

---

# Mean Estimation in High-Dimensional Binary Markov Gaussian Mixture Models

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We consider a high-dimensional mean estimation problem over a binary hidden  
2 Markov model, which illuminates the interplay between memory in data, sample  
3 size, dimension, and signal strength in statistical inference. In this model, an  
4 estimator observes  $n$  samples of a  $d$ -dimensional parameter vector  $\theta_* \in \mathbb{R}^d$ ,  
5 multiplied by a random sign  $S_i$  ( $1 \leq i \leq n$ ), and corrupted by isotropic standard  
6 Gaussian noise. The sequence of signs  $\{S_i\}_{i \in [n]} \in \{-1, 1\}^n$  is drawn from a  
7 stationary homogeneous Markov chain with flip probability  $\delta \in [0, 1/2]$ . As  $\delta$   
8 varies, this model smoothly interpolates two well-studied models: the Gaussian  
9 Location Model for which  $\delta = 0$  and the Gaussian Mixture Model for which  $\delta =$   
10  $1/2$ . Assuming that the estimator knows  $\delta$ , we establish a nearly minimax optimal  
11 (up to logarithmic factors) estimation error rate, as a function of  $\|\theta_*\|$ ,  $\delta$ ,  $d$ ,  $n$ . We  
12 then provide an upper bound to the case of estimating  $\delta$ , assuming a (possibly  
13 inaccurate) knowledge of  $\theta_*$ . The bound is proved to be tight when  $\theta_*$  is an  
14 accurately known constant. These results are then combined to an algorithm which  
15 estimates  $\theta_*$  with  $\delta$  unknown a priori, and theoretical guarantees on its error are  
16 stated.

## 17 1 Introduction

18 To what extent does memory between samples from a high dimensional model affect the possible  
19 estimation error? In general, this estimation error depends in an intricate way on the interplay  
20 between the number of samples, the dimension of the vector parameters to be estimated, the noise  
21 level (signal-to-noise ratio), and the amount of memory between the samples. In this paper, we study  
22 this question in the context of an elementary hidden Markov high-dimensional Gaussian model, with  
23 the goal of establishing a tight characterization of the estimation error in terms of this trade-off. We  
24 next turn to formally define this model and the estimation problem, describe known results, and then  
25 present our contributions.

### 26 1.1 Problem formulation

27 Let  $S_0^n := (S_0, S_1, \dots, S_n)$  be the following homogeneous binary symmetric Markov chain,  $S_i \in$   
28  $\{-1, 1\}$ ,  $\mathbb{P}[S_0 = 1] = 1/2$  and

$$S_i = \begin{cases} S_{i-1}, & \text{w.p. } 1 - \delta \\ -S_{i-1}, & \text{w.p. } \delta \end{cases} \quad (1)$$

29 for  $i \in [n] := \{1, \dots, n\}$ , and where  $\delta \in [0, 1]$  is the flip probability of the binary Markov chain. We  
30 also denote  $\rho := 1 - 2\delta \in [-1, 1]$  which is the correlation between adjacent samples  $\rho = \mathbb{E}[S_i S_{i+1}]$ .  
31 At each time point  $i \in [n]$ , a sample of a  $d$ -dimensional Gaussian mixture model is observed

$$X_i = S_i \theta_* + Z_i, \quad (2)$$

32 where  $Z_i \sim N(0, I_d)$  is an i.i.d. sequence, independent of  $S_0^n$ , and where  $\theta_* \in \mathbb{R}^d$ ,  $d \geq 1$ . At its  
 33 two extremes, this model degenerates to one of two fundamental models. When  $\delta = 0$ , the memory  
 34 length is infinite, and the sign  $S_0 = S_1 \cdots = \cdots = S_n$  is fixed. Thus, up to this sign ambiguity, the  
 35 model (2) is the standard *Gaussian location model* (GLM), which is essentially a memoryless model  
 36 (and exactly so if  $S_0$  is known). When  $\delta = \frac{1}{2}$ , the signs  $S_0^n$  are i.i.d. and have no memory at all. The  
 37 model (2) is then a *Gaussian mixture model* (GMM) with two symmetric components, which is also  
 38 a memoryless model. In all other cases,  $0 < \delta < \frac{1}{2}$  (or  $\frac{1}{2} < \delta < 1$ ), the model is a simple version of  
 39 a *hidden Markov model* (HMM).

40 The inference problem we consider in this paper is the estimation of  $\theta_* \in \mathbb{R}^d$ , under the loss function

$$\text{loss}(\hat{\theta}, \theta_*) := \min\{\|\hat{\theta} - \theta_*\|, \|\hat{\theta} + \theta_*\|\}, \quad (3)$$

41 that is, the Euclidean distance error under a possible sign ambiguity.<sup>1</sup> An intermediate goal (or an  
 42 additional problem) is to estimate  $\delta$ , under the regular absolute error loss function  $|\hat{\delta} - \delta|$ .

43 The fundamental limits of this estimation problem will be gauged by the *local* minimax rate, which is  
 44 the maximal decrease rate of the loss possible for any estimator, given  $n$  samples, at dimension  $d \geq 2$ ,  
 45 for signal strength  $\|\theta_*\| = t$ , and under flip probability  $\delta$ . Specifically, for  $d \geq 2$  it is defined as

$$M(n, d, \delta, t) := \inf_{\hat{\theta}(X_1^n)} \sup_{\|\theta_*\|=t} \mathbb{E} \left[ \text{loss}(\theta_*, \hat{\theta}(X_1^n)) \right]. \quad (4)$$

46 For general  $d \geq 1$ , the *global* minimax rate is defined with the condition  $\|\theta_*\| = t$  replaced by  
 47  $\|\theta_*\| \leq t$  (this condition trivializes the estimator for  $d = 1$ ).

## 48 1.2 Known results for GLM and GMM

49 For the high dimensional Gaussian models we consider in this paper, there are two possible phase  
 50 transitions – one as  $t$  increases, and the other one as  $d$  increases. Typically, the regime of main  
 51 theoretical interest is the low-separation regime ( $t \lesssim 1$ ), in which separating with high accuracy  
 52 between the components is impossible (in general), yet parameter estimation with vanishing loss is  
 53 still possible.

54 At the first extreme, the local minimax rate for the GLM ( $\delta = 0$ ) is the usual parametric error rate

$$M_{\text{GLM}}(n, d, t) := M(n, d, 0, t) \asymp \begin{cases} t, & t \leq \sqrt{\frac{d}{n}}, \\ \sqrt{\frac{d}{n}}, & t \geq \sqrt{\frac{d}{n}}, \end{cases} \quad (5)$$

55 This is achieved by the trivial estimator  $\hat{\theta} = 0$  if  $t \leq \sqrt{\frac{d}{n}}$  and the simple empirical average estimator  
 56  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$  if  $t \geq \sqrt{\frac{d}{n}}$ . The rate  $\Theta(\sqrt{\frac{d}{n}})$  is then the *global* minimax rate, i.e., the largest error  
 57 over  $t = \|\theta_*\| > 0$ . This model does not have a phase transition with dimension. At the other extreme,  
 58 the GMM model ( $\delta = \frac{1}{2}$ ) undergoes a phase transition at the dimension  $d = n$ . At low dimension,  
 59  $d \leq n$ , the minimax rate was neatly shown in [Wu and Zhou, 2019, Appendix B] to be

$$M_{\text{GMM}}(n, d, t) \equiv M\left(n, d, \frac{1}{2}, t\right) \asymp \begin{cases} t, & t \leq \left(\frac{d}{n}\right)^{1/4}, \\ \frac{1}{t} \sqrt{\frac{d}{n}}, & \left(\frac{d}{n}\right)^{1/4} \leq t \leq 1, \\ \sqrt{\frac{d}{n}}, & t > 1 \end{cases} \quad (6)$$

60 whereas at high dimension  $d \geq n$ , it is as for the Gaussian location model in (5), i.e.,  $M_{\text{GMM}}(n, d, t) =$   
 61  $M_{\text{GLM}}(n, d, t)$ . However, we note that at this high dimensional regime, the loss is also lower bounded  
 62 by a constant even if the signal strength  $t$  is unbounded. The loss in (6) is achieved by the trivial  
 63 estimator  $\hat{\theta} = 0$  if  $t \leq \left(\frac{d}{n}\right)^{1/4}$  and by an estimator given by a properly scaled and shifted principal  
 64 component of the empirical covariance matrix of  $X_1^n := (X_1, X_2, \dots, X_n)$ , if  $t \geq \left(\frac{d}{n}\right)^{1/4}$ . For the  
 65 GMM at low dimension,  $d \leq n$ , the global minimax rate is  $\left(\frac{d}{n}\right)^{1/4}$ , which is worse than the minimax  
 66 rate of the Gaussian location model  $\sqrt{\frac{d}{n}}$ .

<sup>1</sup>Similar bounds can be derived for the squared loss by trivial extensions.

67 Therefore, the GMM has worse estimation performance compared to the GLM from three aspects:  
 68 First, at low dimension,  $d \leq n$ , it has a larger global minimax rate  $(\frac{d}{n})^{1/4}$  compared to the parametric  
 69 error rate of the GLM,  $\sqrt{\frac{d}{n}}$ ; Second, at low dimension,  $d \leq n$ , parametric error rate is achieved only  
 70 for constant separation  $t \geq 1$ ; Third, the transition to the high dimension regime occurs at  $d = n$ .

71 As is intuitively appealing from a “data-processing” reasoning, a Markov model with flip probability  
 72  $\delta'$  should allow for lower estimation error of  $\theta_*$  compared to a Markov model with  $\delta > \delta'$ . Indeed,  
 73 and as a specific simple example, any Markov model with  $\delta < \frac{1}{2}$  can be easily transformed to a GMM  
 74 model by randomizing the signs of each of the samples by an independent Rademacher variable. Thus  
 75 we may easily deduce, for instance, that since at high dimension ( $d \geq n$ ) the GLM and the GMM  
 76 have the same minimax rates, the minimax rates for  $d \geq n$  are in fact as in (5) for any  $\delta \in [0, 1]$ .  
 77 We thus exclusively focus in the rest of the paper on the regime  $d \leq n$ . As we will show, that is a  
 78 general phenomenon, and the reduction in estimation error when  $\delta$  is reduced is less profound as the  
 79 dimension increases.

### 80 1.3 Contributions

81 We first consider the case in which  $\delta$  is known to the estimator of  $\theta_*$ . For this case, we show (Theorem  
 82 1) that an estimator of  $\theta_*$  based on a computation of the principal component of a properly chosen  
 83 empirical covariance matrix smoothly interpolates the rates of the GLM in (5) and the GMM in (6).  
 84 At low dimension,  $d \leq \delta n$ , it achieves local minimax rate of

$$M(n, d, \delta, t) \lesssim \begin{cases} t, & t \leq (\frac{\delta d}{n})^{1/4} \\ \frac{1}{t} \sqrt{\frac{\delta d}{n}}, & (\frac{\delta d}{n})^{1/4} \leq t \leq \sqrt{\delta}, \\ \sqrt{\frac{d}{n}}, & t \geq \sqrt{\delta} \end{cases}, \quad (7)$$

85 at high dimension  $d \geq \delta n$ , it is as for the Gaussian location model in (5). The rate of this estimator is  
 86 then further shown to be asymptotically optimal (up to a logarithmic factor) via a minimax lower  
 87 bound (Theorem 2). Evidently, its performance smoothly interpolates between the performance of the  
 88 GLM and the GMM, and the loss is improved with the decrease of  $\delta$  from all three aspects previously  
 89 mentioned. First, at low dimension,  $d \leq \delta n$ , the global minimax rate is  $(\frac{\delta d}{n})^{1/4}$ . Thus whenever  $\delta$  is  
 90 as low as  $\delta = \frac{d}{n}$  it becomes the parametric error rate  $\Theta(\sqrt{\frac{d}{n}})$  of the GLM; whenever  $\delta = \Theta(1)$  it is  
 91 the same as  $\Theta((\frac{d}{n})^{1/4})$  for the GMM; and it smoothly interpolates between these rates for  $\frac{d}{n} \leq \delta \leq 1$ .  
 92 Second, at low dimension,  $d \leq \delta n$ , parametric error rate is achieved for signal strength as low as  
 93  $t = \Theta(\sqrt{\delta})$  (which again, matches the extremes  $t = \Theta(\sqrt{\frac{d}{n}})$  of GLM and  $t = \Theta(1)$  of GMM).  
 94 Third, the transition to the high dimension occurs at  $d = \delta n$ , which is again lower than the transition  
 95 point of the GMM given by  $d = n$ . This lower transition point allows to achieve the error rate of the  
 96 GLM for any signal strength, even in a regime in which the loss vanishes with  $n \rightarrow \infty$  (unlike for  
 97 GMM,  $\delta = \frac{1}{2}$ ); specifically, this occurs whenever  $\delta n \leq d \leq n$ . Beyond the formal proof, Appendix  
 98 A provides a heuristic justification for why the minimax error is naturally expected to scale as in (7).

99 Second, as a step towards the removal of the assumption that  $\delta$  is known to the estimator, we consider  
 100 the complementary problem of estimating  $\delta$  whenever an estimate  $\theta_{\sharp}$  of  $\theta_*$  is available (which can be  
 101 either exact  $\theta_{\sharp} = \theta_*$ , or inaccurate  $\theta_{\sharp} \neq \theta_*$ ). We propose a simple estimator for  $\delta$ , and analyze its  
 102 error in case of a mismatch (Theorem 3). We then specify this result to the matched case  $\theta_{\sharp} = \theta_*$   
 103 and show that in the non-trivial regime ( $\|\theta_*\| \lesssim 1$ ) its error rate is  $\tilde{O}(\frac{1}{\|\theta_*\|^2} \sqrt{\frac{1}{n}})$ . We then proceed  
 104 to show an impossibility lower bound of  $\Omega(\sqrt{\frac{1}{n}})$  (Proposition 5) for this error rate. We discuss the  
 105 challenges in settling the optimal estimation error rate of  $\delta$  as a function of  $\|\theta_*\|$ .

106 Third, we consider the case in which the estimator of  $\theta_*$  has no prior knowledge of  $\delta$ . We propose  
 107 a three step algorithm for this case (Algorithm 1). First, a (possibly) gross estimate  $\hat{\theta}^{(A)}$  of  $\theta_*$   
 108 is computed based on third of the samples, assuming the worst case of  $\delta = \frac{1}{2}$ . Then, an estimate  
 109  $\hat{\delta}^{(B)}$  of  $\delta$  is computed using another third of the samples, assuming the estimate  $\hat{\theta}^{(A)}$ . Finally, a  
 110 refined estimate  $\hat{\theta}^{(C)}$  of  $\theta_*$  is obtained by (essentially) assuming that  $\delta$  is  $\hat{\delta}^{(B)}$ . At each of the steps  
 111 above, the algorithm may stop and decide to return its current estimate when it determines that no

112 further improvement is possible by moving on to the next steps. We analyze the estimation loss of  
 113 this algorithm (Theorem 6), and show that this algorithm is capable of partially achieving the gains  
 114 associated with the case of known  $\delta$ .

#### 115 1.4 Related work

116 Both the GLM [Johnstone, 2002, Tsybakov, 2008] and GMM [Lindsay, 1995, McLachlan et al.,  
 117 2019] are classic models which were well-explored from numerous perspectives. For GMM, it is well  
 118 known that the maximum likelihood estimator (MLE) is consistent when the mixture components are  
 119 sufficiently separated [Redner and Walker, 1984]. Optimal error rates for GMM were derived, e.g.,  
 120 in [Heinrich and Kahn, 2015], using a minimum distance estimator (between an atomic distribution  
 121 convoluted with the Gaussian density and the empirical distribution of the samples). Nonetheless, the  
 122 GMM is a latent variable model, for which maximizing the likelihood w.r.t. the center parameters  
 123 is a non-convex optimization problem. In the last few years, there is a surge of interest in the non-  
 124 asymptotic performance analysis of computationally efficient estimation algorithms for this estimation  
 125 task. For example, Moitra and Valiant [2010], Kalai et al. [2010], Anandkumar et al. [2014], Hardt  
 126 and Price [2015], Wu and Yang [2020] have analyzed method-of-moments-based algorithms, and  
 127 various other papers considered the expectation-maximization (EM) algorithm [Balakrishnan et al.,  
 128 2017, Xu et al., 2016, Jin et al., 2016, Klusowski and Brinda, 2016, Weinberger and Bresler, 2021,  
 129 Dwivedi et al., 2020b,a, 2018, Zhao et al., 2018, Yan et al., 2017]. Specifically, the local minimax  
 130 rate for GMM in (6) was determined in [Wu and Zhou, 2019] as a benchmark for the operation of the  
 131 EM algorithm.

132 The model (2) is a simple instance of a HMM [Ephraim and Merhav, 2002, van Handel, 2008] in  
 133 high dimension. Parameter estimation in such models is practically performed via the Baum-Welch  
 134 algorithm [Baum et al., 1970], which is a computationally efficient version of EM for HMMs. The  
 135 consistency and asymptotic normality of the MLE for this case were established in [Bickel et al.,  
 136 1998]. To the best of our knowledge, there were little attempts to characterize the minimax rates in  
 137 such models. In [Aiyilam, 2018], a local version of the Baum-Welch algorithm was proposed, and  
 138 vanishing error of the convergence of the estimate to the true parameter was established for both  
 139 a population version as well as a finite-sample version. In [Yang et al., 2015], the analysis made  
 140 for the operation of EM on samples from memoryless latent variable models in [Balakrishnan et al.,  
 141 2017] was substantially extended to HMM. The results were then specified to the Gaussian model  
 142 with Markov signs (2) considered here, and it was shown that the Baum-Welch algorithm achieves  
 143 parametric error rate, and converges in a finite number of iterations [Yang et al., 2015, Corollary  
 144 2]. However, the qualifying condition for this result is that  $t = \|\theta_*\| \gtrsim \log \frac{1}{1-(1-2\delta)^2}$ , that is, a  
 145 non-trivial separation when  $\delta$  is constant, which further blows up as  $\delta \downarrow 0$ . By contrast, in this paper,  
 146 our goal is to characterize the estimation error in the regime of  $\delta$  and  $t = \|\theta_*\|$  in which the minimax  
 147 rate is affected by these parameters, and this requires analyzing vanishing  $\delta$  and  $t$ .

148 More broadly, there is a growing interest in advancing the quantitative understanding of the  
 149 performance of statistical learning and inference with dependent data. Bresler et al. [2020] studied  
 150 linear regression with Markovian covariates and characterized the minimax error rate in terms of  
 151 the mixing time of the Markov chain. A stochastic gradient descent-style algorithm adapted to  
 152 the Markov setting was shown to be minimax optimal. Statistical estimation problems including  
 153 linear and logistic regression with more general network dependencies among response variables  
 154 were studied by Daskalakis et al. [2019] and Kandiros et al. [2021]. Learnability and generalization  
 155 bounds were derived by Dagan et al. [2019] for dependent data satisfying the so-called Dobrushin's  
 156 condition. Finally, we also mention in passing the huge body of work studying reconstruction of  
 157 dependence structures and estimation of dependence parameters from graphical samples, which we  
 158 do not exhaustively list here.

#### 159 1.5 Notation conventions

160 For a vector  $v \in \mathbb{R}^d$ ,  $\|v\|$  is the Euclidean norm. For a positive definite matrix  $A$ , let  $\lambda_{\max}(A)$  and  
 161  $v_{\max}(A)$  be the maximal eigenvalue and the associated eigenvector (of unit norm) of  $A$ . Unless  
 162 otherwise stated, the constants involved in Bachman-Landau notation are numerical constants.  
 163 Specifically, they do not depend on the problem parameters  $(n, d, \delta, t)$ . We write  $a \gtrsim b$  (resp.  $a \lesssim b$ )  
 164 if there exists a numerical constant  $c > 0$  (resp.  $C > 0$ ) such that  $a \geq cb$  (resp.  $a \leq Cb$ ). If  $a \lesssim b$  and  
 165  $a \gtrsim b$ , then  $a \asymp b$ . By assuming that  $n$  is sufficiently large, we omit for brevity integer constraints

166 (ceiling and floor) on large quantities. We use the shorthand notation  $a \vee b := \max\{a, b\}$ ,  $a \wedge b :=$   
167  $\min\{a, b\}$ . For a real number  $a$ ,  $(a)_+ := a \vee 0$ . A sequence of objects  $X_1, \dots, X_n$  is denoted by  
168  $X_1^n$ . Expectation, variance and probability are denoted by  $\mathbb{E}$ ,  $\mathbb{V}$  and  $\mathbb{P}$ , respectively. If two random  
169 variables  $X$  and  $Y$  share the same distribution, then we write  $X \stackrel{d}{=} Y$ . All logarithms log are to the  
170 base  $e$ .

## 171 2 Mean estimation for a known flip probability

172 In this section, we consider the problem of estimating  $\theta_*$  whenever  $\delta$  is exactly known to the estimator.  
173 In that case, it may be assumed w.l.o.g. that  $\delta \in [0, \frac{1}{2}]$ , as otherwise one may negate each of the even  
174 samples to obtain an equivalent model with  $\delta$  replaced with  $1 - \delta$ . Hence also  $\rho \in [0, 1]$ . We next  
175 describe an estimator for this task, state a bound on its performance, and then show that it matches  
176 (up to a logarithmic factor) an impossibility result.

177 The estimator operationally interpolates and therefore simultaneously generalizes the empirical  
178 average estimator (30) and the (properly scaled) principal component estimator (32) analyzed in [Wu  
179 and Zhou, 2019, Appendix B]. It degenerates to the latter estimators if  $\delta \downarrow 0$  or  $\delta \uparrow \frac{1}{2}$ . Specifically,  
180 the estimator partitions the sample into blocks of equal length  $k$  each (which will later be set to  
181  $k = \frac{1}{8\delta}$ , according to the mixing time of the Markov chain  $S_0^n$ ). Let  $\ell$  denote the number of blocks  
182 respectively, so that  $k\ell = n$ . Let  $\mathcal{I}_i = \{(i-1)k+1, (i-1)k+2, \dots, ik\}$  denote the indices of the  
183  $i$ th block. Further, let  $\{R_i\}_{i \in [\ell]}$  be an i.i.d. Rademacher sequence ( $R_i \sim \text{Uniform}\{-1, 1\}$ ), and let

$$\bar{X}_i := R_i \cdot \frac{1}{k} \sum_{j \in \mathcal{I}_i} X_j \stackrel{d}{=} \bar{S}_i \theta_* + \bar{Z}_i \quad (8)$$

184 denote the average of the samples in the  $i$ th block (randomized with a sign  $R_i$ ), where

$$\bar{S}_i := \frac{1}{k} \sum_{j \in \mathcal{I}_i} S_j \quad (9)$$

185 is the *gain* (average of the signs) of the  $i$ th block, and

$$\bar{Z}_i := \frac{1}{k} \sum_{j \in \mathcal{I}_i} Z_j \sim N\left(0, \frac{1}{k} \cdot I_d\right) \quad (10)$$

186 is the average noise of the  $i$ th block. Due to the sign randomization, it holds that  $\{\bar{S}_i\}_{i \in [\ell]}$  is an i.i.d.  
187 sequence. Since  $\{\bar{Z}_i\}_{i \in [\ell]}$  is also an i.i.d. sequence, then so is  $\{\bar{X}_i\}_{i \in [\ell]}$ . For notational simplicity  
188 we will omit the block index  $i$  of a generic block. For block length  $k$ , we denote by

$$\xi_k := \mathbb{E}[\bar{S}^2] = \mathbb{E}\left[\left(\frac{1}{k} \sum_{j=1}^k S_j\right)^2\right] \quad (11)$$

189 the second moment of the gain  $\bar{S}$ . Note that  $\xi_k \in [\frac{1}{k}, 1]$  for any  $\delta \in [0, \frac{1}{2}]$  and is in particular always  
190 positive. For a sequence of samples  $X_1^n = (X_1, \dots, X_n)$ , we define by  $\hat{\Sigma}_{n,k}(X_1^n)$  the empirical  
191 covariance matrix of the averaged samples over blocks  $\{\bar{X}_i\}_{i \in [\ell]}$ , that is

$$\hat{\Sigma}_{n,k}(X_1^n) := \frac{1}{\ell} \sum_{i=1}^{\ell} \bar{X}_i \bar{X}_i^\top, \quad (12)$$

192 whose population average is  $\Sigma_{n,k}(\theta_*)$ , where

$$\Sigma_{n,k}(\theta_*) := \mathbb{E}[\bar{X} \bar{X}^\top] = \xi_k \theta_* \theta_*^\top + \frac{1}{k} I_d. \quad (13)$$

193 We note that  $\theta$  is the principal component of  $\Sigma_{n,k}(\theta)$ , that is,  $\lambda_{\max}(\Sigma_{n,k}(\theta)) = \xi_k \|\theta\|^2 + \frac{1}{k}$  and the  
194 corresponding eigenvector is  $v_{\max}(\Sigma_{n,k}(\theta)) = \theta$ . We thus consider the following estimator for  $\theta_*$ ,  
195 from a sequence  $X_1^n$ , and with a block length of  $k$

$$\hat{\theta}_{\text{cov}}(X_1^n; k) := \sqrt{\frac{\left(\lambda_{\max}(\hat{\Sigma}_{n,k}(X_1^n)) - \frac{1}{k}\right)_+}{\xi_k}} \cdot v_{\max}(\hat{\Sigma}_{n,k}(X_1^n)). \quad (14)$$

196 The estimator is thus constructed from two types of averages: First, a coherent average of the samples  
 197 at each block, to obtain  $\ell$  block-samples  $\bar{X}_i$  with gain  $\bar{S}_i$  and noise variance reduced by a factor of  
 198  $k$ . Second, an incoherent average of the “square” of the  $\ell$  block-samples  $\bar{X}_i \bar{X}_i^\top$ , which resolves the  
 199 remaining sign ambiguity between blocks. This can be seen as a balance between two extreme cases:  
 200 If  $\delta = 0$ , then this reduces the problem to the Gaussian location model (with a sign ambiguity) and  
 201  $k = n$  is an optimal choice. If  $\delta = \frac{1}{2}$ , then this reduces the problem to the two-component Gaussian  
 202 mixture model, in which coherent averaging is non-beneficial and  $k = 1$  leads to optimal error rates.  
 203 It is thus obvious that the optimal choice of  $k$  depends on  $\delta$ , the flip probability of the Markov chain.  
 204 Choosing  $k = \Theta(\frac{1}{\delta})$ , that is, a block length proportional to the mixing time of the Markov chain,  
 205 assures that the random gain  $\bar{S}$  is  $\pm 1$  with a (constant) high probability. In fact, an elementary, yet  
 206 crucial, part of the analysis establishes that the random gain  $\bar{S}$  has constant variance for this choice  
 207 of block length (see Lemma 8 in Appendix B.1). On the other hand, if  $k = \Omega(\frac{1}{\delta})$ , then the random  
 208 gain  $\bar{S}$  will not be  $\pm 1$  (or not even bounded away from zero) with high probability, and such choice  
 209 is never efficient. Specifically, we consider the estimator in (14) with  $k = \frac{1}{8\delta}$ . The above estimation  
 210 procedure is depicted in Figure 2 in Appendix B.1.

211 Let us denote

$$\beta(n, d, \delta) := \sqrt{\frac{d}{n}} \vee \left( \frac{\delta d}{n} \right)^{1/4}, \quad (15)$$

212 which will actually be the global minimax rate.

213 **Theorem 1.** Assume that  $\delta \geq \frac{1}{n}$  and  $d \leq n$ , and set  $\hat{\theta} \equiv \hat{\theta}_{\text{cov}}(X_1^n; k)$  with  $k = \frac{1}{8\delta}$ . Then, there exist  
 214 numerical constants  $c_0, c_1, c_2 > 0$  such that for every  $\theta_* \in \mathbb{R}^d$

$$\mathbb{E} \left[ \text{loss}(\hat{\theta}, \theta_*) \right] \leq c_0 \cdot \begin{cases} \beta(n, d, \delta), & \|\theta_*\| \leq \beta(n, d, \delta) \\ \sqrt{\frac{d}{n}} + \frac{1}{\|\theta_*\|} \sqrt{\frac{\delta d}{n}} + \frac{1}{\|\theta_*\|} \cdot \frac{d}{n}, & \beta(n, d, \delta) \leq \|\theta_*\| \end{cases} \quad (16)$$

215 and

$$\text{loss}(\hat{\theta}, \theta_*) \leq c_1 \cdot \log(n) \cdot \mathbb{E} \left[ \text{loss}(\hat{\theta}, \theta_*) \right] \quad (17)$$

216 with probability larger than  $1 - \frac{c_2}{n}$ .

217 Theorem 1 will be proved in Appendix B.1. A simple consequence of Theorem 1 is that the  
 218 upper bound on the minimax rate stated in (7) above holds in low dimension,  $d \leq \delta n$ , whereas  
 219  $M(n, d, \delta, t) \asymp M_{\text{GLM}}(n, d, t)$  holds in high dimension  $d \geq \delta n$ .<sup>2</sup> We also remark that  $\delta \geq \frac{1}{n}$  is not  
 220 a restrictive condition since otherwise the flip probability is so low that the model (2) is essentially  
 221 equivalent to GLM. See Remark 7 in Appendix B.1. Numerical validation of the performance of the  
 222 estimator  $\hat{\theta}_{\text{cov}}(X_1^n; k)$  is shown in Appendix F.

223 We next consider an impossibility result. As we have seen, at high dimension  $d \geq \delta n$  the minimax  
 224 error rates achieved are the same as for the Gaussian location model, and thus clearly cannot be  
 225 improved. We thus next focus on the low dimensional regime  $d \leq \delta n$ .

226 **Theorem 2.** Assume that  $2 \leq d \leq \delta n$  and  $n \geq \frac{128}{d}$ . Then the local minimax rate is bounded as

$$M(n, d, \delta, t) \gtrsim \frac{1}{\sqrt{\log(n)}} \cdot \begin{cases} t, & t \leq \left( \frac{\delta d}{n} \right)^{1/4} \\ \frac{1}{t} \sqrt{\frac{\delta d}{n}}, & \left( \frac{\delta d}{n} \right)^{1/4} \leq t \leq \sqrt{\delta} \\ \sqrt{\frac{d}{n}}, & t \geq \sqrt{\delta} \end{cases} \quad (18)$$

227 Hence, the minimax rates achieved by the estimator in Theorem 1 are nearly asymptotically optimal,  
 228 up to a  $\sqrt{\log(n)}$  factor. The full proof of Theorem 2 together with a summary of the main ideas used  
 229 in the proof is presented in Appendix B.2.

<sup>2</sup>Note that when  $\|\theta_*\| \leq \beta(n, d, \delta)$ , the estimator  $\hat{\theta}_{\text{cov}}(X_1^n; k)$  in Theorem 1 only achieves a rate  $\beta(n, d, \delta)$  which is larger than the promised rate  $\|\theta_*\|$  in (7). To attain the claimed error rate, consider the trivial estimator  $\hat{\theta}_0(X_1^n) \equiv 0$  which incurs loss  $\|\theta_*\|$  for any  $\theta_* \in \mathbb{R}^d$ . Taking the minimum of this rate and (16) yields the desired rate (7).

### 230 3 Flip probability estimation for a given estimator of $\theta_*$

231 In this section, we consider the problem of estimating  $\delta$  whenever  $\theta_*$  is approximately known to be  
 232  $\theta_{\#}$ . We propose a simple estimator, and then discuss the importance of the accuracy of  $\theta_*$ . We then  
 233 derive an impossibility result for the matched case,  $\theta_{\#} = \theta_*$ .

234 An estimator for  $\delta$  can be easily obtained from an estimator for  $\rho$ , via the plug-in  $\hat{\delta} = \frac{1}{2}(1 - \hat{\rho})$ .  
 235 Assuming for simplicity that  $n$  is even, a natural estimator for  $\rho$  is then

$$\hat{\rho}_{\text{corr}}(X_1^n; \theta_{\#}) = \frac{1}{\|\theta_{\#}\|^2} \cdot \frac{2}{n} \sum_{i=1}^{n/2} X_{2i}^{\top} X_{2i-1}. \quad (19)$$

236 That is, the estimator is based on evaluating the correlation of each of two adjacent samples  $X_{2i}$  and  
 237  $X_{2i-1}$ , whose population version is  $\mathbb{E}[X_{2i}^{\top} X_{2i-1}] = \rho \|\theta_*\|^2$ . We first state a general bound on the  
 238 estimation error of this estimator. We then consider the case in which  $\theta_*$  is known, and show how the  
 239 estimation error is improved in this case.

240 **Theorem 3.** *Assume that  $d \leq n$ . Let  $\theta_* \in \mathbb{R}^d$  and let  $\theta_{\#}$  be an estimate of  $\theta_*$ . Set  $\hat{\rho} \equiv \hat{\rho}_{\text{corr}}(X_1^n; \theta_{\#})$   
 241 and  $\hat{\delta} = \frac{1}{2}(1 - \hat{\rho})$ . Then, it holds with probability  $1 - \frac{\delta}{n}$  that*

$$\left| \hat{\delta} - \delta \right| = \frac{1}{2} |\hat{\rho} - \rho| \leq \frac{\left| \|\theta_*\|^2 - \|\theta_{\#}\|^2 \right|}{\|\theta_{\#}\|^2} + 16 \log(n) \left[ \sqrt{\frac{\delta}{n}} + \frac{1}{\|\theta_{\#}\|} \sqrt{\frac{1}{n}} + \frac{1}{\|\theta_{\#}\|^2} \sqrt{\frac{d}{n}} \right]. \quad (20)$$

242 The proof of Theorem 3 appears in Appendix C.1. Note that Theorem 3 states a high-probability  
 243 bound, suitable to its usage later on in Section 4. A bound on the expectation of the error can be  
 244 obtained by standard methods (integrating tails).

245 **The effect of knowledge of  $\theta_*$**  If  $\theta_*$  is known up to a sign, i.e.,  $\theta_{\#} = \pm\theta_*$ , then for the purpose of  $\rho$   
 246 (or equivalently  $\delta$ ) estimation, the model (2) can be reduced to a one-dimensional model by rotational  
 247 invariance of isotropic Gaussian (See additional details in Appendix C.2). It then immediately follows  
 248 from Theorem 3 that:

249 **Corollary 4.** *Assume that  $d \leq n$ ,  $\|\theta_*\| \leq 1$  and  $\theta_{\#} = \pm\theta_*$ . Let  $U_1^n$  be defined in as  $U_i :=$   
 250  $\|\theta_*\| \cdot S_i + W_i$  where  $W_i \sim N(0, 1)$  i.i.d.,  $\hat{\rho} \equiv \hat{\rho}_{\text{corr}}(U_1^n; \theta_{\#})$  and  $\hat{\delta} = \frac{1}{2}(1 - \hat{\rho})$ . Then it holds with  
 251 probability  $1 - \frac{\delta}{n}$  that*

$$\left| \hat{\delta} - \delta \right| = \frac{1}{2} |\hat{\rho} - \rho| \leq \frac{18 \log(n)}{\|\theta_*\|^2} \sqrt{\frac{1}{n}}. \quad (21)$$

252 Numerical validation of the performance of the estimator  $\hat{\delta}_{\text{corr}}(X_1^n; \theta_{\#}) = \frac{1}{2}(1 - \hat{\rho}_{\text{corr}}(X_1^n; \theta_{\#}))$  in the  
 253 mismatched (Theorem 3) and matched (Corollary 4) cases is provided in Appendix F.

254 We next consider an impossibility lower bound.

255 **Proposition 5.** *Suppose that  $\theta_{\#} = \theta_*$  and  $\|\theta_*\| \leq \frac{1}{\sqrt{2}}$ . Then*

$$\inf_{\hat{\delta}(U_1^n)} \sup_{\delta \in [0, 1]} \mathbb{E} \left[ \left| \delta - \hat{\delta}(U_1^n) \right| \right] \geq \frac{1}{10\sqrt{n}}. \quad (22)$$

256 *Here the infimum is over any estimator  $\hat{\delta}(U_1^n)$  based on the model  $U_i = \|\theta_*\| S_i + W_i$  where each  
 257  $W_i$  is i.i.d.  $N(0, 1)$ .*

258 The proof of Proposition 5 is presented in Appendix C.2. According to Corollary 4 and Proposition  
 259 5, in estimating  $\delta$  with a known  $\theta_*$ , though the dependence  $\Theta(\frac{1}{\sqrt{n}})$  of the minimax error rate on  
 260 the sample size is shown to be nearly optimal, it is unclear what the optimal dependence on the  
 261 signal strength should be. This is left as an interesting open question and we discuss the challenges  
 262 associated with this problem in Appendix C.2.

### 263 4 Mean estimation under an unknown flip probability

264 As we have seen, if an estimator for  $\theta_*$  knows the value of  $\delta$ , and if both  $\delta \leq \frac{1}{2}$  and  $d \leq \delta n$  hold,  
 265 then the estimator can achieve improved error rate over the GMM case ( $\delta = \frac{1}{2}$ ). In this section, we

266 assume that both  $\theta_*$  and  $\delta$  are unknown, and so the estimator is required to estimate  $\delta$  in order to use  
 267 this knowledge for an estimator of  $\theta_*$ . We propose an estimation procedure of three steps based on  
 268 sample splitting of  $3n$  samples. We mention at the outset that the regime in which improvement is  
 269 possible will be for low signal strength  $\|\theta_*\| \lesssim 1$  (low separation between the components), and up  
 270 to a dimension which depends on  $\delta$ . Of course the estimation procedure does not know  $(\theta_*, \delta)$  in  
 271 advance, and so it is required to identify if  $(\theta_*, \delta)$  are in this regime during its operation.

272 We now begin with an overview of the steps of the estimation algorithm. At Step A, the algorithm  
 273 estimates  $\theta_*$  based on  $X_1^n$  assuming a Gaussian mixture model ( $\delta = \frac{1}{2}$ ) to obtain an estimate  $\hat{\theta}^{(A)}$ .  
 274 Then, based on  $\|\hat{\theta}^{(A)}\|$ , the algorithm decides whether improvement is potentially possible had  $\delta$   
 275 was known. There are two cases. The first case is that  $\|\hat{\theta}^{(A)}\|$  is too low, and then its estimate is not  
 276 sufficiently accurate to be used in the next steps. Essentially, this happens when the norm is below  
 277 the global minimax rate  $(\frac{d}{n})^{1/4}$ , and the estimation error of the norm on the same scale as the norm  
 278 of  $\|\theta_*\|$ . A trivial estimator of  $\hat{\theta} = 0$  is then optimal in terms of error rates. It can be already noted at  
 279 this step that the minimax rate for the known  $\delta$  case is  $(\frac{\delta d}{n})^{1/4}$ , whereas here the algorithm stops and  
 280 estimates  $\hat{\theta} = 0$  even if just  $\|\theta_*\| \lesssim (\frac{d}{n})^{1/4}$ , leading to larger global minimax rate. The second case is  
 281 that  $\|\theta_*\|$  is larger than a constant. In this case, the estimation based on a GMM already achieves the  
 282 optimal parametric  $O(\sqrt{\frac{d}{n}})$  error rate of the Gaussian location model, and so no further estimation  
 283 steps are necessary. Otherwise, an improvement in the estimation is possible. The algorithm proceeds  
 284 to Step B, and uses  $X_{n+1}^{2n}$  to obtain an estimate  $\hat{\delta}^{(B)}$  of  $\delta$  based on the mismatched  $\theta_{\#} \equiv \hat{\theta}^{(A)}$ . Then,  
 285 based on the estimate  $\hat{\delta}^{(B)}$  the algorithm decides whether the accuracy of  $\hat{\delta}^{(B)}$  is sufficient to be used  
 286 in an refined estimation of  $\theta_*$ . If the accuracy of  $\hat{\delta}^{(B)}$  is not good enough, then the algorithm outputs  
 287 the estimate from Step A, that is  $\hat{\theta}^{(A)}$ . Otherwise, it proceeds to Step C, in which  $\theta_*$  is re-estimated  
 288 using  $X_{2n+1}^{3n}$ , based on a mismatched choice of  $k$ , that is  $k \asymp \frac{1}{\hat{\delta}^{(B)}}$  instead of  $k = \frac{1}{\delta}$ . Intuitively, the  
 289 estimated value  $\hat{\delta}^{(B)}$  should be larger than  $\delta$  so the resulting block size  $k \asymp \frac{1}{\hat{\delta}^{(B)}}$  will be such that the  
 290 gain in the block is still close to 1 with high probability. On the other hand, it is desired that  $\hat{\delta}^{(B)}$   
 291 will be on the same scale as  $\delta$  so that the estimation rate (296) (see also (7)) – which now essentially  
 292 holds with  $\hat{\delta}^{(B)}$  instead of  $\delta$  – would be as small as possible. Thus, if the algorithm has assured in  
 293 Step B that  $\hat{\delta}^{(B)} \asymp \delta$ , then at Step C it will achieve the error rate indicated in (296).

294 A formal description of the operation of the estimation algorithm is provided in Algorithm 1. We  
 295 note in passing that refining the estimation of  $\delta$  can be easily incorporated as a fourth step of this  
 296 algorithm, but we do not present this in order to keep the statement of the result simple. The error of  
 297 the estimator output by Algorithm 1 is as follows:

298 **Theorem 6.** Assume that  $d \leq \frac{n}{2\lambda_\theta \log^2(n) \wedge 16}$ . Then, there exist numerical constants  $c_1, c_2 \geq 0$  and  
 299  $\lambda_\theta, \lambda_\delta \geq 1$  so that the output  $\hat{\theta}$  of the estimation Algorithm 1 satisfies that for any  $\theta_* \in \mathbb{R}^d$ , with  
 300 probability  $1 - O(\frac{1}{n})$ :

301 • If  $d \leq \frac{1}{64\lambda_\delta \lambda_\theta \log^2(n)} \delta^4 n$  then

$$\text{loss}(\hat{\theta}, \theta_*) \leq c_1 \log n \cdot \begin{cases} \|\theta_*\|, & \|\theta_*\| \leq \lambda_\theta \log(n) (\frac{d}{n})^{1/4} \\ \frac{1}{\|\theta_*\|} \sqrt{\frac{d}{n}}, & \lambda_\theta \log(n) (\frac{d}{n})^{1/4} \leq \|\theta_*\| \leq \sqrt{8\lambda_\delta \lambda_\theta \log(n)} (\frac{d}{\delta^2 n})^{1/4} \\ \frac{1}{\|\theta_*\|} \sqrt{\frac{\delta d}{n}}, & \sqrt{8\lambda_\delta \lambda_\theta \log(n)} (\frac{d}{\delta^2 n})^{1/4} \leq \|\theta_*\| \leq \sqrt{\delta} \\ \sqrt{\frac{d}{n}}, & \|\theta_*\| \geq \sqrt{\delta} \end{cases}; \quad (26)$$

302 • If  $d \geq \frac{1}{64\lambda_\delta \lambda_\theta \log^2(n)} \delta^4 n$  then

$$\text{loss}(\hat{\theta}, \theta_*) \leq c_2 \log n \cdot \begin{cases} \|\theta_*\|, & \|\theta_*\| \leq \lambda_\theta \log(n) (\frac{d}{n})^{1/4} \\ \frac{1}{\|\theta_*\|} \sqrt{\frac{d}{n}}, & \lambda_\theta \log(n) (\frac{d}{n})^{1/4} \leq \|\theta_*\| \leq \sqrt{8\lambda_\delta \lambda_\theta \log(n)} (\frac{d}{\delta^2 n})^{1/4} \\ \sqrt{\frac{d}{n}}, & \|\theta_*\| \geq \sqrt{8\lambda_\delta \lambda_\theta \log(n)} (\frac{d}{\delta^2 n})^{1/4} \end{cases}. \quad (27)$$

---

**Algorithm 1** Mean estimation for unknown  $\delta$ 

---

1: **input:** Parameters  $\lambda_\theta, \lambda_\delta > 0$  (from (296) (297) (299)),  $3n$  data samples  $X_1^{3n}$  from the model (2)

2: **step A:** Estimate  $\theta_*$  assuming a Gaussian mixture model

$$\hat{\theta}^{(A)} \equiv \hat{\theta}_{\text{cov}}(X_1^n; k = 1) \quad (23)$$

3: **if**  $\|\hat{\theta}^{(A)}\| \leq 2\lambda_\theta \cdot \log(n) \cdot (\frac{d}{n})^{1/4}$  **then**

4:     **return**  $\hat{\theta} = 0$

▷ No further improvement can be guaranteed

5: **else if**  $\|\hat{\theta}^{(A)}\| \geq \frac{1}{2}$  **then**

6:     **return**  $\hat{\theta} = \hat{\theta}^{(A)}$

▷ No further improvement is possible

7: **end if**

8: **step B:** Estimate  $\delta$  assuming a mismatched mean value of  $\hat{\theta}^{(A)}$

$$\hat{\delta}^{(B)} = \hat{\delta}_{\text{corr}}(X_{n+1}^{2n}; \hat{\theta}^{(A)}) \quad (24)$$

9: **if**  $\hat{\delta}^{(B)} \leq 64\lambda_\delta \lambda_\theta \frac{\log(n)}{\|\hat{\theta}^{(A)}\|^2} \sqrt{\frac{d}{n}}$  **then**

10:     **return**  $\hat{\theta} = \hat{\theta}^{(A)}$

▷ No further improvement can be guaranteed

11: **end if**

12: **step C:** Estimate  $\theta_*$  assuming a Markov model with a mismatched flip probability  $\hat{\delta}^{(B)}$

$$\hat{\theta}^{(C)} = \hat{\theta}_{\text{cov}} \left( X_{2n+1}^{3n}; k = \frac{1}{16\hat{\delta}^{(B)}} \right) \quad (25)$$

13: **return**  $\hat{\theta} = \hat{\theta}^{(C)}$ .

---

303 As can be observed, ignoring logarithmic factors, if the dimension is low enough  $d \lesssim \delta^4 n$  and  
304  $\|\theta_*\| \gtrsim (\frac{d}{\delta^2 n})^{1/4}$  then the error rates of the known  $\delta$  case are recovered. The analysis of Algorithm 1  
305 appears in Appendix D. Numerical validation of the performance of Algorithm 1 can be found in  
306 Appendix F.

307 **The impact of lack of knowledge of  $\delta$**  The deterioration in the estimation of  $\theta_*$  due to the lack  
308 of knowledge of  $\delta$  includes the following (ignoring logarithmic factors). First, the global minimax  
309 rate is  $(\frac{d}{n})^{1/4}$  as for the GMM, instead of the rate  $(\frac{\delta d}{n})^{1/4}$  for the Markov model case with known  $\delta$ .  
310 Second, at the regime  $(\frac{d}{n})^{1/4} \lesssim \|\theta_*\| \lesssim (\frac{d}{\delta^2 n})^{1/4}$  the error rate is  $O(\frac{1}{\|\theta_*\|} \sqrt{\frac{d}{n}})$  instead of the lower  
311  $O(\frac{1}{\|\theta_*\|} \sqrt{\frac{\delta d}{n}})$ . Third, the algorithm is only effective when the dimension is as low as  $d \lesssim \delta^4 n$ . For  
312 higher dimensions, the rates of the GMM are achieved, which can be achieved even without the  
313 knowledge of  $\delta$ .

## 314 5 Conclusion and future work

315 In this paper, we have considered an elementary, yet fundamental, high-dimensional model with  
316 memory. We have obtained a sharp bound on the minimax rate of estimation in case the underlying  
317 statistical dependency (flip probability) is known, and proposed a three-step estimation algorithm  
318 when it is unknown. This has revealed the gains possible in estimation rates due to the memory  
319 between the samples, and smoothly interpolated between the extreme cases of GLM and GMM. An  
320 interesting open problem is to either characterize the optimality of the algorithm or improving in  
321 the unknown  $\delta$  case; this requires understanding optimal estimation of the flip probability, and the  
322 challenges associated with this problem are illuminated.

323 Naturally, as the model considered in this paper is basic, there is an ample of possibilities to generalize  
324 this model. These include, a larger number of components in the mixture, statistical dependency with  
325 a more complicated graphical structure between the data samples, existence of nuisance parameters  
326 such as the noise variance, sharp finite-sample/finite-iteration analysis of specific practical algorithms  
327 such as Baum-Welch, and so on.

328 **References**

- 329 Dhroova Aiylam. Parameter estimation in HMMs with guaranteed convergence. Master’s thesis,  
330 Massachusetts Institute of Technology, 2018.
- 331 Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M Kakade, and Matus Telgarsky. Tensor  
332 decompositions for learning latent variable models. *The Journal of Machine Learning Research*,  
333 15(1):2773–2832, 2014.
- 334 Sivaraman Balakrishnan, Martin J. Wainwright, and Bin Yu. Statistical guarantees for the EM  
335 algorithm: From population to sample-based analysis. *The Annals of Statistics*, 45(1):77–120,  
336 2017.
- 337 Leonard E. Baum, Ted Petrie, George Soules, and Norman Weiss. A maximization technique  
338 occurring in the statistical analysis of probabilistic functions of Markov chains. *The Annals of*  
339 *Mathematical Statistics*, 41(1):164–171, 1970.
- 340 Peter J. Bickel, Ya’acov Ritov, and Tobias Ryden. Asymptotic normality of the maximum-likelihood  
341 estimator for general hidden Markov models. *The Annals of Statistics*, 26(4):1614–1635, 1998.
- 342 Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic*  
343 *theory of independence*. Oxford university press, 2013.
- 344 Guy Bresler, Prateek Jain, Dheeraj Nagaraj, Praneeth Netrapalli, and Xian Wu. Least squares  
345 regression with markovian data: Fundamental limits and algorithms. In *Proceedings of the 34th*  
346 *International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook, NY,  
347 USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 348 T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2006. ISBN  
349 0471241954.
- 350 Yuval Dagan, Constantinos Daskalakis, Nishanth Dikkala, and Siddhartha Jayanti. Learning from  
351 weakly dependent data under dobrushin’s condition. *CoRR*, abs/1906.09247, 2019. URL <http://arxiv.org/abs/1906.09247>.
- 353 Constantinos Daskalakis, Nishanth Dikkala, and Ioannis Panageas. Regression from dependent  
354 observations. In *STOC’19—Proceedings of the 51st Annual ACM SIGACT Symposium on Theory*  
355 *of Computing*, pages 881–889. ACM, New York, 2019.
- 356 Raaz Dwivedi, Koulik Khamaru, Martin J Wainwright, Michael I. Jordan, et al. Theoretical  
357 guarantees for EM under misspecified Gaussian mixture models. In *Advances in Neural Information*  
358 *Processing Systems*, pages 9681–9689, 2018.
- 359 Raaz Dwivedi, Nhat Ho, Koulik Khamaru, Martin Wainwright, Michael Jordan, and Bin Yu. Sharp  
360 analysis of expectation-maximization for weakly identifiable models. In *International Conference*  
361 *on Artificial Intelligence and Statistics*, pages 1866–1876. PMLR, 2020a.
- 362 Raaz Dwivedi, Nhat Ho, Koulik Khamaru, Martin J Wainwright, Michael I Jordan, and Bin Yu.  
363 Singularity, misspecification and the convergence rate of EM. *The Annals of Statistics*, 48(6):  
364 3161–3182, 2020b.
- 365 Yariv Ephraim and Neri Merhav. Hidden Markov processes. *IEEE Transactions on information*  
366 *theory*, 48(6):1518–1569, 2002.
- 367 Moritz Hardt and Eric Price. Tight bounds for learning a mixture of two Gaussians. In *Proceedings*  
368 *of the forty-seventh annual ACM symposium on Theory of computing*, pages 753–760. ACM, 2015.
- 369 Philippe Heinrich and Jonas Kahn. Optimal rates for finite mixture estimation. *arXiv preprint*  
370 *arXiv:1507.04313*, 2015.
- 371 Chi Jin, Yuchen Zhang, Sivaraman Balakrishnan, Martin J. Wainwright, and Michael I. Jordan.  
372 Local maxima in the likelihood of Gaussian mixture models: Structural results and algorithmic  
373 consequences. In *Advances in neural information processing systems*, pages 4116–4124, 2016.

- 374 Iain M. Johnstone. Function estimation and Gaussian sequence models. *Unpublished manuscript*, 2  
375 (5.3):2, 2002.
- 376 Adam Tauman Kalai, Ankur Moitra, and Gregory Valiant. Efficiently learning mixtures of two  
377 gaussians. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages  
378 553–562, 2010.
- 379 Vardis Kandiros, Yuval Dagan, Nishanth Dikkala, Surbhi Goel, and Constantinos Daskalakis.  
380 Statistical estimation from dependent data. In Marina Meila and Tong Zhang, editors, *Proceedings*  
381 *of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine*  
382 *Learning Research*, pages 5269–5278. PMLR, 18–24 Jul 2021. URL [https://proceedings.](https://proceedings.mlr.press/v139/kandiros21a.html)  
383 [mlr.press/v139/kandiros21a.html](https://proceedings.mlr.press/v139/kandiros21a.html).
- 384 Jason M. Klusowski and W.D. Brinda. Statistical guarantees for estimating the centers of a two-  
385 component Gaussian mixture by EM. *arXiv preprint arXiv:1608.02280*, 2016.
- 386 Bruce G. Lindsay. Mixture models: Theory, geometry, and applications. Ims, 1995.
- 387 Geoffrey J McLachlan, Sharon X Lee, and Suren I Rathnayake. Finite mixture models. *Annual*  
388 *review of statistics and its application*, 6:355–378, 2019.
- 389 Ankur Moitra and Gregory Valiant. Settling the polynomial learnability of mixtures of gaussians. In  
390 *2010 IEEE 51st Annual Symposium on Foundations of Computer Science*, pages 93–102. IEEE,  
391 2010.
- 392 Richard A. Redner and Homer F. Walker. Mixture densities, maximum likelihood and the EM  
393 algorithm. *SIAM review*, 26(2):195–239, 1984.
- 394 Phillippe Rigollet and Jan-Christian Hütter. High dimensional statistics. *Lecture notes for course*  
395 *18S997*, 2019.
- 396 Alexandre B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Publishing Company,  
397 Incorporated, 2008. ISBN 0387790519.
- 398 Ramon van Handel. Hidden Markov models. *Unpublished lecture notes*, 2008.
- 399 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,  
400 volume 47. Cambridge University Press, 2018.
- 401 Martin J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48.  
402 Cambridge University Press, 2019.
- 403 Nir Weinberger and Guy Bresler. The em algorithm is adaptively-optimal for unbalanced symmetric  
404 gaussian mixtures. *arXiv preprint arXiv:2103.15653*, 2021.
- 405 Yihong Wu and Pengkun Yang. Optimal estimation of Gaussian mixtures via denoised method of  
406 moments. *The Annals of Statistics*, 48(4):1981–2007, 2020.
- 407 Yihong Wu and Harrison H. Zhou. Randomly initialized EM algorithm for two-component Gaussian  
408 mixture achieves near optimality in  $o(\sqrt{n})$  iterations. *arXiv preprint arXiv:1908.10935*, 2019.
- 409 Ji Xu, Daniel J. Hsu, and Arian Maleki. Global analysis of expectation maximization for mixtures of  
410 two Gaussians. In *Advances in Neural Information Processing Systems*, pages 2676–2684, 2016.
- 411 Bowei Yan, Mingzhang Yin, and Purnamrita Sarkar. Convergence of gradient EM on multi-component  
412 mixture of Gaussians. In *Advances in Neural Information Processing Systems*, pages 6956–6966,  
413 2017.
- 414 Fanny Yang, Sivaraman Balakrishnan, and Martin J. Wainwright. Statistical and computational  
415 guarantees for the Baum-Welch algorithm. In *2015 53rd Annual Allerton Conference on*  
416 *Communication, Control, and Computing (Allerton)*, pages 658–665. IEEE, 2015.
- 417 Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of  
418 convergence. *Annals of Statistics*, pages 1564–1599, 1999.
- 419 Ruofei Zhao, Yuanzhi Li, and Yuekai Sun. Statistical convergence of the EM algorithm on Gaussian  
420 mixture models. *arXiv preprint arXiv:1810.04090*, 2018.

421 **Checklist**

- 422 1. For all authors...
- 423 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
424 contributions and scope? [Yes]
- 425 (b) Did you describe the limitations of your work? [Yes]
- 426 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 427 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
428 them? [Yes]
- 429 2. If you are including theoretical results...
- 430 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 431 (b) Did you include complete proofs of all theoretical results? [Yes]
- 432 3. If you ran experiments...
- 433 (a) Did you include the code, data, and instructions needed to reproduce the main  
434 experimental results (either in the supplemental material or as a URL)? [N/A]
- 435 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
436 were chosen)? [N/A]
- 437 (c) Did you report error bars (e.g., with respect to the random seed after running  
438 experiments multiple times)? [N/A]
- 439 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
440 of GPUs, internal cluster, or cloud provider)? [N/A]
- 441 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 442 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 443 (b) Did you mention the license of the assets? [N/A]
- 444 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 445
- 446 (d) Did you discuss whether and how consent was obtained from people whose data you're  
447 using/curating? [N/A]
- 448 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
449 information or offensive content? [N/A]
- 450 5. If you used crowdsourcing or conducted research with human subjects...
- 451 (a) Did you include the full text of instructions given to participants and screenshots, if  
452 applicable? [N/A]
- 453 (b) Did you describe any potential participant risks, with links to Institutional Review  
454 Board (IRB) approvals, if applicable? [N/A]
- 455 (c) Did you include the estimated hourly wage paid to participants and the total amount  
456 spent on participant compensation? [N/A]