

Regularized EM Algorithms: A Unified Framework and Statistical Guarantees

Authored by:

Constantine Caramanis
Xinyang Yi

Abstract

Latent models are a fundamental modeling tool in machine learning applications, but they present significant computational and analytical challenges. The popular EM algorithm and its variants, is a much used algorithmic tool; yet our rigorous understanding of its performance is highly incomplete. Recently, work in [1] has demonstrated that for an important class of problems, EM exhibits linear local convergence. In the high-dimensional setting, however, the M-step may not be well defined. We address precisely this setting through a unified treatment using regularization. While regularization for high-dimensional problems is by now well understood, the iterative EM algorithm requires a careful balancing of making progress towards the solution while identifying the right structure (e.g., sparsity or low-rank). In particular, regularizing the M-step using the state-of-the-art high-dimensional prescriptions (e.g., ‘a la [19]) is not guaranteed to provide this balance. Our algorithm and analysis are linked in a way that reveals the balance between optimization and statistical errors. We specialize our general framework to sparse gaussian mixture models, high-dimensional mixed regression, and regression with missing variables, obtaining statistical guarantees for each of these examples.

1 Paper Body

We give general conditions for the convergence of the EM method for high-dimensional estimation. We specialize these conditions to several problems of interest, including high-dimensional sparse and low-rank mixed regression, sparse gaussian mixture models, and regression with missing covariates. As we explain below, the key problem in the high-dimensional setting is the M -step. A natural idea is to modify this step via appropriate regularization, yet choosing the appropriate sequence of regularizers is a critical problem. As we know from the theory of regularized M-estimators (e.g., [19]) the regularizer should be chosen proportional to the target estimation error. For EM, however, the target

estimation error changes at each step. The main contribution of our work is technical: we show how to perform this iterative regularization. We show that the regularization sequence must be chosen so that it converges to a quantity controlled by the ultimate estimation error. In existing work, the estimation error is given by the relationship between the population and empirical M-step operators, but this too is not well defined in the highdimensional setting. Thus a key step, related both to our algorithm and its convergence analysis, is obtaining a different characterization of statistical error for the high-dimensional setting. Background and Related Work EM (e.g., [8, 12]) is a general algorithmic approach for handling latent variable models (including mixtures), popular largely because it is typically computationally highly scalable, and easy to implement. On the flip side, despite a fairly long history of studying EM in theory (e.g., [12, 17, 21]),

very little has been understood about general statistical guarantees until recently. Very recent work in [1] establishes a general local convergence theorem (i.e., assuming initialization lies in a local region around true parameter) and statistical guarantees for EM, which is then specialized to obtain near-optimal rates for several specific low-dimensional problems. low-dimensional in the sense of the classical statistical setting where the samples outnumber the dimension. A central challenge in extending EM (and as a corollary, the analysis in [1]) to the high-dimensional regime is the M-step. On the algorithm side, the M-step will not be stable (or even well-defined in some cases) in the high-dimensional setting. To make matters worse, any analysis that relies on showing that the finite-sample M-step is somehow close to the M-step performed with infinite data (the population-level M-step) simply cannot apply in the high-dimensional regime. Recent work in [20] treats high-dimensional EM using a truncated M-step. This works in some settings, but also requires specialized treatment for every different setting, precisely because of the difficulty with the M-step. In contrast to work in [20], we pursue a high-dimensional extension via regularization. The central challenge, as mentioned above, is in picking the sequence of regularization coefficients, as this must control the optimization error (related to the special structure of θ^*), as well as the statistical error. Finally, we note that for finite mixture regression, Stadler et al. [16] consider an ℓ_1 regularized EM algorithm for which they develop some asymptotic analysis and oracle inequality. However, this work doesn't establish the theoretical properties of local optima arising from regularized EM. Our work addresses this issue from a local convergence perspective by using a novel choice of regularization.

2

Classical EM and Challenges in High Dimensions

The EM algorithm is an iterative algorithm designed to combat the non-convexity of max likelihood due to latent variables. For space concerns we omit the standard derivation, and only give the definitions we need in the sequel. Let Y, Z be random variables taking values in \mathcal{Y}, \mathcal{Z} , with joint distribution $f(y, z)$ depending on model parameter $\theta \in \mathbb{R}^p$. We observe samples of Y but not of the latent variable Z . EM seeks to maximize a lower bound on the maximum likelihood function for θ . Letting $f(z|y)$ denote the conditional distribution

of Z given $Y = y$, letting $y_{??}(y)$ denote the marginal distribution of Y , and defining the function $n \sum_{i=1}^n \int_0^1 Q_n(\theta - y_i) \log f(\theta | y_i, z) dz$, (2.1) $n \sum_{i=1}^n \int_0^1$ one iteration of the EM algorithm, mapping $\theta(t)$ to $\theta(t+1)$, consists of the following two steps: θ E-step: Compute function $Q_n(\theta - \theta(t))$ given $\theta(t)$. θ M-step: $\theta(t+1) = \arg \max_{\theta} \int_0^1 Q_n(\theta - \theta(t))$. We can define the population (infinite sample) versions of Q_n and M_n in a natural manner: $\int_0^1 Q(\theta - y) \log f(\theta | y, z) dz dy$

$$\begin{aligned} M(\theta) &= \\ &= \int_0^1 \arg \max_{\theta} Q(\theta - y) dy \end{aligned} \quad (2.2) \quad (2.3)$$

This paper is about the high-dimensional setting where the number of samples n may be far less than the dimensionality p of the parameter θ , but where θ exhibits some special structure, e.g., it may be a sparse vector or a low-rank matrix. In such a setting, the M-step of the EM algorithm may be highly problematic. In many settings, for example sparse mixed regression, the M-step may not even be well defined. More generally, when $n \ll p$, $M_n(\theta)$ may be far from the population version, $M(\theta)$, and in particular, the minimum estimation error $\|M_n(\theta) - M(\theta)\|_k$ can be much larger than the signal strength k . This quantity is used in [1] as well as in follow-up work in [20], as a measure of statistical error. In the high dimensional setting, something else is needed.

3 Algorithm

The basis of our algorithm is the by-now well understood concept of regularized high dimensional estimators, where the regularization is tuned to the underlying structure of θ , thus defining a regularized M-step via

$$\begin{aligned} M_{r,n}(\theta) &:= \arg \max_{\theta} \int_0^1 Q_n(\theta - \theta(t)) dy + \lambda R(\theta), \quad (3.1) \\ &\theta(t) \end{aligned}$$

where $R(\theta)$ denotes an appropriate regularizer chosen to match the structure of θ . The key challenge is how to choose the sequence of regularizers $\{\lambda_n\}$ in the iterative process, so as to control optimization and statistical error. As detailed in Algorithm 1, our sequence of regularizers attempts to match the target estimation error at each step of the EM iteration. For an intuition of what this might look like, consider the estimation error at step t : $\|M_{r,n}(\theta(t)) - M(\theta)\|_k$. By the triangle inequality, we can bound this by a sum of two terms: the optimization error and the final estimation error: $\|M_{r,n}(\theta(t)) - M(\theta)\|_k \leq \|M_{r,n}(\theta(t)) - M_n(\theta)\|_k + \|M_n(\theta) - M(\theta)\|_k$. (3.2) (t)

Since we expect (and show) linear convergence of the optimization, it is natural to update λ_n via a $(t+1)$ recursion of the form $\lambda_n = \eta \lambda_{n-1} + \epsilon$ as in (3.3), where the first term represents the optimization error, and ϵ represents the final statistical error, i.e., the last term above in (3.2). A key part of our analysis shows that this error (and hence ϵ) is controlled by $k^2 Q_n(\theta - \theta(t))$, which in turn can be bounded uniformly for a variety

of important applications of EM, including the three discussed in this paper (see Section 5). While a technical point, it is this key insight that enables the right choice of algorithm and its analysis. In the cases we consider, we obtain min-max optimal rates of convergence, demonstrating that no algorithm, let alone another variant of EM, can perform better. Algorithm 1 Regularized EM

Algorithm Input Samples $\{y_i\}_{i=1}^n$, regularizer R , number of iterations T , initial parameter $\theta(0)$, initial regularization parameter λ_n , estimated statistical error ϵ , contractive factor $\beta \leq 1$.

- 1: For $t = 1, 2, \dots, T$ do
- 2: Regularization parameter update: $\lambda(t) = \beta \lambda(t-1) + \epsilon/n$

3: 4:

E-step: Compute function $Q_n(\theta | \lambda(t))$ Regularized M-step:

$\theta(t)$

(3.3)

) according to (2.1).

$\theta(t) = \text{Mrn}(\lambda(t)) := \arg \max_{\theta} Q_n(\theta | \lambda(t)) - \lambda(t) R(\theta)$.

5: End For Output $\theta(T)$.

4

Statistical Guarantees

We now turn to the theoretical analysis of regularized EM algorithm. We first set up a general analytical framework for regularized EM where the key ingredients are decomposable regularizer and several technical conditions on the population based $Q(\theta | \lambda)$ and the sample based $Q_n(\theta | \lambda)$. In Section 4.3, we provide our main result (Theorem 1) that characterizes both computational and statistical performance of the proposed variant of regularized EM algorithm.

Decomposable Regularizers

Decomposable regularizers (e.g., [3, 6, 14, 19]), have been shown to be useful both empirically and theoretically for high dimensional structural estimation, and they also play an important role in our analytical framework. Recall that for $R : \mathbb{R}^p \rightarrow \mathbb{R}_+$ a norm, and a pair of subspaces (S, S^\perp) in \mathbb{R}^p such that $S \perp S^\perp$, we have the following definition: Definition 1 (Decomposability). Regularizer $R : \mathbb{R}^p \rightarrow \mathbb{R}_+$ is decomposable with respect to (S, S^\perp) if $R(u + v) = R(u) + R(v)$, for any $u \in S, v \in S^\perp$. Typically, the structure of model parameter θ can be characterized by specifying a subspace S such that $\theta \in S$. The common use of a regularizer is thus to penalize the compositions of solution that

live outside S . We are interested in bounding the estimation error in some norm $\|\cdot\|_k$. The following quantity is critical in connecting R to $\|\cdot\|_k$. Definition 2 (Subspace Compatibility Constant). For any subspace $S \subset \mathbb{R}^p$, a given regularizer R and some norm $\|\cdot\|_k$, the subspace compatibility constant of S with respect to $R, \|\cdot\|_k$ is given by $R(u) \wedge(S) := \sup_{u \in S \setminus \{0\}} \frac{R(u)}{\|u\|_k}$

As is standard, the dual norm of R is defined as $R^*(v) := \sup_{u \in \mathbb{R}^p} \langle v, u \rangle - R(u)$. To simplify notation, we let $\|u\|_{kR} := R(u)$ and $\|u\|_{kR^*} := R^*(u)$.

Conditions on $Q(\theta | \lambda)$ and $Q_n(\theta | \lambda)$

Next, we review three technical conditions, originally proposed by [1], on the population level $Q(\theta | \lambda)$ function, and then we give two important conditions that the empirical function $Q_n(\theta | \lambda)$ must satisfy, including one that characterizes the statistical error. It is well known that performance of EM algorithm

is sensitive to initialization. Following the lowdimensional development in [1], our results

are local, and apply to an r -neighborhood region around $\theta^* : B(r; \theta^*) := \{u \in \mathbb{R}^p : \|u - \theta^*\| \leq r\}$. We first require that $Q(\cdot - \theta^*)$ is self consistent as stated below. This is satisfied, in particular, when θ^* maximizes the population log likelihood function, as happens in most settings of interest [12]. Condition 1 (Self Consistency). Function $Q(\cdot - \theta^*)$ is self consistent, namely $\theta^* = \arg \max_{\theta} Q(\cdot - \theta)$.

We also require that the function $Q(\cdot - \theta^*)$ satisfies a certain strong concavity condition and is smooth over \mathbb{R}^p . Condition 2 (Strong Concavity and Smoothness (κ, ℓ, r)). $Q(\cdot - \theta^*)$ is κ -strongly concave over \mathbb{R}^p , i.e.,

$$(4.1) \quad Q(\theta_2 - \theta^*) - Q(\theta_1 - \theta^*) \leq \kappa \|\theta_1 - \theta_2\|, \quad \theta_1, \theta_2 \in B(r; \theta^*),$$

and $Q(\cdot - \theta^*)$ is ℓ -smooth over \mathbb{R}^p , i.e.,

$$(4.2) \quad Q(\theta_2 - \theta^*) - Q(\theta_1 - \theta^*) \leq \ell \|\theta_1 - \theta_2\|, \quad \theta_1, \theta_2 \in B(r; \theta^*).$$

The next condition is key in guaranteeing the curvature of $Q(\cdot - \theta^*)$ is similar to that of $Q(\cdot - \theta)$ when θ is close to θ^* . It has also been called First Order Stability in [1]. Condition 3 (Gradient Stability (κ, r)). For any $\theta \in B(r; \theta^*)$, we have

$\nabla Q(M(\theta) - \theta^*) = \nabla Q(M(\theta) - \theta)$ for $\|\theta - \theta^*\| \leq r$. The above condition only requires that the gradient be stable at one point $M(\theta)$. This is sufficient for our analysis. In fact, for many concrete examples, one can verify a stronger version of Condition

3 that is $\nabla Q(\theta - \theta^*) = \nabla Q(\theta - \theta)$ for $\|\theta - \theta^*\| \leq r$. Next we require two conditions on the empirical function $Q_n(\cdot - \theta^*)$, which is computed from finite number of samples according to (2.1). Our first condition, parallel to Condition 2, imposes a curvature constraint on $Q_n(\cdot - \theta^*)$. In order to guarantee that the estimation error $\|\theta_n(t) - \theta^*\|$ in step t of the EM algorithm is well controlled, we would like $Q_n(\cdot - \theta^*)$ to be strongly concave at θ^* . However, in the setting where $n \leq p$, there might exist directions along which $Q_n(\cdot - \theta^*)$ is flat, e.g., as in mixed linear regression and missing covariate regression. In contrast with Condition 2, we only require $Q_n(\cdot - \theta^*)$ to be strongly concave over a particular set $C(S, S; R)$ that is defined in terms of the subspace pair (S, S) and regularizer R . This set is defined as follows:

$$C(S, S; R) := \{u \in \mathbb{R}^p : \nabla S(u) \in S, \|\nabla S(u) - \nabla S(u)\| \leq R\}, \quad (4.3)$$
 where the projection operator $\nabla S : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is defined as $\nabla S(u) := \arg \min_{v \in \mathbb{R}^p} \|v - u\|_S$. The restricted strong concavity (RSC) condition is as follows.

Condition 4 (RSC (κ, S, S, r, T)). For any fixed $\theta \in B(r; \theta^*)$, with probability at least $1 - \delta$, we have that for all $\theta \in C(S, S; R)$,

$$\kappa_n Q_n(\theta - \theta^*) - Q_n(\theta - \theta^*) \leq \kappa_n Q_n(\theta - \theta^*), \quad \|\theta - \theta^*\| \leq r.$$

The above condition states that $Q_n(\cdot - \theta^*)$ is strongly concave in directions $\theta - \theta^*$ that belong to $C(S, S; R)$. It is instructive to compare Condition 4 with a related condition proposed by [14] for analyzing high dimensional M-estimators. They require the loss function to be strongly convex over the cone $\{u \in \mathbb{R}^p : \|\nabla S(u)\|_R \leq \kappa \|\nabla S(u)\|_S\}$. Therefore our

restrictive set (4.3) is similar to the cone but has the additional term $2\eta(S)\mathbf{k}$. The main purpose of the term $2\eta(S)\mathbf{k}$ is to allow the regularization parameter η to jointly control optimization and statistical error. We note that while Condition 4 is stronger than the usual

RSC condition in M-estimator, in typical settings

the difference is immaterial. This is because $\eta(S) \|\mathbf{u}\|_R$ is within a constant factor of $\eta(S) \|\mathbf{u}\|$, and hence checking RSC over \mathcal{C} amounts to checking it over $\mathbf{k} \in S \cap \{\mathbf{u} \mid \|\mathbf{u}\|_R \leq \eta(S)\mathbf{k}\}$, which is indeed what is typically also done in the M-estimator setting. Finally, we establish the condition that characterizes the achievable statistical error. Condition 5 (Statistical Error (η, r, γ)). For any fixed $\eta \in B(r; \gamma)$, with probability at least $1 - \gamma$, we have

$\eta Q_n(\eta \cdot) \leq \eta Q(\eta \cdot) + \eta \gamma$. (4.4) R This quantity replaces the term $\mathbf{k} M_n(\eta) \mathbf{k}$ which appears in [1] and [20], and which presents problems in the high dimensional regime. 4.3

Main Results

In this section, we provide the theoretical guarantees for a resampled version of our regularized EM algorithm: we split the whole dataset into T pieces and use a fresh piece of data in each iteration of regularized EM. As in [1], resampling makes it possible to check that Conditions 4-5 are satisfied without requiring them to hold uniformly for all $\eta \in B(r; \gamma)$ with high probability. Our empirical results indicate that it is not in fact required and is an artifact of the analysis. We refer to this resampled version as Algorithm 2. In the sequel, we let $m := n/T$ to denote the sample complexity in each iteration. We let $\eta := \sup_{\mathbf{p} \in \mathcal{R}_p} \{\mathbf{k} \mid \eta(\mathbf{p}) \leq \mathbf{k}\}$, where $\mathbf{k} \in \mathcal{R}^*$ is the dual norm of $\mathbf{k} \in \mathcal{R}$. For Algorithm 2, our main result is as follows. The proof is deferred to the Supplemental Material. Theorem 1. Assume the model parameter $\mathbf{p} \in S$ and regularizer R is decomposable with respect to (S, S) where $S \subseteq \mathcal{R}_p$. Assume $r > 0$ is such that $B(r; \gamma) \neq \emptyset$. Further, assume function $Q(\cdot)$, defined in (2.2), is self consistent and satisfies Conditions 2-3 with parameters (η, γ, r) and (η, r) . Given n samples and T iterations, let $m := n/T$. Assume $Q_m(\cdot)$, computed from any m i.i.d. samples according to (2.1), satisfies Conditions 4-5 with parameters $(\eta_m, S, S, r, 0.5\gamma/T)$ and $(\eta_m, r, 0.5\gamma/T)$. Let $\gamma := 5\gamma$, and assume $0 < \gamma < 1$ and $0 < \gamma < 3/4$. Define $m_\gamma := r\eta_m / [60\eta(S)]$ and assume η_m is such that $\eta_m \leq \gamma$. Consider Algorithm 2 with initialization $\mathbf{p}^{(0)} \in B(r; \gamma)$ and with regularization parameters given by $\eta_t = \eta_m \mathbf{k}^{(t)}$ $\mathbf{k}^{(0)} = \mathbf{k} + \gamma$, $t = 1, 2, \dots, T$ (4.5) $m = \eta^{-1} \eta(S)$ for any $\eta \in [3\eta_m, 3\gamma]$, $\gamma \in [\gamma, 3/4]$. Then with probability at least $1 - \gamma$, we have that for any $t \in [T]$, $5\eta \leq \eta_t \leq \eta(S)$. (4.6) $\mathbf{k}^{(t)} \leq \mathbf{k} + \eta_t \mathbf{k}^{(0)} \leq \mathbf{k} + \eta_m \mathbf{1}$ γ The estimation error is bounded by a term decaying linearly with number of iterations t , which we can think of as the optimization error and a second term that characterizes the ultimate estimation error of our algorithm. With $T = O(\log n)$ and suitable choice of γ such that $\gamma = O(\eta n/T)$, we bound the ultimate estimation error as $\mathbf{k}^{(T)} \leq \mathbf{k} + \eta(S)\eta n/T$. (4.7) $(1 - \gamma)\eta n/T \leq$

We note that overestimating the initial error, $\mathbf{k}^{(0)} \leq \mathbf{k}$ is not important, as it may slightly increase the overall number of iterations, but will not impact the ultimate estimation error. The constraint $\eta_m \leq r\eta_m / \eta(S)$ ensures that η

(t) is contained in $B(r; \frac{1}{2} \epsilon)$ for all $t \in [T]$. This constraint is quite mild in the sense that if $\hat{m} = \hat{r}(r; m / \epsilon(S))$, $\hat{r}(0)$ is a decent estimator with estimation error $O(\epsilon(S) \hat{m} / \epsilon)$ that already matches our expectation.

5

Examples: Applying the Theory

Now we introduce three well known latent variable models. For each model, we first review the standard EM algorithm formulations, and discuss the extensions to the high dimensional setting. Then we apply Theorem 1 to obtain the statistical guarantee of the regularized EM with data splitting (Algorithm 2). The key ingredient underlying these results is to check the technical conditions in Section 4 hold for each model. We postpone these tedious details to the Supplemental Material. 5.1

Gaussian Mixture Model

We consider the balanced isotropic Gaussian mixture model (GMM) with two components where the distribution of random variables $(Y, Z) \in \mathbb{R}^p \times \{1, 1\}$ is characterized as $\Pr(Y = y, Z = z) = \frac{1}{2} \phi(y; z, \Sigma, \mu)$, $\Pr(Z = 1) = \Pr(Z = -1) = 1/2$. Here we use $\phi(\cdot; \cdot, \cdot, \cdot)$ to denote the probability density function of $N(\cdot, \Sigma)$. In this example, Z is the latent variable that indicates the cluster id of each sample. Given n i.i.d. samples $\{y_i\}_{i=1}^n$, function $Q_n(\cdot)$ defined in (2.1) corresponds to

$$M_{\text{GMM}}(\cdot) = \frac{1}{n} \sum_{i=1}^n Q(y_i; \cdot)$$

$$Q(y; \cdot) = \frac{1}{2} \exp(-\frac{1}{2} \|y - \mu_1\|^2) + \frac{1}{2} \exp(-\frac{1}{2} \|y - \mu_2\|^2), \quad \mu_1, \mu_2 \in \mathbb{R}^p$$

(5.1)

$$Q(y; \cdot) = \frac{1}{2} \exp(-\frac{1}{2} \|y - \mu_1\|^2) + \frac{1}{2} \exp(-\frac{1}{2} \|y - \mu_2\|^2)$$

where $w(y; \cdot) := \exp(-\frac{1}{2} \|y - \mu_1\|^2) / [\exp(-\frac{1}{2} \|y - \mu_1\|^2) + \exp(-\frac{1}{2} \|y - \mu_2\|^2)]$. We assume $\mu_1, \mu_2 \in B_0(s; p) := \{u \in \mathbb{R}^p : \text{supp}(u) \subseteq S\}$. Naturally, we choose the regularizer $R(\cdot)$ to be the ‘1 norm. We define the signal-to-noise ratio $\text{SNR} := k^2 / \sigma^2$. Corollary 1 (Sparse Recovery in GMM). There exist constants ϵ, C such that if $\text{SNR} \geq \epsilon$, $n/T \geq C(k^2 / \sigma^2 + \epsilon) / (k^2 / \sigma^2) \log p$, $\hat{r}(0) \in B(k^2 / \sigma^2 / 4; \epsilon)$; then with probability at least $1 - \epsilon / p$, Algorithm 2 with parameters $\epsilon = C(k^2 / \sigma^2 + \epsilon) T \log p / n$, $\hat{r}_n / T = 0.2k^2 / \sigma^2$, any $\epsilon \in [1/2, 3/4]$ and ‘1 regularization generates $\hat{r}(t)$ that has estimation error $\leq C(k^2 / \sigma^2 + \epsilon) \log p / (t \epsilon)$, for all $t \in [T]$. 1?? n

(5.2)

Note that by setting $T \log(n / \log p)$, the order of final estimation error turns out to be $(k^2 / \sigma^2 + \epsilon) (s \log p) / n$. The minimax rate for estimating s -sparse vector in a single Gaussian cluster is $s \log p / n$, thereby the rate is optimal on (n, p, s) up to a log factor. 5.2

Mixed Linear Regression

Mixed linear regression (MLR), as considered in some recent work [5, 7, 22], is the problem of recovering two or more linear vectors from mixed linear measurements. In the case of mixed linear regression with two symmetric and balanced components, the response-covariate pair $(Y, X) \in \mathbb{R}^n \times \mathbb{R}^{n \times p}$ is linked through $Y = HX, Z \in \mathbb{R}^n$, where W is the noise term and Z is the latent

variable that has Rademacher distribution over $\{-1, 1\}$. We assume $X \sim N(0, I_p)$, $W \sim N(0, \frac{1}{2}I_2)$. In this setting, with n i.i.d. samples $\{y_i, x_i\}_{i=1}^n$ of pair (Y, X) , function $Q_n(\cdot, \cdot)$ then corresponds to n

$$\frac{1}{2n} \sum_{i=1}^n Q_M(y_i, x_i; \theta) = \frac{1}{2n} \sum_{i=1}^n \left[w(y_i, x_i; \theta) (y_i - h(x_i, \theta))^2 + (1 - w(y_i, x_i; \theta)) (y_i + h(x_i, \theta))^2 \right], \quad (5.3)$$

$$w(y, x; \theta) := \exp(-\theta y h(x, \theta)) [\exp(-\theta y h(x, \theta)) + \exp(\theta (y + h(x, \theta)))]^{-1}.$$

We consider two kinds of structure on θ : Sparse Recovery. Assume $\theta \in B_0(s, p)$. Then let R be the ℓ_1 norm, as in the previous section. We define $\text{SNR} := k^2 / \|\theta\|_1$. Corollary 2 (Sparse recovery in MLR). There exist constant C, C_0 such that if $\text{SNR} \geq C$, $n/T \geq 2C_0 [(k^2 + \|\theta\|_1^2) / k^2] s \log p$, $\theta(0) \in B(k^2 / 240, \|\theta\|_1)$; then with probability at least $1 - T^{-p} / p^p(0)$ Algorithm 2 with parameters $\gamma = C(k^2 + \|\theta\|_1^2) T \log p / n$, $\eta n / T = k^2(0) / (15s)$, any $\gamma \in [1/2, 3/4]$ and ℓ_1 regularization generates $\hat{\theta}(t)$ that has estimation error $\|\hat{\theta}(t) - \theta\|_1 \leq 15C(k^2 + \|\theta\|_1^2) s \log p(t) / k^2(0) / k^2 + T$, for all $t \in [T]$. Performing T

$\log(n/(s \log p))$ iterations gives us estimation rate $(k^2 + p^2) (s \log p / n) \log(n/(s \log p))$ which is near-optimal on (s, p, n) . The dependence on $k^2 + p^2$, which also appears in the analysis of EM in the classical (low dimensional) setting [1], arises from fundamental limits of EM. Removing such dependence for MLR is possible by convex relaxation [7]. It is interesting to study how to remove it in the high dimensional setting. Low Rank Recovery. Second we consider the setting where the model parameter is a matrix $\theta \in \mathbb{R}^{p_1 \times p_2}$ with $\text{rank}(\theta) \leq \min(p_1, p_2)$. We further assume $X \in \mathbb{R}^{p_1 \times p_2}$ is an i.i.d. Gaussian matrix, i.e., entries of X are independent random variables with distribution 1). We apply PPN1, $p(0, 2)$ nuclear norm regularization to serve the low rank structure, i.e., $R(\theta) = \sum_{i=1}^{\min(p_1, p_2)} \sigma_i(\theta)$, where $\sigma_i(\theta)$ is the i th singular value of θ . Similarly, we let $\text{SNR} := k^2 / \|\theta\|_F$. Corollary 3 (Low rank recovery in MLR). There exist constant C, C_0 such that if $\text{SNR} \geq C$, $2n/T \geq C_0 [(k^2 + \|\theta\|_F^2) / k^2] (p_1 + p_2)$, $\theta(0) \in B(k^2 / 1600, \|\theta\|_F)$; then with probability at least $1 - T^{-\exp(p_1 + p_2)}$ Algorithm 2 with parameters $\gamma = C(k^2 + \|\theta\|_F^2) T (p_1 + p_2) / n$, $\eta(0) n / T = 0.01k^2(0) / \|\theta\|_F^2$, any $\gamma \in [1/2, 3/4]$ and nuclear norm regularization generates $\hat{\theta}(t)$ that has estimation error $\|\hat{\theta}(t) - \theta\|_F \leq k^2(0) / k^2 + T$, for all $t \in [T]$. n

The standard low rank matrix recovery with a single component, including other sensing matrix designs beyond the Gaussianity, has been studied extensively (e.g., [2, 4, 13, 15]). To the best of our knowledge, the theoretical study of the mixed low rank matrix recovery has not been considered. 5.3

Missing Covariate Regression

As our last example, we consider the missing covariate regression (MCR) problem. To parallel standard linear regression, $\{y_i, x_i\}_{i=1}^n$ are samples of (Y, X) linked through $Y = hX, \epsilon_i + W$. However, we assume each entry of x_i is missing independently with probability $\epsilon \in (0, 1)$. Thereby takes the form

fore, the observed covariate vector x_i with probability $1 - \epsilon$ $x_{i,j} = \epsilon$ otherwise. We assume the model is under Gaussian design $X \sim N(0, I_p)$, $W \sim N(0, \sigma^2)$. We refer the reader to our Supplementary Material for the specific $Q_n(\epsilon)$ function. In high dimensional case, we assume $\epsilon \leq B_0(s; p)$. We define $\epsilon := k^2 / \epsilon$ to be the SNR and $\epsilon := r/k^2$ to be the relative contractivity radius. In particular, let $\epsilon := (1 + \epsilon)$. Corollary 4 (Sparse Recovery in MCR). There exist constants C, C_0, C_1 such that if $(1 + \epsilon) \leq C_0$, $\epsilon \leq C_1$, $n/T \geq C_0 \max\{\epsilon^2, 1\} s \log p$, $\epsilon \in (0, 1)$ $B(k^2, \epsilon)$; then with probp (0) ability at least $1 - T/p$ Algorithm 2 with parameters $\epsilon = C T \log p/n$, $\epsilon_n/T = k^2 (0) \epsilon \leq k^2/(45 s)$, any $\epsilon \in [1/2, 3/4]$ and ‘1 regularization generates $\epsilon(t)$ that has estimation error $r \leq 45 C \epsilon s \log p(t) \epsilon(t) \leq k^2 \epsilon \leq k^2 + T$, for all $t \in [T]$, $1 \leq n \leq 7$

Unlike the previous two models, we require an upper bound on the signal to noise ratio. This unusual constraint p is in fact unavoidable [10]. By optimizing T , the order of final estimation error turns out to be $\epsilon s \log p/n \log(n/(s \log p))$.

6

Simulations

We now provide some simulation results to back up our theory. Note that while Theorem 1 requires resampling, we believe in practice this is unnecessary. This is validated by our results, where we apply Algorithm 1 to the four latent variable models discussed in Section 5. Convergence Rate. We first evaluate the convergence of Algorithm 1 assuming only that the initialization is a bounded distance from ϵ . For a given error $\epsilon_k \leq k^2$, the initial parameter $\epsilon(0)$ is picked randomly from the sphere centered around ϵ with radius $\epsilon_k \leq k^2$. We use Algorithm 1 with $T = 7$, $\epsilon(0) = 0.7$, ϵ_n in Theorem 1. The choice of the critical parameter ϵ is given in the Supplementary Material. For every single trial, we report estimation error $k^2(t) \leq k^2$ and optimization error $k^2(t) \leq (T) k^2$ in every iteration. We plot the log of errors over iteration t in Figure 1. 2 Est error Opt error

Log error
-1
Log error
1 Est error Opt error
0
-2 -3

3 Est error Opt error
 0
 -2
 Log error
 0
 -4
 -1 -2
 -6
 -4 -5
 -8
 -6
 -10
 0
 1
 2
 3
 4
 5
 6
 7
 Est error Opt error
 2
 Log error
 1
 1 0 -1
 -3
 0
 1
 2
 Number of iterations
 3
 4
 5
 6
 -4
 7
 -2 0
 1
 2
 Number of iterations
 (a) GMM
 3
 4
 5
 6
 -3

7
0
1
2
Number of iterations
(b) MLR(sparse)
3
4
5
6
7
Number of iterations
(c) MLR(low rank)
(d) MCR

Figure 1: Convergence of regularized EM algorithm. In each panel, one curve is plotted from single independent trial. Settings: (a,b,d) $(n, p, s) = (500, 800, 5)$; (d) $(n, p, ?) = (600, 30, 3)$; (a-c) $\text{SNR} = 5$; (d) $(\text{SNR},) = (0.5, 0.2)$; (a-d) $? = 0.5$. Statistical Rate. We now evaluate the statistical rate. We set $T = 7$ and compute estimation error on $?b := ?(T)$. In Figure 2, we plot $k?b ? ? ? k^2$ over normalized sample complexity, i.e., $n/(s \log p)$ for s-sparse parameter and $n/(?p)$ for rank $? p$ -by- p parameter. We refer the reader to Figure 1 for other settings. We observe that the same normalized sample complexity leads to almost identical estimation error in practice, which thus supports the corresponding statistical rate established in Section 5. $0.4 p = 200 p = 400 p = 800$

1.4 $p = 200 p = 400 p = 800$
0.18 0.16 0.14
0.3 0.25
1
0.12
0.6
0.1
0.15
0.4
10
15
20
25
30
 $n/(s \log p)$
(a) GMM
5
10
15
20
25

30
 1.4
 1 4
 5
 6
 7
 $n/(\epsilon p)$
 (b) MLR(sparse)
 1.6
 1.2
 3
 $n/(s \log p)$
 $p = 200 \quad p = 400 \quad p = 800$
 1.8
 0.8
 0.2
 5
 $2 \quad p = 25 \quad p = 30 \quad p = 35$
 1.2
 $\epsilon \quad \epsilon \quad \epsilon \quad k_F \quad k?$
 0.35
 $k?? \quad \epsilon \quad \epsilon \quad k_2$
 $k?? \quad \epsilon \quad \epsilon \quad k_2$
 0.2
 $k?? \quad \epsilon \quad \epsilon \quad k_2$
 0.22
 (c) MLR(low rank)
 8
 5
 10
 15
 20
 25
 30
 $n/(s \log p)$
 (d) MCR

Figure 2: Statistical rates. Each point is an average of 20 independent trials. Settings: (a,b,d) $s = 5$; (c) $\epsilon = 3$.

Acknowledgments The authors would like to acknowledge NSF grants 1056028, 1302435 and 1116955. This research was also partially supported by the U.S. Department of Transportation through the Data-Supported Transportation Operations and Planning (D-STOP) Tier 1 University Transportation Center.

2 References

- [1] Sivaraman Balakrishnan, Martin J Wainwright, and Bin Yu. Statistical guarantees for the EM algorithm: From population to sample-based analysis. arXiv preprint arXiv:1408.2156, 2014.
- [2] T Tony Cai and Anru Zhang. Rob: Matrix recovery via rank-one projections. *The Annals of Statistics*, 43(1):102?138, 2015.
- [3] Emmanuel Candes and Terence Tao. The Dantzig selector: statistical estimation when p is much larger than n . *The Annals of Statistics*, pages 2313?2351, 2007.
- [4] Emmanuel J Candès and Yaniv Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *Information Theory, IEEE Transactions on*, 57(4):2342? 2359, 2011.
- [5] Arun Tejasvi Chaganty and Percy Liang. Spectral experts for estimating mixtures of linear regressions. arXiv preprint arXiv:1306.3729, 2013.
- [6] Yudong Chen, Sujay Sanghavi, and Huan Xu. Improved graph clustering. *Information Theory, IEEE Transactions on*, 60(10):6440?6455, Oct 2014.
- [7] Yudong Chen, Xinyang Yi, and Constantine Caramanis. A convex formulation for mixed regression with two components: Minimax optimal rates. In *Conf. on Learning Theory*, 2014.
- [8] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1?38, 1977.
- [9] Po-Ling Loh and Martin J Wainwright. High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. In *Advances in Neural Information Processing Systems*, pages 2726? 2734, 2011.
- [10] Po-Ling Loh and Martin J Wainwright. Corrupted and missing predictors: Minimax bounds for highdimensional linear regression. In *Information Theory Proceedings (ISIT)*, 2012 IEEE International Symposium on, pages 2601?2605. IEEE, 2012.
- [11] Jinwen Ma and Lei Xu. Asymptotic convergence properties of the em algorithm with respect to the overlap in the mixture. *Neurocomputing*, 68:105?129, 2005.
- [12] Geoffrey McLachlan and Thiriyambakam Krishnan. *The EM algorithm and extensions*, volume 382. John Wiley & Sons, 2007.
- [13] Sahand Negahban, Martin J Wainwright, et al. Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, 39(2):1069?1097, 2011.
- [14] Sahand Negahban, Bin Yu, Martin J Wainwright, and Pradeep K Ravikumar. A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers. In *Advances in Neural Information Processing Systems*, pages 1348?1356, 2009.
- [15] Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 52(3):471?501, 2010.
- [16] Nicolas St?adler, Peter B?uhlmann, and Sara Van De Geer. L1-penalization for mixture regression models. *Test*, 19(2):209?256, 2010.
- [17] Paul Tseng. An analysis of the em algorithm and entropy-like proximal point methods. *Mathematics of Operations Research*, 29(1):27?44, 2004.
- [18] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. arXiv preprint arXiv:1011.3027, 2010.
- [19] Martin J Wainwright. Structured regularizers for high-dimensional problems: Statistical and computational issues. *Annual Review of Statistics and Its Application*, 1:233?253, 2014.
- [20] Zhaoran Wang, Quanquan Gu, Yang Ning,

and Han Liu. High dimensional expectation-maximization algorithm: Statistical optimization and asymptotic normality. arXiv preprint arXiv:1412.8729, 2014. [21] C.F.Jeff Wu. On the convergence properties of the em algorithm. The Annals of statistics, pages 95?103, 1983. [22] Xinyang Yi, Constantine Caramanis, and Sujay Sanghavi. Alternating minimization for mixed linear regression. arXiv preprint arXiv:1310.3745, 2013.