. (2007). *The EM algorithm and extensions*, vol. 382. Wiley.
- [19] MEINSHAUSEN, N. and BÜHLMANN, P. (2010). Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72** 417–473.
- [20] MEINSHAUSEN, N., MEIER, L. and BÜHLMANN, P. (2009). p -values for high-dimensional regression. *Journal of the American Statistical Association* **104** 1671–1681.
- [21] NESTEROV, Y. (2004). *Introductory lectures on convex optimization: A basic course*, vol. 87. Springer.
- [22] NICKL, R. and VAN DE GEER, S. (2013). Confidence sets in sparse regression. *Annals of Statistics* **41** 2852–2876.
- [23] STÄDLER, N., BÜHLMANN, P. and VAN DE GEER, S. (2010). ℓ_1 -penalization for mixture regression models. *TEST* **19** 209–256.
- [24] TAYLOR, J., LOCKHART, R., TIBSHIRANI, R. J. and TIBSHIRANI, R. (2014). Post-selection adaptive inference for least angle regression and the Lasso. *arXiv preprint arXiv:1401.3889*.
- [25] TSENG, P. (2004). An analysis of the EM algorithm and entropy-like proximal point methods. *Mathematics of Operations Research* **29** 27–44.
- [26] VAN DE GEER, S., BÜHLMANN, P., RITOV, Y. and DEZEURE, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Annals of Statistics* **42** 1166–1202.
- [27] VAN DER VAART, A. W. (2000). *Asymptotic statistics*, vol. 3. Cambridge University Press.
- [28] VERSHYNIN, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.
- [29] WASSERMAN, L. and ROEDER, K. (2009). High-dimensional variable selection. *Annals of Statistics* **37** 2178–2201.
- [30] WU, C. F. J. (1983). On the convergence properties of the EM algorithm. *Annals of Statistics* **11** 95–103.
- [31] YI, X., CARAMANIS, C. and SANGHAVI, S. (2013). Alternating minimization for mixed linear regression. *arXiv preprint arXiv:1310.3745*.
- [32] ZHANG, C.-H. and ZHANG, S. S. (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76** 217–242.