

ML ESTIMATION FROM BITS

Anonymous authors

Paper under double-blind review

ABSTRACT

Estimating statistical parameters from quantized signals has received significant attention in recent years, as recovering information from quantized measurements has numerous applications across signal processing, communications, and data analysis. In this work, we focus on maximum likelihood (ML) estimation of statistical parameters from quantized samples. Directly solving the ML problem is challenging, as the likelihood function involves multiple integrals that are difficult to evaluate. To address this challenge, we propose an expectation-conditional-maximization (ECM) algorithm under a general distributional framework. Our approach generalizes the quantization model to multi-bit settings and allows the underlying signal to follow any distribution within the normal mean-variance mixture family. By designing suitable surrogate functions, the ECM algorithm ensures that all model parameters can be updated in closed form at each iteration. Leveraging the ECM framework, we provide convergence guarantees, and under specific distributional assumptions, we further derive bounds on the convergence rate and the statistical error. Extensive experiments demonstrate the effectiveness of our method in recovering statistical parameters from quantized data.

1 INTRODUCTION

Quantization, which represents signals using a finite number of bits, offers a hardware-efficient approach for data storage and transmission, reducing power consumption while maintaining acceptable accuracy (Roberts, 1962; Widrow & Kollár, 2008). These properties have led to the widespread adoption of quantization in applications that require efficient processing of large-scale data.

In recommendation systems, for example, user ratings are often quantized (e.g., binary preferences or 1-5 star ratings), motivating extensive research on recovering the underlying rating matrix (Davenport et al., 2014; Cai & Zhou, 2013; Bhaskar, 2016; Bottegal & Suykens, 2017; Gao et al., 2018). In compressed sensing, quantization arises when estimating sparse signals from limited measurements, with a growing focus on one-bit and multi-bit quantized sensing (Boufounos & Baraniuk, 2008a; Zymnis et al., 2009; Plan & Vershynin, 2013; Ai et al., 2014; Shao et al., 2024). Quantization is also widely used in wireless communication and sensing, including channel estimation, array detection, and radar signal processing (Choi et al., 2016; Stöckle et al., 2016; Plabst et al., 2018; Ren & Li, 2017; Ameri et al., 2019; Jin et al., 2020; Bar-Shalom & Weiss, 2002b; Lu et al., 2024). Beyond signal recovery, quantized data naturally appear in regression and subspace learning, giving rise to quantized regression and subspace estimation problems (Mayne, 1967; Gyorfi & Wegkamp, 2008; Chen et al., 2023; Chi & Fu, 2017; Dirksen et al., 2025).

In many applications, the primary interest is not the signal itself, but its underlying statistical properties. Researchers accomplish downstream tasks such as classification and anomaly detection by estimating the mean and covariance matrix of the signal. In many applications, the primary interest is not the signal itself, but its underlying statistical properties, such as the mean (Papadopoulos et al., 2001; Dabeer & Masry, 2008) for distributed detection (Ribeiro & Giannakis, 2006a;b; Fang & Li, 2008), and the covariance matrix (Van Vleck & Middleton, 1966; Gray & Stockham, 1993; Liu & Lin, 2021; Dirksen et al., 2022; Xiao et al., 2023) for direction-of-arrival estimation (Eamaz et al., 2022; 2023; Bar-Shalom & Weiss, 2002a), power estimation in wireless sensor networks (Mo et al., 2017), spectrum sensing (Yang et al., 2025), and networked sensing (Chi & Fu, 2017).

However, quantization inevitably discards part of the original information, making the recovery of signals or their statistical properties from quantized measurements a challenging and practically

important problem. In this paper, we study parameter estimation of \mathbf{x} from quantized measurements. Given a random signal \mathbf{x} , the quantized measurement \mathbf{y} is obtained in the following way

$$\mathbf{y} = \mathcal{Q}(\mathbf{x}), \quad (1)$$

where $\mathcal{Q} : \mathbb{R}^d \rightarrow \mathcal{K}^d$ is a quantization function, or quantizer, with the set \mathcal{K} containing finite e elements, i.e., $\mathcal{K} = \{k_1, \dots, k_e\}$. The definition of the i -th element $[\mathcal{Q}(\mathbf{a})]_i$ is given by

$$[\mathcal{Q}(\mathbf{a})]_i = k_l, \quad \text{if } a_i \in [\tau_{l-1}, \tau_l), \quad (2)$$

where $\bigcup_{l=1}^e [\tau_{l-1}, \tau_l) = \mathbb{R}$ with $\tau_0 = -\infty$ and $\tau_e = \infty$. When the number of bits used to represent the quantized value, i.e., $\log_2 e$, is small, we call it a coarse quantization. The coarsest quantization method, which is also the most commonly adopted one in practice, is the one-bit quantization, i.e., $e = 2$. With the additional conditions $k_1 = -1$, $k_2 = 1$, and $\tau_1 = 0$, the quantization function \mathcal{Q} becomes the element-wise signum function.

1.1 RELATED WORKS

Parameter estimation from quantized signals has been extensively studied. Among the many directions, two fundamental lines of inquiry concern the estimation of the mean and the covariance of \mathbf{x} , which we take as our starting point. Furthermore, we will demonstrate that our proposed framework is versatile and can be extended to various related tasks, including quantized regression, quantized matrix completion, and quantized compressed sensing.

Quantized mean estimation In Papadopoulos et al. (2001), in order to estimate signals from wireless sensor networks, an optimization problem is formulated for estimating the mean under a one-dimensional Gaussian distribution assumption, based on multi-bit quantized data. To address this optimization problem, the authors propose an ECM-based algorithm. Subsequently, Ribeiro & Giannakis (2006b) focuses on the specific scenario of one-bit quantization under a Gaussian distribution and derives a closed-form expression for mean estimation in this case. Further, Fang & Li (2008) extends this line of research by introducing an adaptive threshold into the one-bit quantization function, thereby improving mean estimation accuracy. In addition, Dabeer & Masry (2008) generalizes the problem to the multidimensional setting. Ribeiro & Giannakis (2006a) extends the model distribution to the generalized Gaussian distribution under the one-dimensional assumption.

Quantized covariance/correlation estimation In Van Vleck & Middleton (1966), the authors investigate the estimation of a correlation matrix from one-bit quantized zero-mean Gaussian measurements, based on the arcsine law (Lévy, 1940). However, in such one-bit zero-mean settings, the individual variance of each dimension cannot be determined, which prevents recovery of the full covariance matrix. To estimate the variance, the “dithering technique” is introduced, which is to set a threshold in the signum function (Liu & Lin, 2021). Based on the dithering technique, the variance can be obtained in closed form (Fang & Li, 2008), but the correlation matrix should be solved analytically. Consequently, subsequent studies (Dirksen et al., 2022; Eamaz et al., 2023; Xiao et al., 2023; Liu & Chou, 2025) have proposed various strategies for correlation estimation. The estimation approaches for correlation can be classified into two categories. The first category relies on the correlation coefficient function of the one-bit Gaussian distribution, with analyses usually restricted to pairwise interactions between dimensions in the multivariate setting. Since this function is generally computationally intractable, different approximate functions have been proposed to compute the correlation coefficient (Eamaz et al., 2023; Xiao et al., 2023; Liu & Chou, 2025). In Eamaz et al. (2022; 2023), the authors introduce a known, time-varying threshold for the signum function and propose a modified arcsine law to express correlation coefficients as integrals, which are then approximated using Gauss–Legendre quadratures. The method in (Xiao et al., 2023) applies maximum likelihood (ML) estimation to compute correlation coefficients, which is equivalent to iteratively evaluating their approximate functions. The work (Liu & Chou, 2025) represents the correlation coefficient function as an infinite series via the one-bit Hermite law (Liu & Lin, 2021) and then approximates it using harmonic approximation. The second category (Dirksen et al., 2022; Chen et al., 2024) directly employs the sample covariance matrix computed using multiple thresholds. These methods are not restricted to one-bit quantization and can be generalized to multi-bit cases, though typically at the expense of reduced estimation accuracy.

1.2 CONTRIBUTION

In this paper, we study the maximum likelihood (ML) estimation of the parameters of a random signal \mathbf{x} with quantization. We assume that \mathbf{x} follows a general normal mean–variance mixture model (Barndorff-Nielsen et al., 1982), which is a flexible family encompassing several well-known distributions, including Gaussian, t , generalized hyperbolic skew- t (GHST), hyperbolic, and generalized hyperbolic (GH) distributions. Unlike previous approaches that are restricted to one-bit or Gaussian settings, our framework accommodates arbitrary bit quantization functions and any distribution within the normal mean–variance mixture family. To solve the ML estimation problem, we develop an expectation conditional maximization (ECM) algorithm. The method alternates between two steps: (i) an expectation step (E-step), where we construct a surrogate of the likelihood via Jensen’s inequality, and (ii) a conditional maximization step (CM-step), where the surrogate is maximized with respect to a subset of parameters while holding the others fixed. This procedure yields closed-form updates for the location, skewness, and scatter parameters. The proposed ECM algorithm inherits the interpretability and stability of the ECM framework, while we further establish that it converges globally at a linear rate.

2 PROBLEM FORMULATION

A random variable $\mathbf{x} \in \mathbb{R}^d$ following a normal mean-variance mixture model is represented by

$$\mathbf{x} = \boldsymbol{\mu} + z\boldsymbol{\xi} + (z\boldsymbol{\Sigma})^{\frac{1}{2}} \boldsymbol{\epsilon}, \quad (3)$$

where $\boldsymbol{\mu}$ is the location parameter, $\boldsymbol{\xi}$ is the skewness parameter, $\boldsymbol{\Sigma}$ is the scatter parameter, z is a nonnegative random variable with density function $p(z)$, and $\boldsymbol{\epsilon}$ is a standard normal random variable with mean zero and covariance matrix identity. z and $\boldsymbol{\epsilon}$ are independent from each other. The mean and the covariance matrix of \mathbf{x} are computed as $\boldsymbol{\mu} + \mathbb{E}[z]\boldsymbol{\xi}$ and $(\mathbb{E}[z^2] - \mathbb{E}[z]^2) \boldsymbol{\xi}\boldsymbol{\xi}^\top + \mathbb{E}[z]\boldsymbol{\Sigma}$, respectively. The conditional density function of \mathbf{x} given z is

$$p(\mathbf{x} | z; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{\frac{d}{2}} \det(z\boldsymbol{\Sigma})^{\frac{1}{2}}} \exp\left(-\frac{1}{2z} \|\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi}\|_{\boldsymbol{\Sigma}^{-1}}^2\right), \quad (4)$$

where $\boldsymbol{\theta} = [\boldsymbol{\mu}^\top, \boldsymbol{\xi}^\top, \text{vec}(\boldsymbol{\Sigma})^\top]^\top$ and $\|\mathbf{x}\|_{\mathbf{A}} = \mathbf{x}^\top \mathbf{A} \mathbf{x}$. Following the definitions of \mathbf{y} in (1), the density function of \mathbf{y} is given as follows:

$$p(\mathbf{y}; \boldsymbol{\theta}) = \int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^\infty p(\mathbf{x} | z; \boldsymbol{\theta}) p(z) dz d\mathbf{x}, \quad (5)$$

where $\mathcal{Q}^{-1}(\mathbf{y})$ maps the quantized data \mathbf{y} into a hyper-rectangle whose projection on each dimension i is an interval corresponding to $[\mathcal{Q}(\mathbf{y})]_i$. Given n independent and identically distributed samples $\mathbf{y}_1, \dots, \mathbf{y}_n$, the ML estimation problem is given by

$$\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \sum_{t=1}^n \log p(\mathbf{y}_t; \boldsymbol{\theta}). \quad (6)$$

When $d = 1$ and z is deterministic, i.e., \mathbf{x} is a univariate normal random variable, and one-bit quantization is applied, the estimation problem admits a closed-form solution (Ribeiro & Giannakis, 2006b). In contrast, as long as $d \geq 2$, the density function in (5) involves multiple integrals, making the optimization problem substantially more challenging.

3 THE ECM ALGORITHM

The optimization problem in (6) is challenging due to the multiple integrals in the density function (5). Nevertheless, the ECM framework can be applied by leveraging the models in (1) and (3), where \mathbf{x} and z can be naturally treated as latent variables. In the following, we introduce an ECM procedure for (6).

3.1 E-STEP

In the E-step, we derive a surrogate function for the log-likelihood function $L(\boldsymbol{\theta})$. Based on (1), the observed signal \mathbf{y} is conditioned on the hidden variable \mathbf{x} . Hence, given $[\mathbf{y}_1 \dots, \mathbf{y}_n]^\top$ and the corresponding hidden variables $[\mathbf{x}_1 \dots, \mathbf{x}_n]^\top$, we have¹

$$L(\boldsymbol{\theta}) = \sum_{t=1}^n \log \int_{\mathbb{R}} p(\mathbf{y}_t | \mathbf{x}_t; \boldsymbol{\theta}) p(\mathbf{x}_t; \boldsymbol{\theta}) d\mathbf{x}_t \geq \sum_{t=1}^n \mathbb{E}_{\mathbf{x}_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} \log p(\mathbf{y}_t, \mathbf{x}_t; \boldsymbol{\theta}) + \text{const.}, \quad (7)$$

where the inequality is based on Jensen's inequality and const. is some constant term independent of $\boldsymbol{\theta}$. Under the models specified in (1), the joint density function is given by

$$p(\mathbf{x}_t, \mathbf{y}_t; \boldsymbol{\theta}) = p(\mathbf{x}_t; \boldsymbol{\theta}) \delta(\mathcal{Q}(\mathbf{x}_t) - \mathbf{y}_t), \quad (8)$$

where $\delta(\cdot)$ denotes the multivariate Dirac delta function, defined as $\delta(\mathbf{x}_t) = \begin{cases} 0, & \text{if } \mathbf{x}_t \neq \mathbf{0} \\ +\infty, & \text{if } \mathbf{x}_t = \mathbf{0} \end{cases}$.

Since the term $\delta(\mathcal{Q}(\mathbf{x}_t) - \mathbf{y}_t)$ is independent of $\boldsymbol{\theta}$, we further obtain

$$L(\boldsymbol{\theta}) \geq \sum_{t=1}^n \mathbb{E}_{\mathbf{x}_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} \log p(\mathbf{x}_t; \boldsymbol{\theta}) + \text{const.} = \sum_{t=1}^n \mathbb{E}_{\mathbf{x}_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} \log \int_0^\infty p(\mathbf{x}_t, z_t; \boldsymbol{\theta}) dz_t + \text{const.}$$

Regarding z as a hidden variable, applying Jensen's inequality leads to

$$L(\boldsymbol{\theta}) \geq \sum_{t=1}^n \mathbb{E}_{\mathbf{x}_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [\mathbb{E}_{z_t | \mathbf{x}_t; \underline{\boldsymbol{\theta}}} \log p(\mathbf{x}_t, z_t; \boldsymbol{\theta})] + \text{const.} \triangleq S(\boldsymbol{\theta}; \underline{\boldsymbol{\theta}}). \quad (9)$$

Since $z | \mathbf{x}$ is independent of $\mathbf{x} | \mathbf{y}$, we have $p(\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}) p(z | \mathbf{x}; \boldsymbol{\theta}) = p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta})$. Given the $p(\mathbf{x} | z; \boldsymbol{\theta})$ in (4), the surrogate function $S(\boldsymbol{\theta}; \underline{\boldsymbol{\theta}})$ can be further expressed as follows:

$$S(\boldsymbol{\theta}; \underline{\boldsymbol{\theta}}) = \sum_{t=1}^n \left[\mathbb{E}_{z_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [\log p(z_t)] - \frac{1}{2} \log \det \boldsymbol{\Sigma} - \frac{1}{2} \text{Tr} \left((\mathbf{U}_t - 2\mathbf{v}_t \boldsymbol{\mu}^\top + \iota_t \boldsymbol{\mu} \boldsymbol{\mu}^\top) - 2(\mathbf{w}_t - \boldsymbol{\mu}) \boldsymbol{\xi}^\top + \zeta_t \boldsymbol{\xi} \boldsymbol{\xi}^\top \right) \boldsymbol{\Sigma}^{-1} \right] + \text{const.}, \quad (10)$$

where

$$\mathbf{U}_t = \mathbb{E}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [z_t^{-1} \mathbf{x} \mathbf{x}^\top], \quad \mathbf{v}_t = \mathbb{E}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [z_t^{-1} \mathbf{x}_t], \quad \mathbf{w}_t = \mathbb{E}_{\mathbf{x}_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [\mathbf{x}_t], \\ \iota_t = \mathbb{E}_{z_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [z_t^{-1}], \quad \zeta_t = \mathbb{E}_{z_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [z_t].$$

The details for computing these expectations are given in Appendix A. In the term $\mathbb{E}_{z_t | \mathbf{y}_t; \underline{\boldsymbol{\theta}}} [\log p(z_t)]$, some scalar parameters (such as the shape parameter ν in Student's t distribution) are contained. In the practical implementation (Galarza et al., 2021), they are typically treated as given or estimated through the one-dimensional search.

3.2 CM-STEP

Based on the surrogate function (10), in the CM-step, we solve the following optimization problem:

$$\max_{\boldsymbol{\theta}} S(\boldsymbol{\theta}; \underline{\boldsymbol{\theta}}), \quad (11)$$

where parameters $\boldsymbol{\mu}$, $\boldsymbol{\xi}$, and $\boldsymbol{\Sigma}$ can be solved with closed-form solutions:

$$\boldsymbol{\mu} = \frac{\sum_{t=1}^n (\mathbf{v}_t - \boldsymbol{\xi})}{\sum_{t=1}^n \iota_t}, \quad \boldsymbol{\xi} = \frac{\sum_{t=1}^n (\mathbf{w}_t - \boldsymbol{\mu})}{\sum_{t=1}^n \zeta_t}, \quad (12) \\ \boldsymbol{\Sigma} = \frac{1}{n} \sum_{t=1}^n \left((\mathbf{U}_t - 2\mathbf{v}_t \boldsymbol{\mu}^\top + \iota_t \boldsymbol{\mu} \boldsymbol{\mu}^\top) - 2(\mathbf{w}_t - \boldsymbol{\mu}) \boldsymbol{\xi}^\top + \zeta_t \boldsymbol{\xi} \boldsymbol{\xi}^\top \right).$$

In the context of quantization model estimation, a prototypical scenario is the one-bit Gaussian case (i.e., $e = 1$ and z is a constant). However, this model suffers from an inherent identifiability issue: for each dimension $i = 1, \dots, d$, only the ratio $\frac{\mu_i}{\sigma_i^2}$ can be determined, rather than the individual

¹Throughout this paper, underlined variables denote those whose values are given as constants.

parameters. Hence, the estimated values of μ_i and σ_i^2 differ from their desirable values by an arbitrary scaling factor. To resolve this, existing estimation methods in the one-bit Gaussian setting typically fix one parameter (either the location vector or the scatter matrix) to identify the other. Moreover, the same ambiguity persists under one-bit quantization whenever \mathbf{x} follows any elliptical distribution.

4 CONVERGENCE ANALYSIS

In this section, we analyze the convergence properties of the proposed ECM algorithm. The ECM algorithm is guaranteed to converge to a specific stationary point given an initial point (McLachlan & Krishnan, 2008). Define $\boldsymbol{\theta}^* = [\boldsymbol{\mu}^{*\top}, \boldsymbol{\xi}^{*\top}, \text{vec}(\boldsymbol{\Sigma}^*)^\top]^\top$ as a stationary point of the optimization problem in (6). We first establish the convergence of each parameter individually, while keeping the other two fixed, i.e., the convergence property of the conditional maximization step.

Proposition 1. *Denote the sequence generated by the ECM algorithm as $\{\boldsymbol{\theta}^{(k)}\}$. For any $k \geq 0$, we have*

$$\begin{aligned} \|\boldsymbol{\mu}^{(k+1)} - \boldsymbol{\mu}^*\|_2 &\leq c_\mu \|\boldsymbol{\mu}^{(k)} - \boldsymbol{\mu}^*\|_2, \\ \|\boldsymbol{\xi}^{(k+1)} - \boldsymbol{\xi}^*\|_2 &\leq c_\xi \|\boldsymbol{\xi}^{(k)} - \boldsymbol{\xi}^*\|_2, \\ \|\boldsymbol{\Sigma}^{(k+1)} - \boldsymbol{\Sigma}^*\|_F &\leq c_\Sigma \|\boldsymbol{\Sigma}^{(k)} - \boldsymbol{\Sigma}^*\|_F, \end{aligned}$$

where

$$\begin{aligned} c_\mu &= \max_{\boldsymbol{\theta}} \left\| \frac{\sum_{t=1}^n \boldsymbol{\Sigma}^{-1} \text{Cov}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z_t^{-1} \mathbf{x}_t]}{\sum_{t=1}^n \mathbb{E}_{z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z_t^{-1}]} \right\|_2 \in (0, 1), \\ c_\xi &= \max_{\boldsymbol{\theta}} \left\| \frac{\sum_{t=1}^n \boldsymbol{\Sigma}^{-1} \text{Cov}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \boldsymbol{\theta}} [\mathbf{x}_t - z_t \boldsymbol{\xi}]}{\sum_{t=1}^n \mathbb{E}_{z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z_t]} \right\|_2 \in (0, 1), \\ c_\Sigma &= \max_{\boldsymbol{\theta}} \left\| \frac{1}{2n} \sum_{t=1}^n (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{Cov}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z_t^{-1} \text{vec}(\mathbf{x}_t \mathbf{x}_t^\top)] \right\|_2 \in (0, 1). \end{aligned}$$

Based on the result in Proposition 1 and Meng & Rubin (1994), we establish that the proposed ECM algorithm converges globally at a linear rate, as detailed in the following result.

Theorem 2. *Denote the sequence generated by the ECM algorithm as $\{\boldsymbol{\theta}^{(k)}\}$. For any $k \geq 0$, we have*

$$\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2 \leq c^k \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2, \quad (13)$$

where $c = \max\{c_\mu, c_\xi, c_\Sigma\} \in (0, 1)$.

5 EXTENDED APPLICATIONS

In the previous sections, an algorithm to the ML estimation problem was proposed for modeling quantized data by a normal mean-variance mixture model. Modeling data with a probabilistic model and then performing ML estimation is a common practice in the field of machine learning. Hence, the proposed algorithm can be readily extended to typical machine learning applications from quantized data. In this section, two applications of the proposed model and estimation methods to quantized matrix completion and quantized compressive sensing will be introduced, and experimental validations will be presented in a subsequent section.

5.1 QUANTIZED MATRIX COMPLETION

The goal of the low-rank matrix completion problem (Candes & Recht, 2012) is to recover an unknown low-rank matrix $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$ from an observed, yet incomplete, matrix. Let \mathbf{X} denote a matrix whose entries are drawn from a normal mean-variance mixture distribution. We use μ_{ij} and x_{ij} to represent the (i, j) -th entries of \mathbf{M} and \mathbf{X} , respectively. Based on the normal mean-variance model, the relationship between μ_{ij} and x_{ij} is given by

$$x_{ij} = \mu_{ij} + z\xi + z^{1/2}\sigma\epsilon, \quad (14)$$

where μ_{ij} serves as the location parameter of x_{ij} . In the quantization scenario (Davenport et al., 2014), however, the matrix \mathbf{X} is not directly accessible; instead, we observe its quantized counterpart $\mathbf{Y} = \mathcal{Q}(\mathbf{X})$. Moreover, the entire matrix \mathbf{Y} is not available for observation. Let Ω denote the index set of the observed entries; specifically, if $(i, j) \in \Omega$, then Y_{ij} is observed, otherwise Y_{ij} is missing. Our goal is to recover M from incomplete \mathbf{Y} . The study of matrix completion was popularized following the Netflix Million Dollar Challenge, which posed the task: accurately predicting the values of those entries with a user–item matrix in which entries represent item ratings.

To recover M from quantized and corrupted observations, a seminal study introduced an ML framework for one-bit low-rank matrix completion (Davenport et al., 2014). Building on this work, random dithering was incorporated into the quantization function to improve recovery performance (Eamaz et al., 2024). Both approaches employed the nuclear norm to relax the low-rank constraint and used projected gradient descent for optimization. An alternative strategy reformulated the low-rank constraint via matrix factorization, followed by projected gradient descent (Bhaskar & Javanmard, 2015). Similarly, factorization combined with a majorization–minimization method was used to derive a surrogate objective, which was solved via the Gauss–Newton method (Liu et al., 2025). The framework was further extended from one-bit to multi-bit quantization using low-rank factorization together with projected gradient descent (Bhaskar, 2016).

5.2 QUANTIZED COMPRESSIVE SENSING

The one-bit compressive sensing (Boufounos & Baraniuk (2008b)) aims to estimate a sparse signal lying in a known measurement subspace based on observed quantized data (Jacques et al., 2013; Chen et al., 2024). Existing methods for this problem generally fall into three categories. The first two primarily focus on 1-bit quantization without incorporating additive noise into the original measurements (Li et al., 2018). The first category, termed regularizer-class algorithms (Laska et al., 2011), introduces additional regularization terms to the classical compressive sensing recovery problem to enforce consistency between the sparse signal and the measurements. The second category, known as penalty-class algorithms (Yan et al., 2012), models sign flips between quantized and recovered measurements as penalty terms. The third category assumes the presence of additive noise in the original measurements (Zymnis et al., 2009; Knudson et al., 2016; Shao et al., 2024). Our proposed model is developed as a further extension of the third approach. To address potential outliers in the additive noise and enhance the robustness of the model, we employ a normal mean–variance mixture model for noise modeling. Denote $\Phi \in \mathbb{R}^{d_1 \times d_2}$ as the known measurement matrix, and $\vartheta \in \mathbb{R}^{d_2}$ as the sparse signal to be estimated. Based on the normal mean-variance model, the quantized compressed sensing model is given by

$$\mathbf{y} = \mathcal{Q}(\mathbf{x}), \quad \mathbf{x} = \Phi\vartheta + z\xi + z^{1/2}\sigma\epsilon, \quad (15)$$

where $\Phi\vartheta$ serves as the location parameter of \mathbf{x} .

6 EXPERIMENTS

6.1 QUANTIZED COVARIANCE/CORRELATION ESTIMATION

In this section, we address the problem of one-bit quantized covariance estimation. Existing approaches can be broadly categorized as follows: non-dithered methods, which include the arcsine law method with a zero threshold (Van Vleck & Middleton, 1966) (Zero Threshold), and methods utilizing a non-zero threshold, such as one-bit autocorrelation estimation (Liu & Lin, 2021) (One-bit Autocorrelation), the one-bit Hermite law (Liu & Chou, 2025) (One-bit Hermite Law), and the covariance matrix recovery method (Xiao et al., 2023) (One-bit MLE). Another category is dithered methods, where the dithering signal follows either a uniform distribution (Dirksen et al., 2022) (Dithering Threshold) or a Gaussian distribution (Eamaz et al., 2022) (One-bit Time-varying); the dithering threshold can also be adaptive based on the data (Dirksen & Maly, 2024) (Adaptive Dithering). For the multi-bit problem, there are the multi-bit covariance estimator (Multi-bit Estimator) and the multi-bit parameter-free estimator (Parameter-free Estimator) (Chen & Ng, 2025).

We begin by comparing the estimation accuracy of our method against existing approaches since the previous work could only estimate the correlation matrix of one-bit data with zero mean. We generated samples from a Student’s t distribution with $\nu = 3$ degrees of freedom (Figure 2). We

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

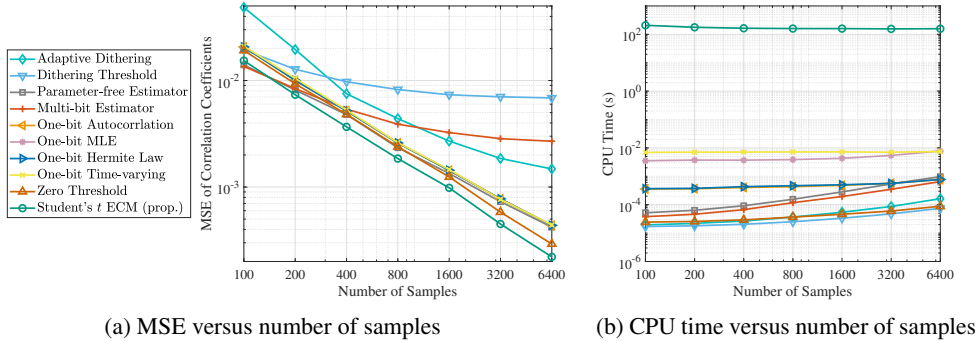


Figure 1: Algorithm performance comparison for the Student’s t distribution

estimated the correlation matrix using the following methods: the Student’s t ECM, the Gaussian ECM, the Parameter-free Estimator, and the Multi-bit Estimator. We also compared the results against those obtained by applying the EM algorithm to perform MLE on the unquantized data (Unquantized MLE). The proposed Student’s t MLE and Gaussian MLE show that as the number of bits increases, the MSE of estimating correlation will decrease and approach the result of the MLE without quantization. However, as the number of bits increases, the time required by the proposed algorithms will also increase.

Figures 3 and 4 present a comparative analysis of the convergence rates of the algorithm across four statistical distributions: Gaussian, Student’s t , GHST, and GH. The results are organized into four subfigures (a–d), each corresponding to one distribution. The convergence is measured by the distance to the optimum, defined as $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2$ for vectors \mathbf{x} and \mathbf{x}^* , and as $\|\mathbf{X}^{(k)} - \mathbf{X}^*\|_F$ for matrices \mathbf{X} and \mathbf{X}^* . All proposed methods achieve linear convergence rates. The numerical values of these rates, which vary per iteration, are detailed in Figure 4. Finally, Table 1 reports experimental results where samples are randomly generated from the GH distribution with parameters $\lambda = -0.5$, $\delta = 2$, and $\gamma = 2$, with $e = 1$ and threshold $\tau = 2$. The MSE of the correlation matrix estimated via the GH method is consistently lower than that obtained by alternative approaches.

Table 1: Correlation MSE comparison for different methods under the synthetic data

Method	10^2 samples	10^3 samples	10^4 samples
Gaussian ECM	4.7960×10^{-2}	7.5326×10^{-3}	4.9231×10^{-3}
Student’s t ECM	5.8602×10^{-2}	9.5297×10^{-3}	6.6826×10^{-3}
GHST ECM	1.5503×10^{-1}	1.0318×10^{-1}	9.5506×10^{-2}
GH ECM	4.5377×10^{-2}	7.3532×10^{-4}	9.0011×10^{-4}

6.2 QUANTIZED MATRIX COMPLETION

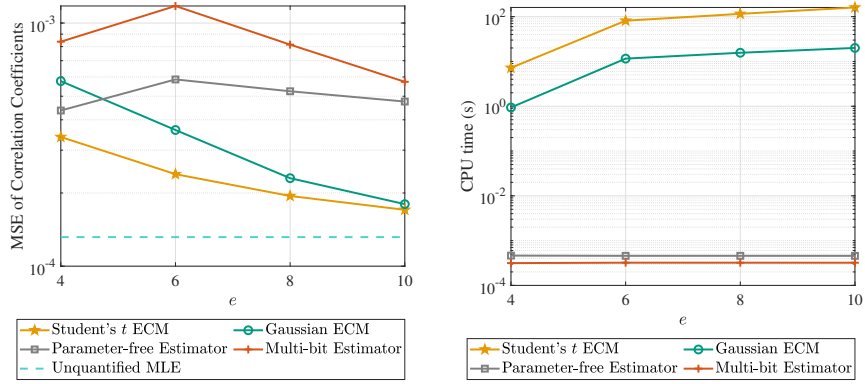
Given (14) and $\mathbf{Y} = \mathcal{Q}(\mathbf{X})$, the ML estimation problem for \mathbf{M} is given by

$$\begin{aligned} \max_{\mathbf{M}} \quad & \sum_{(i,j) \in \Omega} \log p(y_{ij} | \mu_{ij}) \\ \text{s.t.} \quad & \text{rank}(\mathbf{M}) \leq r. \end{aligned} \tag{16}$$

By applying a low-rank factorization to the matrix \mathbf{M} , we express it as $\mathbf{M} = \mathbf{A}\mathbf{B}^\top$, where $\mathbf{A} \in \mathbb{R}^{d_1 \times r}$ and $\mathbf{B} \in \mathbb{R}^{d_2 \times r}$. Through the E-step in our proposed ECM algorithm, we have the surrogate function of (16) as

$$\sum_{(i,j) \in \Omega} (\mathbf{a}_i \mathbf{b}_j^\top - e_{ij})^2, \tag{17}$$

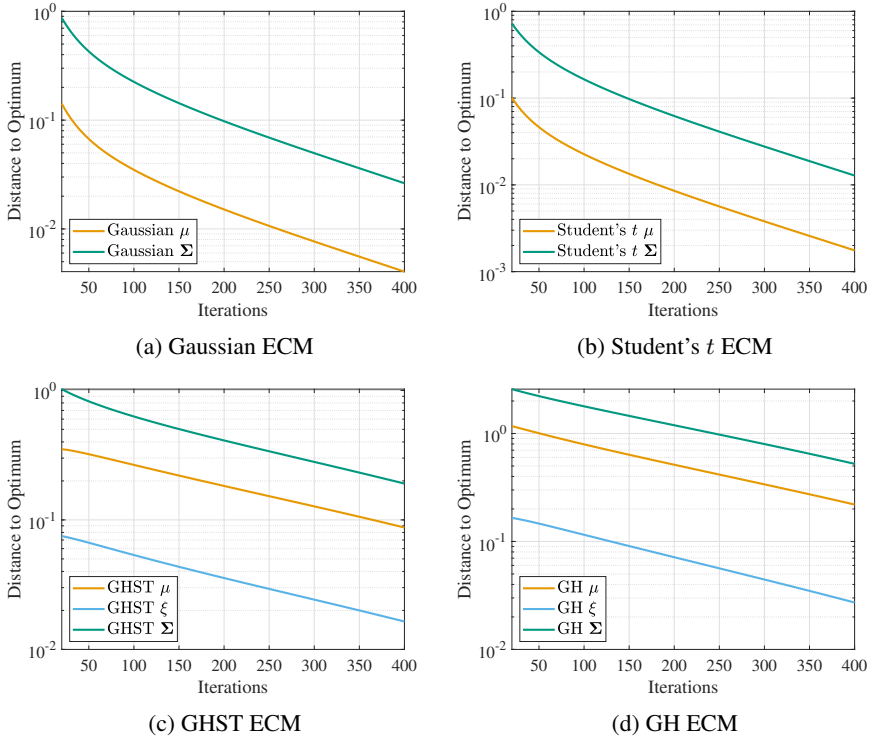
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431



(a) MSE versus e

(b) CPU time versus e

Figure 2: Algorithms performance comparison versus the number of bits



(a) Gaussian ECM

(b) Student's t ECM

(c) GHST ECM

(d) GH ECM

Figure 3: Convergence rate versus iterations

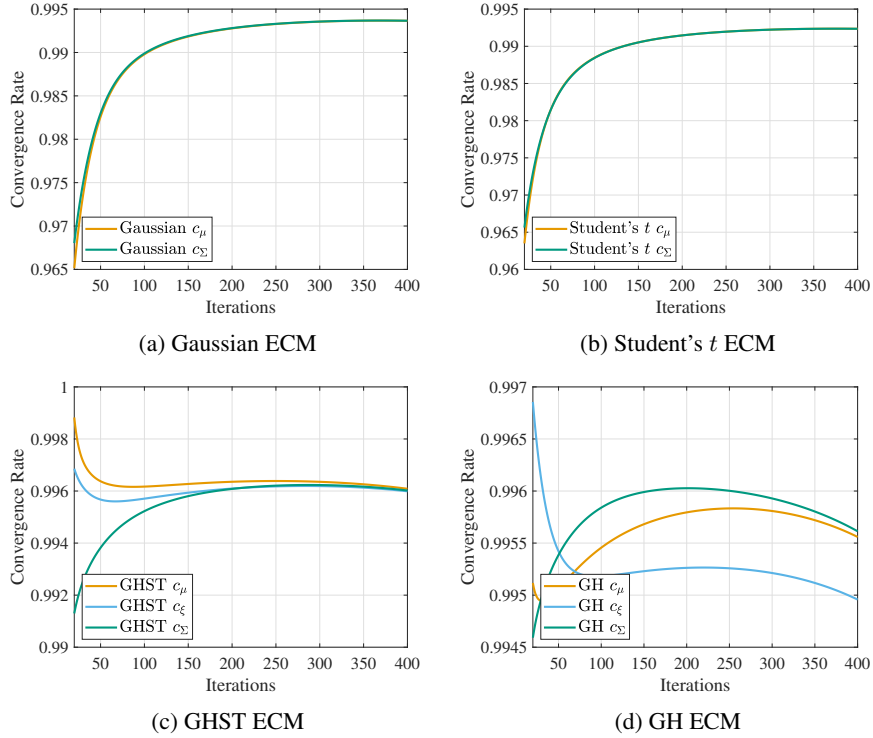


Figure 4: Convergence rate versus iterations

where \mathbf{a}_i and \mathbf{b}_i are i -th row of the matrix \mathbf{A} and \mathbf{B} , respectively, and $e_{ij} = \frac{v_{ij} - \xi}{v_{ij}}$. The details of obtaining the surrogate (17) are given in the Appendix D. Then we can use the ECM algorithm to alternatively update \mathbf{A} and \mathbf{B} .

For the experimental design, the MovieLens 1M dataset (Harper & Konstan, 2015) is employed, which contains 1000000 movie ratings provided by 6040 users for 3952 movies, with each rating ranging from 1 to 5. All ratings are randomly partitioned into training and test sets, accounting for 40% and 60% of the data, respectively. The training set is used to predict the completed rating matrix, and the entries overlapping between the completed rating matrix and the test set are then compared to evaluate prediction accuracy and root mean square error (RMSE), which are defined in Appendix E.

Existing methods encompass classic machine learning approaches for matrix completion, including the singular value decomposition (SVD (Sarwar et al., 2000)) model, the ℓ_2 -regularized matrix factorization (ℓ_2 -regularized (Paterek, 2007)) model, and the nuclear norm minimization (nuclear norm (Cai et al., 2010)) model. Furthermore, several approaches are established upon random variable model assumptions. These include: directly modeling the Gaussian without quantization (Gaussian (Candes & Plan, 2010)); 1-bit quantization with a Gaussian distribution (1-bit Gaussian (Davenport et al., 2014)); and multi-bit quantization under a Gaussian assumption (multi-bit Gaussian (Bhaskar, 2016)). Our proposed method is regarded as an extension of the multi-bit Gaussian approach by extending the Gaussian assumption to a normal mean-variance mixture model, comprising Student's t , GHST, and GH distributions.

The results are demonstrated in Table 2. As shown in the table, the proposed method consistently achieves superior performance in terms of both accuracy and RMSE. Among the evaluated models, the approach based on the GH distribution attains the best results, attributable to its flexibility and robustness in modeling skewness and heavy tails. Details regarding the experimental parameter settings are provided in the Appendix E.

Table 2: Accuracy and RMSE comparisons of matrix completion on MovieLens 1M

Method	Accuracy	RMSE	Time (s)
SVD	0.4261	0.9148	364.36
L2-regularization	0.4388	0.9355	165.08
Nuclear Norm	0.3802	1.1662	867.39
Gaussian	0.4363	0.9486	25.26
1-bit Gaussian	0.4217	0.9659	110.44
Multi-bit Gaussian	0.4453	0.9151	90.84
Multi-bit Student’s t ECM (prop.)	0.4464	0.9091	232.16
Multi-bit GHST ECM (prop.)	0.4495	0.9076	252.23
Multi-bit GH ECM (prop.)	0.4510	0.8892	312.32

6.3 QUANTIZED COMPRESSIVE SENSING

Given (15) and the ℓ_1 regularization model from Zymnis et al. (2009), we have the optimization problem

$$\max_{\vartheta} \log p(\mathbf{y} | \vartheta) + \eta \|\vartheta\|_1. \quad (18)$$

With the E-step in our proposed ECM algorithm, we have the surrogate function of (18) as

$$\frac{1}{2} \|\Phi\vartheta - e\|_2^2 + \eta \|\vartheta\|_1,$$

where the details of deriving the surrogate function are given in the Appendix D. We can use the FISTA algorithm (Beck & Teboulle, 2009) to solve the surrogate.

Performance comparisons are conducted using synthetic data, where the measurement dimension is denoted as $d_1 = 3000$ and the 30-sparse signal dimension as $d_2 = 1000$. To evaluate algorithmic robustness, an additive noise term exhibiting both skewed and heavy-tailed characteristics is introduced to the original measurements. The existing methods included in our comparison consist of restricted step shrinkage (RSS) (Laska et al., 2011), adaptive outlier pursuit (AOP) (Yan et al., 2012), and 1-bit Gaussian MLE (Zymnis et al., 2009). For each signal-to-noise ratio (SNR) level, experiments are repeated 10 times, and the average cosine similarity (Cos Sim) and computational time are reported. The results are summarized in the Table 3, which demonstrates that the 1-bit GH ECM algorithm achieves the largest average cosine similarity. The definitions of SNR and cosine similarity are given in Appendix E.

Table 3: Cosine similarity comparisons of 1-bit compressive sensing under different SNR levels.

Method	SNR = 0 dB		SNR = -5 dB		SNR = -10 dB	
	Cos Sim	Time (s)	Cos Sim	Time (s)	Cos Sim	Time (s)
RSS	0.7872	0.2639	0.6598	0.2651	0.3027	0.2546
AOP	0.7052	0.2689	0.4041	0.2951	0.2209	0.2738
1-bit Gaussian	0.8220	2.7704	0.5487	3.0875	0.2896	2.7202
1-bit Student’s t ECM (prop.)	0.8396	2.8573	0.5678	3.1643	0.3031	3.3174
1-bit GHST ECM (prop.)	0.8893	2.7435	0.7701	3.4204	0.5237	4.0281
1-bit GH ECM (prop.)	0.8915	2.7982	0.7717	3.3562	0.5407	3.9535

7 CONCLUSION AND DISCUSSION

In this paper, we propose an ECM-based algorithm for parameter estimation in quantized models. The method exhibits excellent scalability, handling problems from one-bit to multi-bit quantization and accommodating distributions ranging from Gaussian to the broader class of normal mean–variance mixtures. Experiments show that our approach yields more accurate estimates than existing methods in certain cases (e.g., the one-bit Gaussian setting), while maintaining high accuracy across a wide range of extended scenarios. Furthermore, it demonstrates strong potential in the typical machine learning tasks including quantized matrix completion and quantized compressive sensing.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

ETHICS STATEMENT

All authors of this paper have read and agree to comply with the ICLR Code of Ethics (<https://iclr.cc/public/CodeOfEthics>). This work does not involve human subjects, sensitive data, or applications that could directly cause harm. The datasets employed in our experiments are publicly available and do not contain personally identifiable information. We have taken measures to ensure that our methodology does not introduce or exacerbate unfair bias, and potential limitations as well as societal impacts are discussed in the main text. There are no conflicts of interest or external sponsorships that could influence the results or their interpretation.

REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our results. All experimental details, including model architectures, hyperparameters, training procedures, and evaluation metrics, are described in Section 6 and Appendix E. The datasets used are publicly accessible, and we include the version information in Section 6.

REFERENCES

- 594
595
596 Albert Ai, Alex Lapanowski, Yaniv Plan, and Roman Vershynin. One-bit compressed sensing with
597 non-gaussian measurements. *Linear Algebra and its Applications*, 441:222–239, 2014. 1
- 598
599 Aria Ameri, Arindam Bose, Jian Li, and Mojtaba Soltanalian. One-bit radar processing with time-
600 varying sampling thresholds. *IEEE Transactions on Signal Processing*, 67(20):5297–5308, 2019.
601 1
- 602
603 Ofer Bar-Shalom and Anthony J Weiss. DOA estimation using one-bit quantized measurements.
IEEE Transactions on Aerospace and Electronic Systems, 38(3):868–884, 2002a. 1
- 604
605 Ofer Bar-Shalom and Anthony J Weiss. DOA estimation using one-bit quantized measurements.
IEEE Transactions on Aerospace and Electronic Systems, 38(3):868–884, 2002b. 1
- 606
607 Ole Barndorff-Nielsen, John Kent, and Michael Sørensen. Normal variance-mean mixtures and z
608 distributions. *International Statistical Review*, 50(2):145, 1982. 3
- 609
610 Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse
611 problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009. 10
- 612
613 Sonia A. Bhaskar. Probabilistic low-rank matrix completion from quantized measurements. *Journal
of Machine Learning Research*, 17(60):1–34, 2016. 1, 6, 9
- 614
615 Sonia A Bhaskar and Adel Javanmard. 1-bit matrix completion under exact low-rank constraint. In
616 *2015 49th Annual Conference on Information Sciences and Systems (CISS)*, pp. 1–6. IEEE, 2015.
617 6
- 618
619 Giulio Bottegal and Johan A. K. Suykens. Probabilistic matrix factorization from quantized mea-
620 surements. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pp. 270–277,
2017. 1
- 621
622 Petros T. Boufounos and Richard G. Baraniuk. 1-bit compressive sensing. In *2008 42nd Annual
Conference on Information Sciences and Systems*, pp. 16–21, 2008a. 1
- 623
624 Petros T Boufounos and Richard G Baraniuk. 1-bit compressive sensing. In *2008 42nd Annual
625 Conference on Information Sciences and Systems*, pp. 16–21. IEEE, 2008b. 6
- 626
627 Herm Jan Brascamp and Elliott H Lieb. On extensions of the brunn-minkowski and prékopa-leindler
628 theorems, including inequalities for log concave functions, and with an application to the diffusion
equation. *Journal of functional analysis*, 22(4):366–389, 1976. 20
- 629
630 Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A singular value thresholding algorithm for
631 matrix completion. *SIAM Journal on Optimization*, 20(4):1956–1982, 2010. 9
- 632
633 Tony Cai and Wen-Xin Zhou. A max-norm constrained minimization approach to 1-bit matrix
completion. *Journal of Machine Learning Research*, 14(1):3619–3647, 2013. 1
- 634
635 Emmanuel Candes and Yaniv Plan. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):
636 925–936, 2010. 9
- 637
638 Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Com-
munications of the ACM*, 55(6):111–119, 2012. 5
- 639
640 George Casella and Roger Berger. *Statistical inference*. CRC press, 2024. 18
- 641
642 Junren Chen and Michael K Ng. A parameter-free two-bit covariance estimator with improved
operator norm error rate. *Applied and Computational Harmonic Analysis*, pp. 101774, 2025. 6
- 643
644 Junren Chen, Yueqi Wang, and Michael K Ng. Quantized low-rank multivariate regression with
645 random dithering. *IEEE Transactions on Signal Processing*, 71:3913–3928, 2023. 1
- 646
647 Junren Chen, Michael K. Ng, and Di Wang. Quantizing heavy-tailed data in statistical estimation:
(near) minimax rates, covariate quantization, and uniform recovery. *IEEE Transactions on Infor-
mation Theory*, 70(3):2003–2038, 2024. 2, 6

- 648 Yuejie Chi and Haoyu Fu. Subspace learning from bits. *IEEE Transactions on Signal Processing*,
649 65(17):4429–4442, 2017. 1
- 650
- 651 Junil Choi, Jianhua Mo, and Robert W. Heath. Near maximum-likelihood detector and channel
652 estimator for uplink multiuser massive MIMO systems with one-bit ADCs. *IEEE Transactions*
653 *on Communications*, 64(5):2005–2018, 2016. 1
- 654
- 655 Onkar Dabeer and Elias Masry. Multivariate signal parameter estimation under dependent noise
656 from 1-bit dithered quantized data. *IEEE Transactions on Information Theory*, 54(4):1637–1654,
657 2008. 1, 2
- 658 Mark A. Davenport, Yaniv Plan, Ewout Van Den Berg, and Mary Wootters. 1-bit matrix completion.
659 *Information and Inference*, 3(3):189–223, 2014. 1, 6, 9
- 660
- 661 Sjoerd Dirksen and Johannes Maly. Tuning-Free One-Bit Covariance Estimation Using Data-Driven
662 Dithering. *IEEE Transactions on Information Theory*, 70(7):5228–5247, July 2024. 6
- 663
- 664 Sjoerd Dirksen, Johannes Maly, and Holger Rauhut. Covariance estimation under one-bit quantiza-
665 tion. *The Annals of Statistics*, 50(6):3538–3562, 2022. 1, 2, 6
- 666
- 667 Sjoerd Dirksen, Weilin Li, and Johannes Maly. Subspace estimation under coarse quantization. In
668 *2025 International Conference on Sampling Theory and Applications (SampTA)*, pp. 1–5. IEEE,
2025. 1
- 669
- 670 Arian Eamaz, Farhang Yeganegi, and Mojtaba Soltanalian. Covariance recovery for one-bit sam-
671 pled non-stationary signals with time-varying sampling thresholds. *IEEE Transactions on Signal*
672 *Processing*, 70:5222–5236, 2022. 1, 2, 6
- 673
- 674 Arian Eamaz, Farhang Yeganegi, and Mojtaba Soltanalian. Covariance recovery for one-bit sampled
675 stationary signals with time-varying sampling thresholds. *Signal Processing*, 206:108899, 2023.
1, 2
- 676
- 677 Arian Eamaz, Farhang Yeganegi, and Mojtaba Soltanalian. Matrix completion from one-bit dither
678 samples. *IEEE Transactions on Signal Processing*, pp. 1–14, 2024. 6
- 679
- 680 Jun Fang and Hongbin Li. Distributed adaptive quantization for wireless sensor networks: From
681 delta modulation to maximum likelihood. *IEEE Transactions on Signal Processing*, 56(10):5246–
5257, 2008. 1, 2
- 682
- 683 Christian E. Galarza, Tsung-I Lin, Wan-Lun Wang, and Víctor H. Lachos. On moments of folded
684 and truncated multivariate student-t distributions based on recurrence relations. *Metrika*, 84(6):
685 825–850, 2021. 4
- 686
- 687 Pengzhi Gao, Ren Wang, Meng Wang, and Joe H. Chow. Low-rank matrix recovery from noisy,
688 quantized, and erroneous measurements. *IEEE Transactions on Signal Processing*, 66(11):2918–
2932, 2018. 1
- 689
- 690 Robert M. Gray and Thomas G. Stockham. Dithered quantizers. *IEEE Transactions on Information*
691 *Theory*, 39(3):805–812, 1993. 1
- 692
- 693 László Györfi and Marten Wegkamp. Quantization for nonparametric regression. *IEEE Transactions*
694 *on Information Theory*, 54(2):867–874, 2008. 1
- 695
- 696 F. Maxwell Harper and Joseph A. Konstan. The movielens datasets: History and context. *ACM*
transactions on interactive intelligent systems, 5(4):1–19, 2015. 9
- 697
- 698 Hsiu J Ho, Tsung-I Lin, Hsuan-Yu Chen, and Wan-Lun Wang. Some results on the truncated multi-
699 variate t distribution. *Journal of Statistical Planning and Inference*, 142(1):25–40, 2012. 16
- 700
- 701 Laurent Jacques, Jason N. Laska, Petros T. Boufounos, and Richard G. Baraniuk. Robust 1-bit
compressive sensing via binary stable embeddings of sparse vectors. *IEEE Transactions on In-*
formation Theory, 59(4):2082–2102, 2013. 6

- 702 Benzhou Jin, Jiang Zhu, Qihui Wu, Yuhong Zhang, and Zhiwei Xu. One-bit lfm-cw radar: Spectrum
703 analysis and target detection. *IEEE Transactions on Aerospace and Electronic Systems*, 56(4):
704 2732–2750, 2020. [1](#)
- 705
706 Karin Knudson, Rayan Saab, and Rachel Ward. One-bit compressive sensing with norm estimation.
707 *IEEE Transactions on Information Theory*, 62(5):2748–2758, 2016. [6](#)
- 708
709 Jason N Laska, Zaiwen Wen, Wotao Yin, and Richard G Baraniuk. Trust, but verify: Fast and
710 accurate signal recovery from 1-bit compressive measurements. *IEEE Transactions on Signal*
711 *Processing*, 59(11):5289–5301, 2011. [6](#), [10](#)
- 712
713 Paul Lévy. Sur certains processus stochastiques homogènes. *Compositio mathematica*, 7:283–339,
1940. [2](#)
- 714
715 Zhilin Li, Wenbo Xu, Xiaobo Zhang, and Jiaru Lin. A survey on one-bit compressed sensing:
716 Theory and applications. *Frontiers of Computer Science*, 12(2):217–230, 2018. [6](#)
- 717
718 Chun-Lin Liu and Yi-Hung Chou. Approximation and analysis of the one-bit Hermite law. In *2025*
719 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5.
IEEE, 2025. [2](#), [6](#)
- 720
721 Chun-Lin Liu and Zi-Min Lin. One-bit autocorrelation estimation with non-zero thresholds. In
722 *2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp.
723 4520–4524. IEEE, 2021. [1](#), [2](#), [6](#)
- 724
725 Xiaoqian Liu, Xu Han, Eric C. Chi, and Boaz Nadler. A majorization-minimization gauss-newton
726 method for 1-bit matrix completion. *Journal of Computational and Graphical Statistics*, 34(3):
1017–1029, 2025. [6](#)
- 727
728 Xicheng Lu, Wei Liu, and Akram Alomainy. A 1.5-bit quantization scheme and its application to
729 sparse array direction estimation. In *2024 19th International Symposium on Wireless Communi-*
730 *cation Systems (ISWCS)*, pp. 1–5, 2024. [1](#)
- 731
732 David Q. Mayne. The use of quantization in regression analysis. *International Journal of Control*,
6(6):573–577, 1967. [1](#)
- 733
734 Geoffrey J. McLachlan and T. Krishnan. *The EM algorithm and extensions*. John Wiley & Sons,
735 2008. [5](#), [18](#)
- 736
737 Xiao-Li Meng and Donald B Rubin. On the global and componentwise rates of convergence of the
em algorithm. *Linear Algebra and its Applications*, 199:413–425, 1994. [5](#), [28](#)
- 738
739 Jianhua Mo, Philip Schniter, and Robert W Heath. Channel estimation in broadband millimeter wave
740 MIMO systems with few-bit ADCs. *IEEE Transactions on Signal Processing*, 66(5):1141–1154,
2017. [1](#)
- 741
742 Haralabos C. Papadopoulos, Gregory W Wornell, and Alan V Oppenheim. Sequential signal encod-
743 ing from noisy measurements using quantizers with dynamic bias control. *IEEE Transactions on*
744 *Information Theory*, 47(3):978–1002, 2001. [1](#), [2](#)
- 745
746 Arkadiusz Paterek. Improving regularized singular value decomposition for collaborative filtering.
747 In *Proceedings of the KDD Cup and Workshop*, pp. 394–401. ACM, 2007. [9](#)
- 748
749 Daniel Plabst, Jawad Munir, Amine Mezghani, and Josef A. Nossek. Efficient non-linear equaliza-
750 tion for 1-bit quantized cyclic prefix-free massive MIMO systems. In *2018 15th International*
Symposium on Wireless Communication Systems (ISWCS), pp. 1–6, 2018. [1](#)
- 751
752 Yaniv Plan and Roman Vershynin. Robust 1-bit compressed sensing and sparse logistic regression:
753 A convex programming approach. *IEEE Transactions on Information Theory*, 59(1):482–494,
2013. [1](#)
- 754
755 Jiaying Ren and Jian Li. One-bit digital radar. In *2017 51st Asilomar Conference on Signals,*
Systems, and Computers, pp. 1142–1146, 2017. [1](#)

- 756 Alejandro Ribeiro and Georgios B Giannakis. Bandwidth-constrained distributed estimation for
757 wireless sensor networks-part II: Unknown probability density function. *IEEE Transactions on*
758 *Signal Processing*, 54(7):2784–2796, 2006a. [1](#), [2](#)
759
- 760 Alejandro Ribeiro and Georgios B Giannakis. Bandwidth-constrained distributed estimation for
761 wireless sensor networks-part I: Gaussian case. *IEEE Transactions on Signal Processing*, 54(3):
762 1131–1143, 2006b. [1](#), [2](#), [3](#)
- 763 Lawrence Roberts. Picture coding using pseudo-random noise. *IRE Transactions on Information*
764 *Theory*, 8(2):145–154, 1962. [1](#)
765
- 766 B. M. Sarwar, G. Karypis, J. A. Konstan, and J. T. Riedl. Application of dimensionality reduction in
767 recommender systems: A case study. In *WebKDD Workshop at the ACM SIGKDD*. ACM, 2000.
768 [9](#)
- 769 Mingjie Shao, Wing-Kin Ma, Junbin Liu, and Zihao Huang. Accelerated and deep expectation
770 maximization for one-bit MIMO-OFDM detection. *IEEE Transactions on Signal Processing*, 72:
771 1094–1113, 2024. [1](#), [6](#)
772
- 773 Christoph Stöckle, Jawad Munir, Amine Mezghani, and Josef A. Nossek. Channel estimation in
774 massive MIMO systems using 1-bit quantization. In *2016 IEEE 17th International Workshop on*
775 *Signal Processing Advances in Wireless Communications (SPAWC)*, pp. 1–6, 2016. [1](#)
- 776 J.H. Van Vleck and D. Middleton. The spectrum of clipped noise. *Proceedings of the IEEE*, 54(1):
777 2–19, 1966. [1](#), [2](#), [6](#)
778
- 779 Bernard Widrow and István Kollár. *Quantization noise: roundoff error in digital computation, signal*
780 *processing, control, and communications*. Cambridge University Press, 2008. [1](#)
- 781 Yu-Hang Xiao, Lei Huang, David Ramírez, Cheng Qian, and Hing Cheung So. Covariance matrix
782 recovery from one-bit data with non-zero quantization thresholds: Algorithm and performance
783 analysis. *IEEE Transactions on Signal Processing*, 71:4060–4076, 2023. [1](#), [2](#), [6](#)
784
- 785 Ming Yan, Yi Yang, and Stanley Osher. Robust 1-bit compressive sensing using adaptive outlier
786 pursuit. *IEEE Transactions on Signal Processing*, 60(7):3868–3875, 2012. [6](#), [10](#)
- 787 Jian Yang, Zihang Song, Han Zhang, and Yue Gao. Compressive spectrum sensing with 1-bit ADCs.
788 *IEEE Transactions on Vehicular Technology*, 2025. [1](#)
- 789 Argyrios Zymnis, Stephen Boyd, and Emmanuel Candes. Compressed sensing with quantized mea-
790 surements. *IEEE Signal Processing Letters*, 17(2):149–152, 2009. [1](#), [6](#), [10](#), [29](#)
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

LLM USAGE STATEMENT

During the preparation of this manuscript, we employed large language models (LLMs) as general-purpose writing and editing tools. LLMs were utilized to enhance the clarity, grammar, and structure of the text, but they did not contribute to research ideation, experimental design, or scientific content. No LLM was used to generate new ideas, mathematical proofs, or experimental results. All content, including any text suggested by LLMs, was thoroughly reviewed and verified by the authors, who assume full responsibility for the final manuscript.

A EXPECTATION EVALUATIONS ON SPECIFIC DISTRIBUTION ASSUMPTIONS

From Section 3, we have established that when the random variable \mathbf{x} in model (3) belongs to the normal mean–variance mixture family, the surrogate function (10) can be derived, along with its closed-form solutions with respect to $\boldsymbol{\mu}$, $\boldsymbol{\xi}$, and $\boldsymbol{\Sigma}$. Nevertheless, these closed-form solutions (12) depend on the expected values \mathbf{U}_t , \mathbf{v}_t , \mathbf{w}_t , ι_t , and ζ_t , which are determined by the distribution of the random variable \mathbf{x} . Therefore, in the following, we analyze the integrals corresponding to these expected values under various distributions of \mathbf{x} , and examine the convergence properties and statistical characteristics of the proposed algorithm under specific distributional assumptions. We first consider the fundamental case of the Gaussian distribution (i.e., $z = 1$ and $\boldsymbol{\xi}$). Let $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. In this case, the first- and second-order moments are given by Ho et al. (2012) as

$$\mathbb{E}_{\mathbf{x}|\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}] = \boldsymbol{\mu} + p^{-1}(\mathbf{y}; \boldsymbol{\theta}) \boldsymbol{\Sigma} \mathbf{q}, \quad (19)$$

$$\mathbb{E}_{\mathbf{x}|\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}\mathbf{x}^\top] = \boldsymbol{\mu}\boldsymbol{\mu}^\top + \boldsymbol{\Sigma} + p^{-1}(\mathbf{y}; \boldsymbol{\theta}) (2(\boldsymbol{\Sigma}\mathbf{q})\boldsymbol{\mu}^\top + \boldsymbol{\Sigma}(\mathbf{H} + \mathbf{D})\boldsymbol{\Sigma}). \quad (20)$$

Let the interval $\mathcal{Q}^{-1}(y_i)$ be $[l_i, u_i]$. The i -th elements of vector \mathbf{q} is given by

$$q_i = p(x_i = l_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) - p(x_i = u_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}), \quad (21)$$

The matrix \mathbf{H} is a matrix with all diagonal entries being zero and the (i, j) -th off-diagonal element being

$$\mathbf{H}_{ij} = h(l_i, l_j) - h(l_i, u_j) - h(u_i, l_j) + h(u_i, u_j), \quad (22)$$

with

$$h(c_i, c_j) = p(x_i = c_i, x_j = c_j, \mathbf{y}_{\setminus i,j}; \boldsymbol{\theta}).$$

The matrix \mathbf{D} is a diagonal matrix with the diagonal entries:

$$\mathbf{D}_{ii} = \frac{l_i - \mu_i}{\sigma_i^2} p(x_i = l_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) - \frac{u_i - \mu_i}{\sigma_i^2} p(x_i = u_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) - \frac{[\boldsymbol{\Sigma}\mathbf{H}]_{ii}}{\sigma_i^2}. \quad (23)$$

Now we consider the general normal mean-variance mixture with $\mathbf{x} = \boldsymbol{\mu} + z\boldsymbol{\xi} + (z\boldsymbol{\Sigma})^{\frac{1}{2}} \boldsymbol{\epsilon}$. In this case, the first and second moments satisfy

$$\mathbb{E}_{\mathbf{x},z|\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}] = \mathbb{E}_{z|\mathbf{y};\boldsymbol{\theta}}[\mathbb{E}_{\mathbf{x}|z,\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}]] \text{ and } \mathbb{E}_{\mathbf{x},z|\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}\mathbf{x}^\top] = \mathbb{E}_{z|\mathbf{y};\boldsymbol{\theta}}[\mathbb{E}_{\mathbf{x}|z,\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}\mathbf{x}^\top]].$$

Based on $\mathbf{x} | z \sim \mathcal{N}(\boldsymbol{\mu} + z\boldsymbol{\xi}, z\boldsymbol{\Sigma})$ and the moments in (19) and (20), we have

$$\mathbb{E}_{\mathbf{x}|z,\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}] = \boldsymbol{\mu} + z\boldsymbol{\xi} + p^{-1}(\mathbf{y} | z; \boldsymbol{\theta}) z \boldsymbol{\Sigma} \mathbf{q}_z, \quad (24)$$

$$\begin{aligned} \mathbb{E}_{\mathbf{x}|z,\mathbf{y};\boldsymbol{\theta}}[\mathbf{x}\mathbf{x}^\top] &= \boldsymbol{\mu}\boldsymbol{\mu}^\top + z^2\boldsymbol{\xi}\boldsymbol{\xi}^\top + 2z\boldsymbol{\mu}\boldsymbol{\xi}^\top + z\boldsymbol{\Sigma} \\ &\quad + p^{-1}(\mathbf{y} | z; \boldsymbol{\theta}) (2z(\boldsymbol{\Sigma}\mathbf{q}_z)\boldsymbol{\mu}^\top + 2z^2(\boldsymbol{\Sigma}\mathbf{q}_z)\boldsymbol{\xi}^\top + z^2\boldsymbol{\Sigma}(\mathbf{H}_z + \mathbf{D}_z)\boldsymbol{\Sigma}), \end{aligned} \quad (25)$$

where the i -th elements of vector \mathbf{q}_z is denoted as

$$q_{z,i} = p(x_i = l_i, \mathbf{y}_{\setminus i} | z; \boldsymbol{\theta}) - p(x_i = u_i, \mathbf{y}_{\setminus i} | z; \boldsymbol{\theta}), \quad (26)$$

the matrix \mathbf{H}_z is a matrix with all diagonal entries being zero and the (i, j) -th off-diagonal element being

$$\mathbf{H}_{z,ij} = h_z(l_i, l_j) - h_z(l_i, u_j) - h_z(u_i, l_j) + h_z(u_i, u_j), \quad (27)$$

with

$$h_z(c_i, c_j) = p(x_i = c_i, x_j = c_j, \mathbf{y}_{\setminus i,j} | z; \boldsymbol{\theta}),$$

and the matrix D_z is a diagonal matrix with the diagonal entries:

$$D_{z,ii} = \frac{l_i - \mu_i - z\xi_i}{\sigma_i^2} p(x_i = l_i, \mathbf{y}_{\setminus i} | z; \boldsymbol{\theta}) - \frac{u_i - \mu_i - z\xi_i}{\sigma_i^2} p(x_i = u_i, \mathbf{y}_{\setminus i} | z; \boldsymbol{\theta}) - \frac{[\boldsymbol{\Sigma} \mathbf{H}_z]_{ii}}{\sigma_i^2}. \quad (28)$$

Based on (24) and (25), we can obtain

$$\mathbf{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x}] = \boldsymbol{\mu} + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\xi} + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z p^{-1}(\mathbf{y} | z; \boldsymbol{\theta}) \boldsymbol{\Sigma} \mathbf{q}_z], \quad (29)$$

$$\begin{aligned} \mathbf{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} \mathbf{x}^\top] &= \boldsymbol{\mu} \boldsymbol{\mu}^\top + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^2] \boldsymbol{\xi} \boldsymbol{\xi}^\top + 2 \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\mu} \boldsymbol{\xi}^\top + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\Sigma} \\ &\quad + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [p^{-1}(\mathbf{y} | z; \boldsymbol{\theta}) (2z (\boldsymbol{\Sigma} \mathbf{q}_z) \boldsymbol{\mu}^\top + 2z^2 (\boldsymbol{\Sigma} \mathbf{q}_z) \boldsymbol{\xi}^\top + z^2 \boldsymbol{\Sigma} (\mathbf{H}_z + D_z) \boldsymbol{\Sigma})]. \end{aligned} \quad (30)$$

We first compute $\mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^k]$ as follows:

$$\mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^k] = \int_0^{+\infty} p(z | \mathbf{y}; \boldsymbol{\theta}) z^k dz = \int_0^{+\infty} \frac{p(\mathbf{y} | z; \boldsymbol{\theta}) p(z)}{p(\mathbf{y}; \boldsymbol{\theta})} z^k dz.$$

Here, we introduce the size-biased distribution of order k of the positive random variable z , which has the density function $p_k(z) = \frac{z^k p(z)}{\mathbf{E}_z [z^k]}$. Hence, we have

$$\mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^k] = \mathbf{E}_z [z^k] \frac{\int_0^{+\infty} p(\mathbf{y} | z; \boldsymbol{\theta}) p_k(z) dz}{p(\mathbf{y}; \boldsymbol{\theta})} = \mathbf{E}_z [z^k] \frac{p_k(\mathbf{y}; \boldsymbol{\theta})}{p(\mathbf{y}; \boldsymbol{\theta})}, \quad (31)$$

where we denote $p_k(\mathbf{y}; \boldsymbol{\theta}) = \int_0^{+\infty} p(\mathbf{y} | z; \boldsymbol{\theta}) p_k(z) dz$.

Similarly, we can obtain

$$\mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^k p^{-1}(\mathbf{y} | z; \boldsymbol{\theta}) \mathbf{q}_z] = \mathbf{E}_z [z^k] p^{-1}(\mathbf{y}; \boldsymbol{\theta}) \mathbf{q}_k, \quad (32)$$

where $\mathbf{q}_k = \int_0^{+\infty} \mathbf{q}_z p_k(z) dz$ with the i -th element

$$q_{k,i} = p_k(x_i = l_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) - p_k(x_i = u_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}). \quad (33)$$

For the expectation of $z^k \mathbf{H}_z$ we have

$$\mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^k p^{-1}(\mathbf{y} | z; \boldsymbol{\theta}) (\mathbf{H}_z + D_z)] = p^{-1}(\mathbf{y}; \boldsymbol{\theta}) (\mathbf{H}_k + D_k), \quad (34)$$

with the (i, j) entry of

$$\mathbf{H}_{k,ij} = \mathbf{E}_z [z^k] (h_k(l_i, l_j) - h_k(l_i, u_j) - h_k(u_i, l_j) + h_k(u_i, u_j)), \quad (35)$$

$$h_k(c_i, c_j) = p_k(x_i = c_i, x_j = c_j, \mathbf{y}_{\setminus i, j}; \boldsymbol{\theta}), \quad (36)$$

and the matrix D_k is a diagonal matrix with the diagonal entries:

$$\begin{aligned} D_{k,ii} &= \mathbf{E}_z [z^{k-1}] \frac{l_i - \mu_i}{\sigma_i^2} p_{k-1}(x_i = l_i, \mathbf{y}_{\setminus i} | z; \boldsymbol{\theta}) - \mathbf{E}_z [z^{k-1}] \frac{u_i - \mu_i}{\sigma_i^2} p_{k-1}(x_i = u_i, \mathbf{y}_{\setminus i} | z; \boldsymbol{\theta}) \\ &\quad - \mathbf{E}_z [z^k] \frac{\xi_i}{\sigma_i^2} q_{k,i} - \frac{[\boldsymbol{\Sigma} \mathbf{H}_k]_{ii}}{\sigma_i^2}. \end{aligned} \quad (37)$$

Therefore, the expectations in (29) and (30) become

$$\mathbf{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x}] = \boldsymbol{\mu} + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\xi} + \mathbf{E}_z [z] p^{-1}(\mathbf{y}; \boldsymbol{\theta}) \boldsymbol{\Sigma} \mathbf{q}_1, \quad (38)$$

$$\begin{aligned} \mathbf{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} \mathbf{x}^\top] &= \boldsymbol{\mu} \boldsymbol{\mu}^\top + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^2] \boldsymbol{\xi} \boldsymbol{\xi}^\top + 2 \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\mu} \boldsymbol{\xi}^\top + \mathbf{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\Sigma} \\ &\quad + p^{-1}(\mathbf{y}; \boldsymbol{\theta}) (2 \mathbf{E}_z [z] (\boldsymbol{\Sigma} \mathbf{q}_1) \boldsymbol{\mu}^\top + 2 \mathbf{E}_z [z^2] (\boldsymbol{\Sigma} \mathbf{q}_2) \boldsymbol{\xi}^\top + \boldsymbol{\Sigma} (\mathbf{H}_2 + D_2) \boldsymbol{\Sigma}). \end{aligned} \quad (39)$$

918 Meanwhile, we can also obtain

$$919 \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \mathbf{x}] = \mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [\mathbb{E}_{\mathbf{x} | z, \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \mathbf{x}]] = \mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1}] \boldsymbol{\mu} + \boldsymbol{\xi} + p^{-1}(\mathbf{y}; \boldsymbol{\theta}) \boldsymbol{\Sigma} \mathbf{q}, \quad (40)$$

$$921 \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \mathbf{x} \mathbf{x}^\top] = \mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [\mathbb{E}_{\mathbf{x} | z, \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \mathbf{x} \mathbf{x}^\top]] \\ 922 = \mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1}] \boldsymbol{\mu} \boldsymbol{\mu}^\top + \mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\xi} \boldsymbol{\xi}^\top + 2 \boldsymbol{\mu} \boldsymbol{\xi}^\top + \boldsymbol{\Sigma} \quad (41) \\ 923 + p^{-1}(\mathbf{y}; \boldsymbol{\theta}) (2(\boldsymbol{\Sigma} \mathbf{q}) \boldsymbol{\mu}^\top + 2 \mathbb{E}_z [z] (\boldsymbol{\Sigma} \mathbf{q}_1) \boldsymbol{\xi}^\top + \boldsymbol{\Sigma} (\mathbf{H}_1 + \mathbf{D}_1) \boldsymbol{\Sigma}).$$

926 B THE PROOF OF PROPOSITION 1

927
928 In this section, we prove the global linear convergence of our proposed ECM algorithm with respect
929 to the three parameters: $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$, and $\boldsymbol{\xi}$. Since the ECM algorithm can always converge to the sta-
930 tionary points (McLachlan & Krishnan, 2008), we define $\boldsymbol{\mu}^*$, $\boldsymbol{\Sigma}^*$, and $\boldsymbol{\xi}^*$ as the stationary points
931 of the optimization problem (6). In the following, we give the convergence analysis of these three
932 parameters with the convergence rate, respectively.

934 B.1 CONVERGENCE RATE OF LOCATION PARAMETER

935 Given the update rule of $\boldsymbol{\mu}$ ² in (12), we have

$$937 \boldsymbol{\mu} = \frac{\sum_{t=1}^n (\mathbf{v}_t - \boldsymbol{\xi})}{\sum_{t=1}^n \iota_t} = \frac{\sum_{t=1}^n (\mathbb{E}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z^{-1} \mathbf{x}] - \boldsymbol{\xi})}{\sum_{t=1}^n \mathbb{E}_{z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z^{-1}]}. \quad (42)$$

939 Based on (31) and (40), the update rule of $\boldsymbol{\mu}$ in (42) becomes

$$941 \boldsymbol{\mu} = \underline{\boldsymbol{\mu}} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n p^{-1}(\mathbf{y}_t; \boldsymbol{\theta}) \underline{\boldsymbol{\Sigma}} \mathbf{q}_t \quad (43)$$

944 Hence, the distance between $\boldsymbol{\mu}$ at the current iteration and $\boldsymbol{\mu}^*$ is given by

$$945 \|\boldsymbol{\mu} - \boldsymbol{\mu}^*\|_2 = \left\| \underline{\boldsymbol{\mu}} - \boldsymbol{\mu}^* + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n p^{-1}(\mathbf{y}_t; \boldsymbol{\theta}) \underline{\boldsymbol{\Sigma}} \mathbf{q}_t \right\|_2. \quad (44)$$

949 To further analyze the term on the right side in (44), we first give a result of \mathbf{q} .

950 **Lemma 3.** *When \mathbf{x} follows a normal mean-variance mixture and \mathbf{q} is defined in (21), we have*

$$951 \mathbf{q} = \nabla_{\boldsymbol{\mu}} p(\mathbf{y}; \boldsymbol{\theta}).$$

952 *Proof.* Consider the partial derivative of $p(\mathbf{y}; \boldsymbol{\theta})$ with respect to the i -th element of $\boldsymbol{\mu}$. We have that

$$953 \nabla_{\mu_i} p(\mathbf{y}; \boldsymbol{\theta}) = \nabla_{\mu_i} \int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \nabla_{\mu_i} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x},$$

954 where the interchange of the integration and differentiation operations in the second equal sign is
955 valid as the Leibniz integral rule (Casella & Berger, 2024). Since the partial derivatives of $p(\mathbf{x}; \boldsymbol{\theta})$

$$956 \nabla_{\mu_i} p(\mathbf{x} | z; \boldsymbol{\theta}) = -\nabla_{x_i} p(\mathbf{x} | z; \boldsymbol{\theta}),$$

957 we can obtain that

$$958 \nabla_{\mu_i} p(\mathbf{y}; \boldsymbol{\theta}) = - \int_0^{+\infty} \int_{\mathcal{Q}^{-1}(\mathbf{y}_{\setminus i})} \int_{l_i}^{u_i} \nabla_{x_i} p(\mathbf{x} | z; \boldsymbol{\theta}) p(z) dx_i d\mathbf{x}_{\setminus i} dz \\ 959 = \int_0^{+\infty} \int_{\mathcal{Q}^{-1}(\mathbf{y}_{\setminus i})} p(z) (p(\mathbf{x}_{\setminus i}, x_i = l_i | z; \boldsymbol{\theta}) - p(\mathbf{x}_{\setminus i}, x_i = u_i | z; \boldsymbol{\theta})) d\mathbf{x}_{\setminus i} dz \\ 960 = p(x_i = l_i, \mathbf{y}_{\setminus i}, \boldsymbol{\theta}) - p(x_i = u_i, \mathbf{y}_{\setminus i}, \boldsymbol{\theta}),$$

961 which is equivalent to the definition of the i -th element of \mathbf{q} in (21).

970 ²For simplicity, we use $\boldsymbol{\mu}$, $\boldsymbol{\xi}$, and $\underline{\boldsymbol{\Sigma}}$ to denote the parameters $\boldsymbol{\mu}^{(k)}$, $\boldsymbol{\xi}^{(k)}$, and $\boldsymbol{\Sigma}^{(k)}$ before the k -th update,
971 and use $\boldsymbol{\mu}$, $\boldsymbol{\xi}$, and $\boldsymbol{\Sigma}$ to denote the updated parameters $\boldsymbol{\mu}^{(k+1)}$, $\boldsymbol{\xi}^{(k+1)}$, and $\boldsymbol{\Sigma}^{(k+1)}$.

Based on Lemma 3, we have $p^{-1}(\mathbf{y}; \boldsymbol{\theta})\mathbf{q} = \nabla_{\boldsymbol{\mu}} \log p(\mathbf{y}; \boldsymbol{\theta})$ and $\mathbf{q}^* = \mathbf{0}$. Then the distance (44) becomes

$$\begin{aligned} \|\underline{\boldsymbol{\mu}} - \boldsymbol{\mu}^*\|_2 &= \left\| \underline{\boldsymbol{\mu}} - \boldsymbol{\mu}^* + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n \underline{\boldsymbol{\Sigma}} \nabla_{\underline{\boldsymbol{\mu}}} \log p(\mathbf{y}_t; \boldsymbol{\theta}) \right. \\ &\quad \left. - \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta}^*)}{p(\mathbf{y}_i; \boldsymbol{\theta}^*)} \right)^{-1} \sum_{t=1}^n \underline{\boldsymbol{\Sigma}} \nabla_{\underline{\boldsymbol{\mu}^*}} \log p(\mathbf{y}_t; \boldsymbol{\theta}^*) \right\|_2 \\ &= \left\| \underline{\boldsymbol{\mu}} - \boldsymbol{\mu}^* + \int_0^1 \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}})}{p(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}})} \right)^{-1} \sum_{t=1}^n \underline{\boldsymbol{\Sigma}} \nabla_{\tilde{\boldsymbol{\mu}}}^2 \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}}) (\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}^*) d\beta \right\|_2 \\ &\leq \sup_{\beta \in (0,1)} \left\| \mathbf{I} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}})}{p(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}})} \right)^{-1} \sum_{t=1}^n \underline{\boldsymbol{\Sigma}} \nabla_{\tilde{\boldsymbol{\mu}}}^2 \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}}) \right\|_2 \|\underline{\boldsymbol{\mu}} - \boldsymbol{\mu}^*\|_2, \end{aligned}$$

where the second equation is given by the mean value theorem of integrals and $\tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}} = \{\tilde{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}}, \underline{\boldsymbol{\xi}}\}$ with $\tilde{\boldsymbol{\mu}} = \beta \underline{\boldsymbol{\mu}} + (1 - \beta) \boldsymbol{\mu}^*$ and $\beta \in (0, 1]$.

To bound the term $\sup_{\beta \in (0,1)} \left\| \mathbf{I} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}})}{p(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}})} \right)^{-1} \sum_{t=1}^n \underline{\boldsymbol{\Sigma}} \nabla_{\tilde{\boldsymbol{\mu}}}^2 \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\mu}}) \right\|_2$, we need some results for the second order derivative of $\log p(\mathbf{y}_t; \boldsymbol{\theta})$ at first.

Lemma 4. *The second order derivative of $\log p(\mathbf{y}; \boldsymbol{\theta})$ with respect to $\boldsymbol{\mu}$, i.e., $\nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$, satisfy*

1. $\nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$ is a negative definite matrix;
2. $\mathbf{I} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta})$ is a positive definite matrix,

for all $\boldsymbol{\theta}$ in the feasible set.

Proof. The second order derivative of $\log p(\mathbf{y}; \boldsymbol{\theta})$ with respect to $\boldsymbol{\mu}$ is given as

$$\nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\mu}} (p^{-1}(\mathbf{y}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\mu}} p(\mathbf{y}; \boldsymbol{\theta})) = \frac{\nabla_{\boldsymbol{\mu}}^2 p(\mathbf{y}; \boldsymbol{\theta}) \cdot p(\mathbf{y}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\mu}} p(\mathbf{y}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\mu}}^{\top} p(\mathbf{y}; \boldsymbol{\theta})}{p^2(\mathbf{y}; \boldsymbol{\theta})}. \quad (45)$$

To analyze the above expression, we should compute the first and second order derivatives of $p(\mathbf{y}; \boldsymbol{\theta})$ with respect to $\boldsymbol{\mu}$ at first. The first order one is given as follows:

$$\nabla_{\boldsymbol{\mu}} p(\mathbf{y}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\mu}} \int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \nabla_{\boldsymbol{\mu}} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}. \quad (46)$$

Since $p(\mathbf{x}; \boldsymbol{\theta}) = \int_0^{+\infty} p(z) p(\mathbf{x} | z; \boldsymbol{\theta}) dz$, (46) becomes

$$\begin{aligned} \nabla_{\boldsymbol{\mu}} p(\mathbf{y}; \boldsymbol{\theta}) &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} p(z) \nabla_{\boldsymbol{\mu}} p(\mathbf{x} | z; \boldsymbol{\theta}) dz d\mathbf{x} \\ &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \frac{p(z)}{z} p(\mathbf{x} | z; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu} - z \boldsymbol{\xi}) dz d\mathbf{x}. \end{aligned}$$

Then the second order one can be further computed as

$$\begin{aligned} \nabla_{\boldsymbol{\mu}}^2 p(\mathbf{y}; \boldsymbol{\theta}) &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \frac{p(z)}{z} \nabla_{\boldsymbol{\mu}} (p(\mathbf{x} | z; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu} - z \boldsymbol{\xi})) dz d\mathbf{x} \\ &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \boldsymbol{\Sigma}^{-1} p(\mathbf{x} | z; \boldsymbol{\theta}) p(z) (-z^{-1} \mathbf{I} + z^{-2} (\mathbf{x} - \boldsymbol{\mu} - z \boldsymbol{\xi})(\mathbf{x} - \boldsymbol{\mu} - z \boldsymbol{\xi})^{\top} \boldsymbol{\Sigma}^{-1}) dz d\mathbf{x}. \end{aligned} \quad (47)$$

Lemma 5. For a function $f(\mathbf{x}, z)$, if $\int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x} | z; \boldsymbol{\theta}) f(\mathbf{x}, z) d\mathbf{x}$ is integrable, we have

$$\int_{\mathbb{R}^d} \int_0^{+\infty} p(z) p(\mathbf{x}, \mathbf{y} | z; \boldsymbol{\theta}) f(\mathbf{x}, z) dz d\mathbf{x} = p(\mathbf{y}; \boldsymbol{\theta}) \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [f(\mathbf{x}, z)].$$

Proof. Given that

$$\int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x} | z; \boldsymbol{\theta}) f(\mathbf{x}, z) d\mathbf{x} = \int_{\mathbb{R}^d} p(\mathbf{x}, \mathbf{y} | z; \boldsymbol{\theta}) f(\mathbf{x}, z) d\mathbf{x},$$

we can obtain that

$$\begin{aligned} \int_{\mathbb{R}^d} \int_0^{+\infty} p(z) p(\mathbf{x}, \mathbf{y} | z; \boldsymbol{\theta}) f(\mathbf{x}, z) dz d\mathbf{x} &= p(\mathbf{y}; \boldsymbol{\theta}) \int_{\mathbb{R}^d} \int_0^{+\infty} p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}) f(\mathbf{x}, z) dz d\mathbf{x} \\ &= p(\mathbf{y}; \boldsymbol{\theta}) \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [f(\mathbf{x}, z)]. \end{aligned}$$

□

Based on Lemma 5, the first and second order derivatives of $p(\mathbf{y}; \boldsymbol{\theta})$ becomes

$$\nabla_{\boldsymbol{\mu}} p(\mathbf{y}; \boldsymbol{\theta}) = p(\mathbf{y}; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1}(\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi})],$$

and

$$\nabla_{\boldsymbol{\mu}}^2 p(\mathbf{y}; \boldsymbol{\theta}) = p(\mathbf{y}; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} (-\mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1}] \mathbf{I} + \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-2}(\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi})(\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi})^\top] \boldsymbol{\Sigma}^{-1}).$$

Then the second order derivative of $\log p(\mathbf{y}; \boldsymbol{\theta})$ in (45) becomes

$$\begin{aligned} \nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}; \boldsymbol{\theta}) &= -\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-1} \mathbb{E}_{\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}} [(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbb{E}_{\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} - \boldsymbol{\mu}] \mathbb{E}_{\mathbf{x} | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} - \boldsymbol{\mu}]^\top \boldsymbol{\Sigma}^{-1} \\ &= \boldsymbol{\Sigma}^{-1} (-\mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1}] \boldsymbol{\Sigma} + \text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \mathbf{x}]) \boldsymbol{\Sigma}^{-1}. \end{aligned} \quad (48)$$

To prove the negative definite of $\nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$, we introduce the Brascamp–Lieb inequality.

Lemma 6 (Brascamp & Lieb (1976)). Consider a probability density function $p(\mathbf{x})$ which is log-concave to \mathbf{x} . The Brascamp–Lieb inequality is given by

$$\text{Cov}_{\tilde{\mathbf{x}}}(f(\mathbf{x})) \preceq \mathbb{E}_{\tilde{\mathbf{x}}} [\nabla_{\mathbf{x}}^\top f(\mathbf{x}) [\nabla_{\mathbf{x}}^2 (-\log p(\tilde{\mathbf{x}}))]^{-1} \nabla_{\mathbf{x}} f(\mathbf{x})],$$

where the equality is obtained when $\log p(\tilde{\mathbf{x}})$ is linear with respect to \mathbf{x} .

Let $\tilde{\mathbf{x}} = \mathbf{x}, z | \mathbf{y}$ and $f(\mathbf{x}) = z^{-1} \mathbf{x}$. Since $\nabla_{\tilde{\mathbf{x}}}^2 \log p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}) = -\frac{\boldsymbol{\Sigma}^{-1}}{z}$, $p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta})$ is a log-concave function with respect to \mathbf{x} . Based on the Brascamp–Lieb inequality in Lemma 6, since $\log p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta})$ is non-linear to \mathbf{x} , we have

$$\text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \mathbf{x}] \prec \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} z \boldsymbol{\Sigma} z^{-1}] = \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1}] \boldsymbol{\Sigma},$$

and hence $\nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$ is negative definite. Substituting $\nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$ in (48) into $\mathbf{I} +$

$\left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})}\right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta})$, we have

$$\mathbf{I} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})}\right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta}) = \frac{\sum_{t=1}^n \boldsymbol{\Sigma}^{-1} \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \boldsymbol{\theta}} [z^{-1} \mathbf{x}]}{\sum_{t=1}^n \mathbb{E}_{z | \mathbf{y}_t; \boldsymbol{\theta}} [z^{-1}]}.$$

Since both $\boldsymbol{\Sigma}$ and $\text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \mathbf{x}_t]$ are positive definite matrices, $\mathbf{I} +$

$\left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})}\right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\mu}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta})$ is positive definite. □

Based on Lemma 4, we have that all the eigenvalues of matrix $\mathbf{I} + \Sigma \nabla_{\mu}^2 \log p(\mathbf{y}; \theta)$ for any θ are belongs to $(0, 1)$, hence we have

$$\begin{aligned} & \sup_{\beta \in (0,1]} \left\| \mathbf{I} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \tilde{\theta}_{\mu})}{p(\mathbf{y}_i; \tilde{\theta}_{\mu})} \right)^{-1} \sum_{t=1}^n \underline{\Sigma} \nabla_{\mu}^2 \log p(\mathbf{y}_t; \tilde{\theta}_{\mu}) \right\|_2 \\ &= \sup_{\beta \in (0,1]} \left\| \frac{\sum_{t=1}^n \underline{\Sigma}^{-1} \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \tilde{\theta}_{\mu}} [z^{-1} \mathbf{x}]}{\sum_{t=1}^n \mathbb{E}_{z | \mathbf{y}_t; \tilde{\theta}_{\mu}} [z^{-1}]} \right\|_2 \\ &\leq \max_{\theta} \left\| \frac{\sum_{t=1}^n \underline{\Sigma}^{-1} \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \theta} [z^{-1} \mathbf{x}]}{\sum_{t=1}^n \mathbb{E}_{z | \mathbf{y}_t; \theta} [z^{-1}]} \right\|_2 \\ &\triangleq c_{\mu} \in (0, 1). \end{aligned}$$

□

B.2 CONVERGENCE RATE OF SKEWNESS PARAMETER

Given the update rule of ξ in (12), we have

$$\xi = \frac{\sum_{t=1}^n (\mathbf{w}_t - \mu)}{\sum_{t=1}^n \zeta_t} = \frac{\sum_{t=1}^n (\mathbb{E}_{\mathbf{x} | \mathbf{y}_t; \theta} [\mathbf{x}] - \mu)}{\sum_{t=1}^n \mathbb{E}_{z_t | \mathbf{y}_t; \theta} [z]}. \quad (49)$$

Based on (38), we have

$$\mathbb{E}_{\mathbf{x} | \mathbf{y}; \theta} [\mathbf{x}] - \mu = \mathbb{E}_{z | \mathbf{y}; \theta} [z] \xi + p^{-1}(\mathbf{y}; \theta) \Sigma \mathbf{q}_1.$$

Hence, the update rule of ξ in (49) becomes

$$\xi = \underline{\xi} + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \theta)}{p(\mathbf{y}_i; \theta)} \right)^{-1} \sum_{t=1}^n p^{-1}(\mathbf{y}_t; \theta) \underline{\Sigma} \mathbf{q}_{1,t} \quad (50)$$

Hence, the distance between ξ at the current iteration and ξ^* is given by

$$\|\xi - \xi^*\|_2 = \left\| \underline{\xi} - \xi^* + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \theta)}{p(\mathbf{y}_i; \theta)} \right)^{-1} \sum_{t=1}^n p^{-1}(\mathbf{y}_t; \theta) \underline{\Sigma} \mathbf{q}_{1,t} \right\|_2. \quad (51)$$

To further analyze the term on the right side in (51), we first give a result of \mathbf{q}_1 .

Lemma 7. When \mathbf{x} follows a normal mean-variance mixture and \mathbf{q}_1 is defined in (21), we have

$$\mathbf{q}_1 = \nabla_{\xi} p(\mathbf{y}; \theta).$$

Proof. Consider the partial derivative of $p(\mathbf{y}; \theta)$ with respect to the i -th element of ξ . We have that

$$\nabla_{\xi_i} p(\mathbf{y}; \theta) = \nabla_{\xi_i} \int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x}; \theta) d\mathbf{x} = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \nabla_{\xi_i} p(\mathbf{x}; \theta) d\mathbf{x}.$$

Since the partial derivatives of $p(\mathbf{x}; \theta)$

$$\nabla_{\xi_i} p(\mathbf{x} | z; \theta) = -z \nabla_{x_i} p(\mathbf{x} | z; \theta),$$

we can obtain that

$$\begin{aligned} \nabla_{\xi_i} p(\mathbf{y}; \theta) &= - \int_0^{+\infty} \int_{\mathcal{Q}^{-1}(\mathbf{y}_{\setminus i})} \int_{l_i}^{u_i} z \nabla_{x_i} p(\mathbf{x} | z; \theta) p(z) dx_i d\mathbf{x}_{\setminus i} dz \\ &= \int_0^{+\infty} \int_{\mathcal{Q}^{-1}(\mathbf{y}_{\setminus i})} \frac{p(z)}{z} z (p(\mathbf{x}_{\setminus i}, x_i = l_i | z; \theta) - p(\mathbf{x}_{\setminus i}, x_i = u_i | z; \theta)) d\mathbf{x}_{\setminus i} dz \\ &= p(x_i = l_i, \mathbf{y}_{\setminus i} | z, \theta) - p(x_i = u_i, \mathbf{y}_{\setminus i} | z, \theta). \end{aligned}$$

which is equivalent to the definition of the i -th element of \mathbf{q}_1 from (33).

Based on Lemma 7, we have $p^{-1}(\mathbf{y}; \boldsymbol{\theta})\mathbf{q}_1 = \nabla_{\boldsymbol{\xi}} \log p(\mathbf{y}; \boldsymbol{\theta})$ and $\mathbf{q}_1^* = \mathbf{0}$. Then the distance (44) becomes

$$\begin{aligned} \|\underline{\boldsymbol{\xi}} - \boldsymbol{\xi}^*\|_2 &= \left\| \underline{\boldsymbol{\xi}} - \boldsymbol{\xi}^* + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n \underline{\boldsymbol{\Sigma}} \nabla_{\underline{\boldsymbol{\xi}}} \log p(\mathbf{y}; \boldsymbol{\theta}) \right. \\ &\quad \left. - \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \boldsymbol{\theta}^*)}{p(\mathbf{y}_i; \boldsymbol{\theta}^*)} \right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma}^* \nabla_{\boldsymbol{\xi}^*} \log p(\mathbf{y}; \boldsymbol{\theta}^*) \right\|_2 \\ &= \left\| \underline{\boldsymbol{\xi}} - \boldsymbol{\xi}^* + \int_0^1 \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}})}{p(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}})} \right)^{-1} \sum_{t=1}^n \tilde{\boldsymbol{\Sigma}} \nabla_{\tilde{\boldsymbol{\xi}}}^2 \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}}) (\tilde{\boldsymbol{\xi}} - \boldsymbol{\xi}^*) d\beta \right\|_2 \\ &\leq \sup_{\beta \in (0,1]} \left\| \mathbf{I} + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}})}{p(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}})} \right)^{-1} \sum_{t=1}^n \tilde{\boldsymbol{\Sigma}} \nabla_{\tilde{\boldsymbol{\xi}}}^2 \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}}) \right\|_2 \|\underline{\boldsymbol{\xi}} - \boldsymbol{\xi}^*\|_2, \end{aligned}$$

where the second equation is given by the mean value theorem of integrals and $\tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}} = \{\underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}}, \tilde{\boldsymbol{\xi}}\}$ with $\tilde{\boldsymbol{\xi}} = \beta \underline{\boldsymbol{\xi}} + (1 - \beta)\boldsymbol{\xi}^*$, and $\beta \in (0, 1]$.

To bound the term $\sup_{\beta \in (0,1]} \left\| \mathbf{I} + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}})}{p(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}})} \right)^{-1} \sum_{t=1}^n \tilde{\boldsymbol{\Sigma}} \nabla_{\tilde{\boldsymbol{\xi}}}^2 \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\boldsymbol{\xi}}) \right\|_2$, we need some results for the second order derivative of $\log p(\mathbf{y}_t; \boldsymbol{\theta})$ at first.

Lemma 8. *The second order derivative of $\log p(\mathbf{y}; \boldsymbol{\theta})$ with respect to $\boldsymbol{\xi}$, i.e., $\nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$, satisfy*

1. $\nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$ is a negative definite matrix;
2. $\mathbf{I} + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta})$ is a positive definite matrix.

Proof. The second order derivative of $\log p(\mathbf{y}; \boldsymbol{\theta})$ with respect to $\boldsymbol{\xi}$ is given as

$$\nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\xi}} (p^{-1}(\mathbf{y}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\xi}} p(\mathbf{y}; \boldsymbol{\theta})) = \frac{\nabla_{\boldsymbol{\xi}}^2 p(\mathbf{y}; \boldsymbol{\theta}) \cdot p(\mathbf{y}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\xi}} p(\mathbf{y}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\xi}}^{\top} p(\mathbf{y}; \boldsymbol{\theta})}{p^2(\mathbf{y}; \boldsymbol{\theta})}. \quad (52)$$

To analyze the above expression, we should compute the first and second order derivatives of $p(\mathbf{y}; \boldsymbol{\theta})$ with respect to $\boldsymbol{\xi}$ at first. The first order one is given as follows:

$$\nabla_{\boldsymbol{\xi}} p(\mathbf{y}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\xi}} \int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \nabla_{\boldsymbol{\xi}} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}. \quad (53)$$

Since $p(\mathbf{x}; \boldsymbol{\theta}) = \int_0^{+\infty} p(z) p(\mathbf{x} | z; \boldsymbol{\theta}) dz$, (53) becomes

$$\begin{aligned} \nabla_{\boldsymbol{\xi}} p(\mathbf{y}; \boldsymbol{\theta}) &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} p(z) \nabla_{\boldsymbol{\xi}} p(\mathbf{x} | z; \boldsymbol{\theta}) dz d\mathbf{x} \\ &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} p(z) p(\mathbf{x} | z; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi}) dz d\mathbf{x}. \end{aligned}$$

Then the second order derivative can be further computed as

$$\begin{aligned} \nabla_{\boldsymbol{\xi}}^2 p(\mathbf{y}; \boldsymbol{\theta}) &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} p(z) \nabla_{\boldsymbol{\xi}} (p(\mathbf{x} | z; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi})) dz d\mathbf{x} \\ &= \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \boldsymbol{\Sigma}^{-1} p(\mathbf{x} | z; \boldsymbol{\theta}) p(z) (-z\mathbf{I} + (\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi})(\mathbf{x} - \boldsymbol{\mu} - z\boldsymbol{\xi})^{\top} \boldsymbol{\Sigma}^{-1}) dz d\mathbf{x}. \end{aligned} \quad (54)$$

Based on Lemma 5, the derivatives of $p(\mathbf{y}; \boldsymbol{\theta})$ becomes

$$\nabla_{\boldsymbol{\xi}} p(\mathbf{y}; \boldsymbol{\theta}) = p(\mathbf{y}; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} - \boldsymbol{\mu} - z \boldsymbol{\xi}],$$

and

$$\nabla_{\boldsymbol{\xi}}^2 p(\mathbf{y}; \boldsymbol{\theta}) = p(\mathbf{y}; \boldsymbol{\theta}) \boldsymbol{\Sigma}^{-1} (-\mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \mathbf{I} + \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [(\mathbf{x} - \boldsymbol{\mu} - z \boldsymbol{\xi})(\mathbf{x} - \boldsymbol{\mu} - z \boldsymbol{\xi})^\top] \boldsymbol{\Sigma}^{-1}).$$

Then the second order derivative of $\log p(\mathbf{y}; \boldsymbol{\theta})$ in (52) becomes

$$\nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}; \boldsymbol{\theta}) = \boldsymbol{\Sigma}^{-1} (-\mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\Sigma} + \text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} - z \boldsymbol{\xi}]) \boldsymbol{\Sigma}^{-1}. \quad (55)$$

Setting $\tilde{\mathbf{x}} = \mathbf{x}, z | \mathbf{y}$ and $f(\mathbf{x}) = \mathbf{x} - z \boldsymbol{\xi}$, since $\nabla_{\tilde{\mathbf{x}}}^2 \log p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}) = -\frac{\boldsymbol{\Sigma}^{-1}}{z}$, $p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta})$ is a log-concave function with respect to \mathbf{x} . Based on the Brascamp–Lieb inequality in Lemma 6, we have

$$\text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} - z \boldsymbol{\xi}] \prec \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z \boldsymbol{\Sigma}] = \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z] \boldsymbol{\Sigma},$$

and hence $\nabla_{\boldsymbol{\xi}}^2 p(\mathbf{y}; \boldsymbol{\theta})$ is negative definite.

Substituting $\nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$ in (55) into $\mathbf{I} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta})$, we have

$$\mathbf{I} + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta}) = \frac{\sum_{t=1}^n \boldsymbol{\Sigma}^{-1} \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \boldsymbol{\theta}} [\mathbf{x} - z \boldsymbol{\xi}]}{\sum_{t=1}^n \mathbb{E}_{z | \mathbf{y}_t; \boldsymbol{\theta}} [z]}.$$

Since both $\boldsymbol{\Sigma}$ and $\text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [\mathbf{x} - z \boldsymbol{\xi}]$ are positive definite matrices, $\mathbf{I} + \left(\sum_{i=1}^n \frac{p_{-1}(\mathbf{y}_i; \boldsymbol{\theta})}{p(\mathbf{y}_i; \boldsymbol{\theta})} \right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}_t; \boldsymbol{\theta})$ is positive definite. \square

Based on Lemma 8, we have that all the eigenvalues of matrix $\mathbf{I} + \boldsymbol{\Sigma} \nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}; \boldsymbol{\theta})$ for any $\boldsymbol{\theta}$ are belong to $(0, 1)$, hence we have

$$\begin{aligned} & \sup_{\beta \in (0,1)} \left\| \mathbf{I} + \left(\sum_{i=1}^n \frac{p_1(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\xi})}{p(\mathbf{y}_i; \tilde{\boldsymbol{\theta}}_{\xi})} \right)^{-1} \sum_{t=1}^n \boldsymbol{\Sigma} \nabla_{\boldsymbol{\xi}}^2 \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\xi}) \right\|_2 \\ &= \sup_{\beta \in (0,1)} \left\| \frac{\sum_{t=1}^n \boldsymbol{\Sigma}^{-1} \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\xi}} [\mathbf{x} - z \boldsymbol{\xi}]}{\sum_{t=1}^n \mathbb{E}_{z | \mathbf{y}_t; \tilde{\boldsymbol{\theta}}_{\xi}} [z]} \right\|_2 \\ &\leq \max_{\boldsymbol{\theta}} \left\| \frac{\sum_{t=1}^n \boldsymbol{\Sigma}^{-1} \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \boldsymbol{\theta}} [\mathbf{x} - z \boldsymbol{\xi}]}{\sum_{t=1}^n \mathbb{E}_{z | \mathbf{y}_t; \boldsymbol{\theta}} [z]} \right\|_2 \\ &\triangleq c_{\xi} \in (0, 1). \end{aligned}$$

\square

B.3 CONVERGENCE RATE OF SCATTER MATRIX

Given the update rule of $\boldsymbol{\Sigma}$ in (12), we have

$$\begin{aligned} \boldsymbol{\Sigma} &= \frac{1}{n} \sum_{t=1}^n \left((\mathbf{U}_t - 2\mathbf{v}_t \boldsymbol{\mu}^\top + \iota_t \boldsymbol{\mu} \boldsymbol{\mu}^\top) - 2(\mathbf{w}_t - \boldsymbol{\mu}) \boldsymbol{\xi}^\top + \zeta_t \boldsymbol{\xi} \boldsymbol{\xi}^\top \right) \\ &= \frac{1}{n} \sum_{t=1}^n \left(\mathbb{E}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z^{-1} \mathbf{x} \mathbf{x}^\top] - 2\mathbb{E}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z^{-1} \mathbf{x}_t] \boldsymbol{\mu}^\top + \mathbb{E}_{z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z^{-1}] \boldsymbol{\mu} \boldsymbol{\mu}^\top \right. \\ &\quad \left. - 2(\mathbb{E}_{\mathbf{x} | \mathbf{y}_t; \boldsymbol{\theta}} [\mathbf{x}] - \boldsymbol{\mu}) \boldsymbol{\xi}^\top + \mathbb{E}_{z_t | \mathbf{y}_t; \boldsymbol{\theta}} [z] \boldsymbol{\xi} \boldsymbol{\xi}^\top \right). \end{aligned} \quad (56)$$

In the following, we derive all five terms in the summation in (56).

1242 **Term 1:** Based on (41), we have

$$1243 \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \underline{\theta}} [z^{-1} \mathbf{x} \mathbf{x}^\top] = \mathbb{E}_{z | \mathbf{y}; \underline{\theta}} [z^{-1}] \underline{\boldsymbol{\mu}} \underline{\boldsymbol{\mu}}^\top + \mathbb{E}_{z | \mathbf{y}; \underline{\theta}} [z] \underline{\boldsymbol{\xi}} \underline{\boldsymbol{\xi}}^\top + 2 \underline{\boldsymbol{\mu}} \underline{\boldsymbol{\xi}}^\top + \underline{\boldsymbol{\Sigma}} \\ 1244 + p^{-1}(\mathbf{y}; \underline{\theta}) \left(2(\underline{\boldsymbol{\Sigma}} \mathbf{q}) \underline{\boldsymbol{\mu}}^\top + 2 \mathbb{E}_z [z] (\underline{\boldsymbol{\Sigma}} \mathbf{q}_1) \underline{\boldsymbol{\xi}}^\top + \underline{\boldsymbol{\Sigma}} (\mathbf{H}_1 + \mathbf{D}_1) \underline{\boldsymbol{\Sigma}} \right), \\ 1245$$

1246 where \mathbf{q} , \mathbf{q}_1 , \mathbf{H}_1 and \mathbf{D}_1 are defined in (21), (32), (34) and (37), respectively.

1247 **Term 2:** Based on (40), we can obtain

$$1248 -2 \mathbb{E}_{\mathbf{x}_t, z_t | \mathbf{y}_t; \underline{\theta}} [z_t^{-1} \mathbf{x}_t] \underline{\boldsymbol{\mu}}^\top = -2 \left(\mathbb{E}_{z | \mathbf{y}; \underline{\theta}} [z^{-1}] \underline{\boldsymbol{\mu}} \underline{\boldsymbol{\mu}}^\top + \underline{\boldsymbol{\xi}} + p^{-1}(\mathbf{y}; \underline{\theta}) \underline{\boldsymbol{\Sigma}} \mathbf{q} \right) \underline{\boldsymbol{\mu}}^\top. \\ 1249$$

1250 **Term 4:** Based on (38), we have

$$1251 -2 \left(\mathbb{E}_{\mathbf{x}_t | \mathbf{y}_t; \underline{\theta}} [\mathbf{x}_t] - \underline{\boldsymbol{\mu}} \right) \underline{\boldsymbol{\xi}}^\top = -2 \mathbb{E}_{z | \mathbf{y}; \underline{\theta}} [z] \underline{\boldsymbol{\xi}} \underline{\boldsymbol{\xi}}^\top - 2 p^{-1}(\mathbf{y}; \underline{\theta}) (\underline{\boldsymbol{\Sigma}} \mathbf{q}_1) \underline{\boldsymbol{\xi}}^\top \\ 1252$$

1253 Upon cancellation of the opposing terms, the update rule in (56) reduces to:

$$1254 \underline{\boldsymbol{\Sigma}} = \underline{\boldsymbol{\Sigma}} + \frac{1}{n} \sum_{t=1}^n p^{-1}(\mathbf{y}_t; \underline{\theta}) \underline{\boldsymbol{\Sigma}} (\mathbf{H}_{1,t} + \mathbf{D}_{1,t}) \underline{\boldsymbol{\Sigma}}, \quad (57) \\ 1255$$

1256 Hence, the distance between $\underline{\boldsymbol{\Sigma}}$ at the current iteration and $\underline{\boldsymbol{\Sigma}}^*$ is given by

$$1257 \|\underline{\boldsymbol{\Sigma}} - \underline{\boldsymbol{\Sigma}}^*\|_F = \left\| \underline{\boldsymbol{\Sigma}} - \underline{\boldsymbol{\Sigma}}^* + \frac{1}{n} \sum_{t=1}^n p^{-1}(\mathbf{y}_t; \underline{\theta}) \underline{\boldsymbol{\Sigma}} (\mathbf{H}_{1,t} + \mathbf{D}_{1,t}) \underline{\boldsymbol{\Sigma}} \right\|_F. \quad (58) \\ 1258$$

1259 To further analyze the term on the right side, we first give a result of $\mathbf{H}_{1,t} + \mathbf{D}_{1,t}$.

1260 **Lemma 9.** When \mathbf{x} follows a normal mean variance mixture and \mathbf{H}_1 and \mathbf{D}_1 are defined in (34) and (37), respectively, we have

$$1261 \frac{1}{2} (\mathbf{H}_1 + \mathbf{D}_1) = \nabla_{\boldsymbol{\Sigma}} p(\mathbf{y}; \boldsymbol{\theta}). \\ 1262$$

1263 *Proof.* Considering the partial derivative of $p(\mathbf{y}; \boldsymbol{\theta})$ with respect to the (i, j) -th element of $\boldsymbol{\Sigma}$, we have that

$$1264 \nabla_{\boldsymbol{\Sigma}_{ij}} p(\mathbf{y}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\Sigma}_{ij}} \int_{\mathcal{Q}^{-1}(\mathbf{y})} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \nabla_{\boldsymbol{\Sigma}_{ij}} p(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}. \\ 1265$$

1266 The partial derivatives of $p(\mathbf{x}; \boldsymbol{\theta})$ satisfy the following equation:

$$1267 \nabla_{\boldsymbol{\Sigma}_{ij}} p(\mathbf{x} | z; \boldsymbol{\theta}) = \frac{z}{2} \nabla_{x_i x_j}^2 p(\mathbf{x} | z; \boldsymbol{\theta}), \quad (59) \\ 1268$$

1269 Hence, for $i \neq j$, we can obtain that

$$1270 \nabla_{\boldsymbol{\Sigma}_{ij}} p(\mathbf{y}; \boldsymbol{\theta}) = \int_{\mathcal{Q}^{-1}(\mathbf{y}_{\setminus i, j})} \int_{l_i}^{u_i} \int_{l_j}^{u_j} \int_0^{+\infty} \frac{z}{2} \nabla_{x_i x_j}^2 p(\mathbf{x} | z; \boldsymbol{\theta}) p(z) dz dx_i dx_j d\mathbf{x}_{\setminus i, j} \\ 1271 = \int_{\mathcal{Q}^{-1}(\mathbf{y}_{\setminus i, j})} \int_0^{+\infty} \frac{z p(z)}{2} \left(p(\mathbf{x}_{\setminus i, j}, x_i = l_i, x_j = l_j | z; \boldsymbol{\theta}) - p(\mathbf{x}_{\setminus i, j}, x_i = l_i, x_j = u_j | z; \boldsymbol{\theta}) \right. \\ 1272 \quad \left. - p(\mathbf{x}_{\setminus i, j}, x_i = u_i, x_j = l_j | z; \boldsymbol{\theta}) + p(\mathbf{x}_{\setminus i, j}, x_i = u_i, x_j = u_j | z; \boldsymbol{\theta}) \right) dz d\mathbf{x}_{\setminus i, j} \\ 1273 = \frac{1}{2} \mathbb{E}_z(z) \left(p_1(\mathbf{y}_{\setminus i, j}, x_i = l_i, x_j = l_j; \boldsymbol{\theta}) - p_1(\mathbf{y}_{\setminus i, j}, x_i = l_i, x_j = u_j; \boldsymbol{\theta}) \right. \\ 1274 \quad \left. - p_1(\mathbf{y}_{\setminus i, j}, x_i = u_i, x_j = l_j; \boldsymbol{\theta}) + p_1(\mathbf{y}_{\setminus i, j}, x_i = u_i, x_j = u_j; \boldsymbol{\theta}) \right) \\ 1275 = \frac{1}{2} \mathbf{H}_{1, ij}. \\ 1276$$

1277 For the case $i = j$, we first introduce a result.

1278 **Lemma 10.** When \mathbf{x} follows a normal mean-variance mixture, we have the following equation:

$$1279 \sum_{k=1}^d \boldsymbol{\Sigma}_{ik} \nabla_{x_k x_i}^2 p_1(\mathbf{x} | z; \boldsymbol{\theta}) = \nabla_{x_i} \left(-\frac{p_1(\mathbf{x} | z; \boldsymbol{\theta})}{z} (x_i - \mu_i - z \xi_i) \right). \quad (60) \\ 1280$$

1296 *Proof.* We begin with

$$1297 \sum_{k=1}^d \Sigma_{ik} \nabla_{x_k x_i}^2 p_1(\mathbf{x} | z; \boldsymbol{\theta}) = \nabla_{x_i} \left(\sum_{k=1}^d \Sigma_{ik} \nabla_{x_k} p_1(\mathbf{x} | z; \boldsymbol{\theta}) \right)$$

1301 For the term $\sum_{k=1}^d \Sigma_{ik} \nabla_{x_k} p_1(\mathbf{x} | z; \boldsymbol{\theta})$, we have

$$1302 \sum_{k=1}^d \Sigma_{ik} \nabla_{x_k} p_1(\mathbf{x} | z; \boldsymbol{\theta}) = \sum_{k=1}^d \Sigma_{ik} \left(-\frac{p_1(\mathbf{x} | z; \boldsymbol{\theta})}{z} \sum_{l=1}^d \Sigma_{kl}^{-1} (x_l - \mu_l - z \xi_l) \right)$$

$$1303 = -\frac{p_1(\mathbf{x} | z; \boldsymbol{\theta})}{z} \sum_{l=1}^d \left(\sum_{k=1}^d \Sigma_{ik} \Sigma_{kl}^{-1} (x_l - \mu_l - z \xi_l) \right).$$

1308 Since $\sum_{k=1}^d \Sigma_{ik} \Sigma_{kl}^{-1}$ is equivalent to the (i, l) -entry of the matrix $\Sigma \Sigma^{-1}$, we can obtain that

$$1310 \sum_{k=1}^d \Sigma_{ik} \Sigma_{kl}^{-1} = \begin{cases} 0, & i \neq l, \\ 1, & i = l. \end{cases}$$

1313 Hence, we have

$$1314 -\frac{p_1(\mathbf{x} | z; \boldsymbol{\theta})}{z} \sum_{k=1}^d \left(\sum_{l=1}^d \Sigma_{ik} \Sigma_{kl}^{-1} (x_l - \mu_l - z \xi_l) \right) = -\frac{p_1(\mathbf{x} | z; \boldsymbol{\theta})}{z} (x_i - \mu_i - z \xi_i).$$

1317 Therefore, we prove the equation (60). \square

1321 Computing the integral of the left side of (60), we have

$$1322 \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \sum_{k=1}^d \Sigma_{ik} \nabla_{x_k x_i}^2 p_1(\mathbf{x} | z; \boldsymbol{\theta}) p(z) dz d\mathbf{x}$$

$$1323 = \int_0^{+\infty} \int_{\mathcal{Q}^{-1}(\mathbf{y})} \left(\sigma_i \nabla_{x_i}^2 p_1(\mathbf{x} | z; \boldsymbol{\theta}) + \sum_{k \neq i} \Sigma_{ik} \nabla_{x_k x_i}^2 p_1(\mathbf{x} | z; \boldsymbol{\theta}) \right) p(z) dz d\mathbf{x} \quad (61)$$

$$1324 = 2\mathbb{E}_z[z^{-1}] \sigma_i \nabla_{\sigma_i} p(\mathbf{y}; \boldsymbol{\theta}) + [\Sigma \mathbf{H}_1]_{ii}.$$

1330 Then the integral of the right side of (60) is given by

$$1331 \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \nabla_{x_i} \left(-\frac{p_1(\mathbf{x}; \boldsymbol{\theta})}{z} (x_i - \mu_i - z \xi_i) \right) dz d\mathbf{x}$$

$$1332 = \int_{\mathcal{Q}^{-1}(\mathbf{y}_{\setminus i})} \int_{l_i}^{u_i} \int_0^{+\infty} \nabla_{x_i} \left(-\frac{p_1(\mathbf{x}; \boldsymbol{\theta})}{z} (x_i - \mu_i - z \xi_i) \right) dz dx_i d\mathbf{x}_{\setminus i} \quad (62)$$

$$1333 = \mathbb{E}_z(z^{-1}) \left((l_i - \mu_i) p(x_i = l_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) - \xi_i p_1(x_i = l_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) \right.$$

$$1334 \left. - (u_i - \mu_i) p(x_i = u_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) + \xi_i p_1(x_i = u_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) \right).$$

1339 Since (61) and (62) are equivalent, we have

$$1340 \nabla_{\sigma_i} p(\mathbf{y}; \boldsymbol{\theta}) = \frac{1}{2\sigma_i} \left((l_i - \mu_i) p(x_i = l_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) - \xi_i p_1(x_i = l_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) \right.$$

$$1341 \left. - (u_i - \mu_i) p(x_i = u_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) + \xi_i p_1(x_i = u_i, \mathbf{y}_{\setminus i}; \boldsymbol{\theta}) - \mathbb{E}_z[z] [\Sigma \mathbf{H}_1]_{ii} \right) = \frac{1}{2} \mathbf{D}_{1,ii}.$$

1345 Therefore, $\frac{1}{2}(\mathbf{H}_1 + \mathbf{D}_1) = \nabla_{\Sigma} p(\mathbf{y}; \boldsymbol{\theta})$ is valid. \square

1348 Based on Lemma 9, we have that

$$1349 p^{-1}(\mathbf{y}; \boldsymbol{\theta}) \Sigma (\mathbf{H}_1 + \mathbf{D}_1) \Sigma = -2 \nabla_{\Sigma^{-1}} \log p(\mathbf{y}; \boldsymbol{\theta}).$$

Hence, the distance (58) becomes

$$\begin{aligned}
\|\underline{\Sigma} - \Sigma^*\|_F &= \left\| \underline{\Sigma} - \Sigma^* - \frac{2}{n} \sum_{t=1}^n \nabla_{\underline{\Sigma}^{-1}} \log p(\mathbf{y}; \underline{\theta}) + \frac{2}{n} \sum_{t=1}^n \nabla_{\Sigma^{*-1}} \log p(\mathbf{y}; \theta^*) \right\|_F \\
&= \left\| \underline{\Sigma} - \Sigma^* - \frac{2}{n} \sum_{t=1}^n \int_0^1 \nabla_{\tilde{\Sigma}^{-1}, \tilde{\Sigma}}^2 \log p(\mathbf{y}_t; \tilde{\theta}_{\Sigma}) (\tilde{\Sigma} - \Sigma^*) d\beta \right\|_F \\
&\leq \sup_{\beta \in (0,1]} \left\| \mathbf{I}_{d \times d} - \frac{2}{n} \sum_{t=1}^n \frac{\partial \text{vec}(\nabla_{\tilde{\Sigma}^{-1}} \log p(\mathbf{y}_t; \tilde{\theta}_{\Sigma}))}{\partial \text{vec}(\tilde{\Sigma})} \right\|_2 \|\underline{\Sigma} - \Sigma^*\|_F.
\end{aligned} \tag{63}$$

where the second equation is given by the mean value theorem of integrals and $\tilde{\theta}_{\Sigma} = \{\underline{\mu}, \tilde{\Sigma}, \underline{\xi}\}$ with $\tilde{\Sigma} = \beta \underline{\Sigma} + (1 - \beta) \Sigma^*$ and $\beta \in (0, 1]$.

To bound the term $\sup_{\beta \in (0,1]} \left\| \mathbf{I}_{d \times d} - \frac{2}{n} \sum_{t=1}^n \frac{\partial \text{vec}(\nabla_{\tilde{\Sigma}^{-1}} \log p(\mathbf{y}_t; \tilde{\theta}_{\Sigma}))}{\partial \text{vec}(\tilde{\Sigma})} \right\|_2$, we need some results for the term $\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} \log p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma)}$ at first.

Lemma 11. *The term $\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} \log p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma)}$ satisfy*

1. *it is a positive definite matrix;*
2. $\mathbf{I}_{d \times d} - \frac{2}{n} \sum_{t=1}^n \frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} \log p(\mathbf{y}_t; \theta))}{\partial \text{vec}(\Sigma)}$ *is also a positive definite matrix.*

Proof. Since $\log p(\mathbf{x} | z; \theta)$ is linear with respect to Σ^{-1} , the second order derivative of $\log p(\mathbf{y}; \theta)$ with respect to Σ^{-1} is easier to be obtained. We first compute $\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} \log p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma^{-1})}$ and transform it later.

The derivative term $\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} \log p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma^{-1})}$ can be computed as

$$\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} \log p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma^{-1})} = \frac{\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma^{-1})} \cdot p(\mathbf{y}; \theta) - \nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta) \otimes \nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta)}{p^2(\mathbf{y}; \theta)}. \tag{64}$$

Since the first and second order derivatives of $p(\mathbf{y}; \theta)$ with respect to Σ^{-1} satisfy

$$\nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta) = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta) dz d\mathbf{x},$$

and

$$\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma^{-1})} = \int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta))}{\partial \text{vec}(\Sigma^{-1})} dz d\mathbf{x},$$

based on Lemma 5, the expression (64) becomes

$$\begin{aligned}
&\frac{\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta))}{\partial \text{vec}(\Sigma^{-1})} \cdot p(\mathbf{y}; \theta) - \nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta) \otimes \nabla_{\Sigma^{-1}} p(\mathbf{y}; \theta)}{p^2(\mathbf{y}; \theta)} \\
&= \frac{\int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta))}{\partial \text{vec}(\Sigma^{-1})} dz d\mathbf{x}}{p(\mathbf{y}; \theta)} \\
&\quad - \frac{\int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta) dz d\mathbf{x}}{p(\mathbf{y}; \theta)} \otimes \frac{\int_{\mathcal{Q}^{-1}(\mathbf{y})} \int_0^{+\infty} \nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta) dz d\mathbf{x}}{p(\mathbf{y}; \theta)} \\
&= \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \theta} \left[\frac{\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta))}{\partial \text{vec}(\Sigma^{-1})}}{p(\mathbf{x}, z; \theta)} \right] - \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \theta} \left[\frac{\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta)}{p(\mathbf{x}, z; \theta)} \right] \otimes \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \theta} \left[\frac{\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \theta)}{p(\mathbf{x}, z; \theta)} \right].
\end{aligned} \tag{65}$$

1404 Since we have

$$1405 \quad \nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta}) = \frac{p(\mathbf{x}, z; \boldsymbol{\theta})}{2z} (\boldsymbol{\Sigma} - (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top),$$

1408 and

$$1409 \quad \frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta}))}{\partial \text{vec}(\Sigma^{-1})} = \frac{p(\mathbf{x}, z; \boldsymbol{\theta})}{4z^2} (\boldsymbol{\Sigma} - (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top) \otimes (\boldsymbol{\Sigma} - (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top)$$

$$1410 \quad - \frac{p(\mathbf{x}, z; \boldsymbol{\theta})}{2z} \boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma},$$

1413 the expression (65) can further be derived as

$$1414 \quad \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} \left[\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta}))}{\partial \text{vec}(\Sigma^{-1})} \right] - \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} \left[\frac{\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta})}{p(\mathbf{x}, z; \boldsymbol{\theta})} \right] \otimes \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} \left[\frac{\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta})}{p(\mathbf{x}, z; \boldsymbol{\theta})} \right]$$

$$1415 \quad = \frac{1}{4} \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} \left[\frac{1}{z^2} \text{vec}((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top) \text{vec}((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top)^\top - \frac{1}{2z} \boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma} \right]$$

$$1416 \quad - \frac{1}{4} \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \text{vec}((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top)] \mathbb{E}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \text{vec}((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top)^\top]$$

$$1417 \quad = \frac{1}{4} \text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \text{vec}(\mathbf{x}\mathbf{x}^\top)] - \frac{1}{2} \mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1} \boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}].$$

1425 Setting $\tilde{\mathbf{x}} = \mathbf{x}, z | \mathbf{y}$ and $f(\text{vec}(\mathbf{x}\mathbf{x}^\top)) = z^{-1} \text{vec}(\mathbf{x}\mathbf{x}^\top)$, since $\nabla_{\text{vec}(\mathbf{x}\mathbf{x}^\top)}^2 \log p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}) =$
 1426 $-\frac{\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}}{z}$, $p(\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta})$ is a log-concave function with respect to $\text{vec}(\mathbf{x}\mathbf{x}^\top)$. Based on the
 1427 Brascamp–Lieb inequality in Lemma 6, we have

$$1428 \quad \text{Cov}_{\mathbf{x}, z | \mathbf{y}; \boldsymbol{\theta}} (z^{-1} \text{vec}(\mathbf{x}\mathbf{x}^\top)) \prec 2 \mathbb{E}_{z | \mathbf{y}; \boldsymbol{\theta}} [z^{-1}] \boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}.$$

1429 Since $\boldsymbol{\Sigma} \otimes \boldsymbol{\Sigma}$ is positive definite, $\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{y}; \boldsymbol{\theta}))}{\partial \text{vec}(\Sigma^{-1})}$ is negative definite.

1432 Based on

$$1433 \quad \frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta}))}{\partial \text{vec}(\boldsymbol{\Sigma})} = - \frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta}))}{\partial \text{vec}(\Sigma^{-1})} (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}),$$

1436 we can obtain that $\frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{x}, z; \boldsymbol{\theta}))}{\partial \text{vec}(\boldsymbol{\Sigma})}$ is positive definite and

$$1437 \quad \mathbf{I}_{d \times d} - \frac{2}{n} \sum_{t=1}^n \frac{\partial \text{vec}(\nabla_{\Sigma^{-1}} p(\mathbf{y}_t; \boldsymbol{\theta}))}{\partial \text{vec}(\boldsymbol{\Sigma})}$$

$$1438 \quad = \frac{1}{2n} \sum_{t=1}^n (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \boldsymbol{\theta}} (\text{vec}(\mathbf{x}\mathbf{x}^\top)) \succ \mathbf{0}.$$

1444 \square

1447 Based on Lemma 11, we have

$$1448 \quad \sup_{\beta \in (0, 1]} \left\| \mathbf{I}_{d \times d} - \sum_{t=1}^n \frac{\partial \text{vec}(\nabla_{\tilde{\Sigma}^{-1}} \log p(\mathbf{y}_t; \tilde{\boldsymbol{\theta}}_\Sigma))}{\partial \text{vec}(\tilde{\boldsymbol{\Sigma}})} \right\|_2$$

$$1449 \quad = \sup_{\beta \in (0, 1]} \left\| \frac{1}{2n} \sum_{t=1}^n (\tilde{\boldsymbol{\Sigma}}^{-1} \otimes \tilde{\boldsymbol{\Sigma}}^{-1}) \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \tilde{\boldsymbol{\theta}}_\Sigma} [\text{vec}(\mathbf{x}\mathbf{x}^\top)] \right\|_2$$

$$1450 \quad = \max_{\boldsymbol{\theta}} \left\| \frac{1}{2n} \sum_{t=1}^n (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{Cov}_{\mathbf{x}, z | \mathbf{y}_t; \boldsymbol{\theta}} (\text{vec}(\mathbf{x}\mathbf{x}^\top)) \right\|_2$$

$$1451 \quad \triangleq c_\Sigma \in (0, 1).$$

C PROOF OF THEOREM 2

Based on Proposition 1, we have proved the linear convergence rates of the three parameters in θ . To connect the global convergence rates and the separable convergence rates, we introduce a result.

Lemma 12 (Meng & Rubin (1994)). *The global convergence rate of an ECM algorithm is the maximum of the componentwise rates of convergence.*

Based on Lemma 12, we obtain that

$$\|\theta^{(k)} - \theta^*\|_2 \leq c^k \|\theta^{(0)} - \theta^*\|_2,$$

where $c = \max\{c_\mu, c_\xi, c_\Sigma\} \in (0, 1)$.

D DERIVATIONS OF SURROGATE FUNCTIONS IN QUANTIZED MATRIX COMPLETION AND COMPRESSIVE SENSING

D.1 SURROGATE FUNCTION DERIVATION IN QUANTIZED MATRIX COMPLETION

The goal of the low-rank matrix completion problem is to recover an unknown low-rank matrix $M \in \mathbb{R}^{d_1 \times d_2}$ from an observed, yet incomplete, matrix. Let X denote a matrix whose entries are drawn from a normal mean-variance mixture distribution. We use μ_{ij} and x_{ij} to represent the (i, j) -th entries of M and X , respectively. Based on the normal mean-variance model, the relationship between μ_{ij} and x_{ij} is given by

$$x_{ij} = \mu_{ij} + z\xi + z^{1/2}\sigma\epsilon,$$

where μ_{ij} serves as the location parameter of x_{ij} and both ξ and σ are given constants. In the quantization scenario, we have $Y = Q(X)$. The entire matrix Y is not available for observation. Let Ω denote the index set of the observed entries. Our goal is to recover M from incomplete Y , which is equivalent to estimating all the μ_{ij} in the matrix M .

Since x_{ij} follows a univariate normal mean-variance model and $Y = Q(X)$, the density function of y_{ij} is identical with (5) under the univariate case. Hence, the negative log-likelihood function for y_{ij} with $ij \in \Omega$ is

$$\sum_{(i,j) \in \Omega} \log p(y_{ij} | M),$$

which is the objective function of (16). Applying Jensen's inequality as in (9), we have the surrogate function

$$\begin{aligned} S(M; \underline{M}) = & \sum_{(i,j) \in \Omega} \left[\mathbb{E}_{z_t | y_t; \underline{\theta}} [\log p(z_t)] - \frac{1}{2} \log \det \sigma^2 \right. \\ & \left. - \frac{1}{2\sigma^2} ((u_{ij} - 2v_{ij}\mu + \iota_{ij}\mu^2) - 2(w_{ij} - \mu)\xi + \zeta_{ij}\xi^2) \right] + \text{const.}, \end{aligned}$$

where u_{ij} , v_{ij} , and w_{ij} are the univariate versions of U_{ij} , v_{ij} , and w_{ij} , respectively. Ignoring the terms which are independent with μ_{ij} , the surrogate function becomes

$$S(M; \underline{M}) = \frac{1}{2\sigma^2} \sum_{(i,j) \in \Omega} (2v_t\mu - \iota_t\mu^2 - 2\mu\xi) + \text{const.}$$

Making a low-rank factorization to the matrix M as $M = AB^\top$, we have $\mu_{ij} = \mathbf{a}_i \mathbf{b}_j^\top$. Setting $e_{ij} = \frac{v_{ij} - \xi}{\iota_{ij}}$, we have that maximizing $S(M; \underline{M})$ is equivalent to maximize

$$\sum_{(i,j) \in \Omega} (\mathbf{a}_i \mathbf{b}_j^\top - e_{ij})^2,$$

which is identical with (17).

1512 D.2 SURROGATE FUNCTION DERIVATION IN QUANTIZED COMPRESSIVE SENSING

1513
1514 In quantized compressive sensing, the base model with normal mean-variance mixture noise is given
1515 by

$$1516 \mathbf{y} = \mathcal{Q}(\mathbf{x}), \quad \mathbf{x} = \Phi\boldsymbol{\vartheta} + z\xi + z^{1/2}\sigma\epsilon.$$

1517
1518 In the above model, the term $\Phi\boldsymbol{\vartheta}$ can be regarded as the location parameter of \mathbf{x} . We can use the
1519 ML estimation method to recover the sparse signal $\boldsymbol{\vartheta}$, which has the optimization problem

$$1520 \max_{\boldsymbol{\vartheta}} \log p(\mathbf{y} | \boldsymbol{\vartheta}).$$

1521
1522 Following the suggestions from Zymnis et al. (2009), we add a ℓ_1 -regularization term to force the
1523 solution of $\boldsymbol{\vartheta}$ to be sparse. Hence, we obtain the optimization problem

$$1524 \max_{\boldsymbol{\vartheta}} \log p(\mathbf{y} | \boldsymbol{\vartheta}) + \eta\|\boldsymbol{\vartheta}\|_1.$$

1525
1526 To solve the above optimization problem through ECM algorithm, we can apply the E-step as in (9)
1527 to obtain the surrogate function

$$1528 \begin{aligned} 1529 S(\boldsymbol{\vartheta}; \underline{\boldsymbol{\vartheta}}) &= \mathbb{E}_{z_t | \mathbf{y}_t; \underline{\boldsymbol{\vartheta}}} [\log p(z_t)] - \frac{1}{2} \log \det \sigma^2 \\ 1530 &\quad - \frac{1}{2\sigma^2} ((\mathbf{u} - 2\mathbf{v}^\top \mathbf{A}\boldsymbol{\vartheta} + \iota\boldsymbol{\vartheta}^\top \mathbf{A}^\top \mathbf{A}\boldsymbol{\vartheta}) - 2(\mathbf{w} - \mathbf{A}\boldsymbol{\vartheta})^\top \xi + \zeta\xi^\top \xi) + \eta\|\boldsymbol{\vartheta}\|_1 + \text{const}. \end{aligned}$$

1531 Ignoring the terms which are independent with $\boldsymbol{\vartheta}$, the surrogate function becomes

$$1532 S(\boldsymbol{\vartheta}; \underline{\boldsymbol{\vartheta}}) = \frac{1}{2\sigma^2} (2\mathbf{v}^\top \mathbf{A}\boldsymbol{\vartheta} - \iota\boldsymbol{\vartheta}^\top \mathbf{A}^\top \mathbf{A}\boldsymbol{\vartheta} - 2\boldsymbol{\vartheta}^\top \mathbf{A}^\top \xi) + \eta\|\boldsymbol{\vartheta}\|_1 + \text{const}.$$

1533 Setting $\mathbf{e} = \frac{\mathbf{v} - \xi}{\iota}$, we have that maximizing $S(\boldsymbol{\vartheta}; \underline{\boldsymbol{\vartheta}})$ is equivalent to maximize

$$1534 \|\mathbf{A}\boldsymbol{\vartheta} - \mathbf{e}\|_2^2 + \eta\|\boldsymbol{\vartheta}\|_1,$$

1535 which is identical with the surrogate function in quantized compressive sensing.

1544 E EXPERIMENT DETAILS

1546 E.1 BENCHMARK SETTINGS

1547
1548 **Quantized matrix completion:** For all algorithms presented in Table 2, the complete matrix M is
1549 reconstructed from the training data. Denote by Ω_{test} the set of indices corresponding to entries in
1550 the test data. The definitions of the two benchmarks are given as follows.

- 1551 1. Accuracy: $\sum_{(i,j) \in \Omega_{\text{test}}} I(\mathcal{Q}(\mu_{ij}) = y_{ij});$
- 1552 2. RMSE: $\sqrt{\frac{1}{n_{\text{test}}} \sum_{(i,j) \in \Omega_{\text{test}}} (\mu_{ij} - y_{ij})^2},$

1553 where $I(\cdot)$ is the indicator function.

1554 **Quantized compressive sensing:** The signal-to-noise ratio (SNR) is defined as the ratio of the
1555 expectation of the original measurements to the variance of the noise term, which is

$$1556 \text{SNR} = \frac{\mathbb{E}[\mathbf{A}\boldsymbol{\vartheta}]}{\text{Var}[\mathbf{x} - \mathbf{A}\boldsymbol{\vartheta}]}.$$

1557 Denote the ground truth sparse signal as $\boldsymbol{\vartheta}_{\text{gd}}$, the estimated sparse signal as $\boldsymbol{\vartheta}_{\text{est}}$. The cosine simi-
1558 larity is defined as

$$1559 \text{Cos Sim} = \frac{\boldsymbol{\vartheta}_{\text{gd}}^\top \boldsymbol{\vartheta}_{\text{est}}}{\|\boldsymbol{\vartheta}_{\text{gd}}\|_2 \|\boldsymbol{\vartheta}_{\text{est}}\|_2}.$$

1566 E.2 EXPERIMENTAL SETTINGS
1567

1568 In Figure 4, we set $\mu_i \sim \text{Uniform}(0, 1)$, $\xi_i \sim \text{Uniform}(0, 2)$, and $\mathbf{v}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ for $i = 1, \dots, d$,
1569 and construct the covariance matrix as $\Sigma = \sum_{i=1}^d \mathbf{v}_i \mathbf{v}_i^\top$. For the GH distribution, we use $\lambda = -0.5$,
1570 $\delta = 2$, and $\gamma = 2$; for GHST, we set $\nu = 5$, $\lambda = -\mu/2$, $\delta = \nu$, and $\gamma = 0$. The number of samples
1571 is $n = 1 \times 10^4$ and the threshold is fixed at $\tau = 0.3$.

1572 The specific parameter settings for the algorithms in Table 2 are given by
1573

- 1574 1. Standard Gaussian: $r = 2$ and $\sigma^2 = 0.8$;
 - 1575 2. one-bit Gaussian: $r = 2$ and $\sigma^2 = 0.8$;
 - 1576 3. multi-bit Gaussian: $r = 2$ and $\sigma^2 = 0.8$;
 - 1577 4. multi-bit Student's t : $r = 2$, $\sigma^2 = 0.8$, and $\nu = 6$;
 - 1578 5. multi-bit GHST: $r = 2$, $\sigma^2 = 0.8$, $\xi = -0.1$, and $\nu = 8$;
 - 1579 6. multi-bit GH: $r = 2$, $\sigma^2 = 0.8$, $\xi = -0.1$, $\lambda = -4$, $\delta = 4$, and $\gamma = 0.3$.
- 1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619