Computational Approaches for Exponential-Family Factor Analysis

Liang Wang leonwang@bu.edu

Luis Carvalho lecarval@bu.edu

Department of Mathematics and Statistics **Boston University**

Abstract We study a general factor analysis framework where the *n*-by-*p* data matrix is assumed to follow a general exponential family distribution entry-wise. While this model framework has been proposed before, we here further relax its distributional assumption by using a quasi-likelihood setup. By parameterizing the mean-variance relationship on data entries, we additionally introduce a dispersion parameter and entry-wise weights to model large variations and missing values. The resulting model is thus not only robust to distribution misspecification but also more flexible and able to capture non-Gaussian covariance structures of the data matrix. Our main focus is on efficient computational approaches to perform the factor analysis. Previous modeling frameworks rely on simulated maximum likelihood (SML) to find the factorization solution, but this method was shown to lead to asymptotic bias when the simulated sample size grows slower than the square root of the sample size *n*, eliminating its practical application for data matrices with large *n*. Borrowing from expectation-maximization (EM) and stochastic *n.* Borrowing from expectation-maximization (EM) and stochastic gradient descent (SGD), we investigate three estimation procedures based on iterative factorization updates. Our proposed solution does not show asymptotic biases, and scales even better for large matrix factorizations with error O(1/p). To support our findings, we conduct simulation experiments and discuss its application in three case studies. *Keywords:* matrix factorization, exponential family, factor model. **1. Introduction** Over the past decades, factor analysis has gained tremendous attention in psychology (Ford et al., 1986), computer science (Prince \mathcal{V} *n*. Borrowing from expectation-maximization (EM) and stochastic

tention in psychology (Ford et al., 1986), computer science (Prince et al., 2008), finance (Fama and French, 2015), and biological research (Xu et al., 2021). In particular, when the data $X \in \mathbb{R}^{n \times p}$ is high dimensional $(n \ll p)$, effectively modeling and estimating the covariance structure has been problematic (Basilevsky, 1994). The factor model provides an effective approach to model high dimensional data in which the covariance of the observations is assumed to lie on a lower dimensional manifold.

Despite its popularity in modeling high dimensional data, factor models have several limitations. First and foremost, both data and latent variables are assumed to follow a Gaussian distribution, which is not ideal for modeling binary, count, or other non-constant

variance data. To address the first limitation, there exist some prior works (Wedel and Kamakura, 2001; Wedel et al., 2003) that extend the factor model with more general exponential family assumption. However, even with the improved assumption from Gaussian, the exponential family distribution assumption is often too restrictive for real world, overly dispersed data. Moreover, as we shown later in Section 1.1.3, the proposed maximum likelihood estimation algorithm for such an extended model is problematic with both numerical and asymptotic convergence issues. As a minor issue, the latent factors are only identifiable up to a rotational transformation, potentially causing problems in interpreting the latent factors. Lastly, both these extended works and the traditionally factor analysis framework lack the flexibility to model missing data, preventing several interesting applications such as matrix completion.

This paper thus aims at generalizing the existing works by:

- · Assuming only a mean-variance relationship along with column-wise dispersion parameters to model data covariance;
- Providing interpretability for latent factors via orthogonal identifiability constraints;
- Proposing fast, accurate, and robust optimization algorithms leveraging modern advances in stochastic optimization;
- Facilitating application with an efficient package implementation that allows for entry-wise factor weights and covariance modeling.

To introduce appropriate notations and to understand some of the relevant attempts to address those issues, we elaborate below on the limitations of factor models along with some existing remedies proposed in the literature that motivated our generalization.

1.1. The Factor Model and Its Limitations

For given data $X \in \mathbb{R}^{n \times p}$, the traditional rank q factor model assumes that $q \ll p$ and that the data is generated by latent factors $\Lambda \in \mathbb{R}^{n \times q}$ with $\Lambda^{\dagger} = [\Lambda_1 \cdots \Lambda_n]$, and a deterministic projection matrix $V \in \mathbb{R}^{p \times q}$, the loading matrix. We implicitly assume the following data generating process: for each observation i = 1, ..., n:

:: 4

$$\begin{aligned}
& \Lambda_i \stackrel{\text{norm}}{\sim} N(0, I_q), \\
& X_i \mid \Lambda_i \stackrel{\text{ind}}{\sim} N(V\Lambda_i, \Phi)
\end{aligned} \tag{1}$$

where Φ is a *p*-th order symmetric positive definite matrix, a covariance matrix providing potential heterogeneous noise.

Marginally, $X_i \stackrel{\text{ind}}{\sim} N(0, \Phi + VV^{\top})$, so maximum likelihood estimation of V and Φ is equivalent to covariance estimation. It is common to assume that Φ is diagonal so as to not confound the effect of the loadings V (Bartholomew et al., 2011), and so we adopt the same assumption from now on. In this case, the MLE estimator for V can be obtained in closed form using matrix calculus, based on the eigen-decomposition of the data covariance. Alternatively, V and Φ can be estimated via expectation-maximization, especially if some of the entries in the data matrix X are missing.

The model in general has interesting connections to matrix factorization. For example, probability PCA (Tipping and Bishop, 1998) can be considered as equivalent to the factor model with the only difference that the factor model permits heterogeneous noise structure through the specification of Φ . Under $\Phi = \sigma^2 I_p$, Anderson (1963) established the connection between these two models by demonstrating that the stationary point solution of the factor model likelihood spans the columns of the sample covariance eigenvectors. Drawing further the analogy from the relationship between probability PCA and PCA, the factor model can be considered as the random counterpart of matrix factorization by allowing the factorized components (or latent local factors) Λ to be random.

While finding a wide range of applications, the factor model and its deterministic counterpart (matrix factorization) have, however, some limitations. We discuss them below, including a brief summary on some recent improvements, along with our proposed solutions to further generalize the factor model.

1.1.1. Relaxing the restrictive distributional assumption

In a factor model setup, both the latent variable Λ_i and the data are assumed to (conditionally) follow a Gaussian distribution, yielding a marginal Gaussian distribution for the data. Both assumptions require careful examination when dealing with real data.

For the data distribution, assuming simply a Gaussian data likelihood overlooks many interesting structures in the data. For example, network adjacency matrices take only binary values of 0 and 1, while computer images take a integer values for pixel intensities. Both types of data have been shown to be better modeled with discrete distributions from the exponential family (Wang and Carvalho, 2023). One obvious relaxation is thus to extend the data likelihood assumption from Gaussian to exponential families, or, to accommodate more robust specifications, to specify mean and covariance structures, as in quasi-likelihood approaches. Moreover, those exponential family generalizations do not consider the flexible covariance modeling of the data matrix, which has shown to be one of the most important applications of Gaussian factor model (Fan et al., 2008). Ideally, at least a column-wise idiosyncratic error structure should be modeled for a flexible consideration of the high-dimensional data covariance.

The latent variable assumption is usually considered less restrictive when compared to the likelihood assumption, as evidenced from similar Gaussian latent structures in hierarchical statistical models (e.g. the random effects model (Borenstein et al., 2010) and the state space model (Carter and Kohn, 1996)). This assumption is however frequently studied together with factor identifiability (Shapiro, 1985) to ensure unique latent representations of the data. Specifically, the factor model (1) is not identifiable (or unique) since for any orthogonal matrix $T \in O(q)$, $\Lambda_i^* \doteq T\Lambda_i \stackrel{\text{iid}}{\sim} N(0, I_q)$ and, with $V^* = VT^{\top}$, $X_i | \Lambda_i^* \stackrel{\text{ind}}{\sim} N(V^*\Lambda_i^*, \Phi)$ specify the same model since $V\Lambda_i = VT^{\top}T\Lambda_i = V^*\Lambda_i^*$, that is, $\Lambda^* = \Lambda T^{\top}$ and V^* are not identifiable from Λ and V. For this reason it is common in factor analysis to rotate factors after fitting the model to achieve better sparsity and/or interpretability, e.g. with varimax rotation (Kaiser, 1958). However, it is advantageous to address these identifiability issues from the outset to reduce the space of potential solutions and speed up estimation procedures. In this case, it is helpful to borrow from the matrix factorization research. For example, adding various factorization constraints such as sparsity (Gribonval and Schnass, 2010), positivity (Lee and Choi, 1999) and orthogonality (Li et al., 2010) was shown to provide more representative and unique latent factors. The stochastic counterpart of these factor constraints is closely related to an evolving research field related to data manifolds (Ma and Fu, 2012).

1.1.2. Allowing entry-wise weight and link transformation

Another potential improvement that has remained absent from factor analysis research is the specification of entry-wise likelihood weights and non-linear transformations. In matrix factorization, allowing entry-wise factorization weights and the flexibility of transforming the original data has been shown to be valuable in providing more representative factorized results. For example, in the field of natural language processing, Global Vectors for Word Representation (GloVe; Jeffrey Pennington and Manning, 2014) received great success in obtaining word embeddings. The method essentially applied a log transformation on the word-occurrence matrix with heuristic entry-wise weights to avoid over and underweighting toward rare and common word co-occurrences. In the field of computer vision (Kalayeh et al., 2014), weighting matrices related to classification class frequencies are introduced to alleviate issues with class imbalance. This residual boosting weight matrix provides latent factors that are more suitable for downstream classification. In the field of matrix completion (Davenport et al., 2014), specification of zero factorization weights can eliminate missing entries from the factorization, which in turn allows the latent structure to impute them.

Perhaps due to the computational complexity associated with these enhancements, such flexibility has not been transferred from matrix factorization to factor analysis. Link transformations might appear in the literature, e.g. (Reimann et al., 2002), but are mostly applied as an ad-hoc pre-processing methodology. In practice, it is clear that entry-wise factor weights and link transformations could greatly improve the flexibility of the factor modeling framework. Nevertheless, a unified factor modeling framework that enables such flexibility is still missing from the literature.

1.1.3. Improving on efficient optimization

Lastly, as we seek to improve on the traditional Gaussian factor model, it is natural to consider practical computational concerns: can we scale fitting the improved model to modern large datasets?

While the marginal likelihood under model (1) is available in closed form, deriving the marginal likelihood under a non-Gaussian data assumption is typically difficult and recent research have resorted to simulated maximum likelihood (SML; Wedel and Kamakura, 2001), Markov chain Monte Carlo (MCMC) or variational inference (Gopalan et al., 2015). However, these methods have their own difficulties. Variational inference is based on an approximation to the target marginal distribution and usually relies on oversimplified representations for computational gains at the cost of poor representativity. As for MCMC, due to the identifiability issue introduced earlier, the marginal likelihood is constant along high dimensional quotient spaces on V imposed by equivalence under orthogonal operations (rotations). These equivalent spaces cause significant challenges for both the MCMC sampling and the assessment of convergence. Lastly, although theoretically attractive, MCMC is computationally expensive since it usually requires long running times to achieve convergence up to a desired precision when compared to other approaches such as Laplacian approximations (Rue et al., 2009).

As the original optimization method proposed with the initial exponential factor generalization (Wedel and Kamakura, 2001; Wu and Zhang, 2003), the SML approach is considered as one of the most common estimation methods. Specifically, the maximum like-lihood estimator is obtained by maximizing the following simulated likelihood based on *S* Monte Carlo samples,

$$L(V;X) = \prod_{i=1}^{n} f_{V}(X_{i}) \approx \prod_{i=1}^{n} \frac{1}{S} \sum_{s=1}^{S} f_{V}(X_{i} \mid \Lambda_{i}^{(s)}) \doteq L_{\mathrm{MC}}(V;X),$$

where $\Lambda_i^{(s)} \stackrel{\text{iid}}{\sim} N(0, I_q)$. To obtain the maximizer of log $L_{\text{MC}}(V; X)$, the gradient is needed (up to a constant):

$$\nabla_V \sum_{i=1}^n \log \sum_{s=1}^S f_V(X_i \mid \Lambda_i^{(s)}) = \sum_{i=1}^n \frac{\sum_{s=1}^S \nabla_V f_V(X_i \mid \Lambda_i^{(s)})}{\sum_{s=1}^S f_V(X_i \mid \Lambda_i^{(s)})}.$$
 (2)

Despite the fact that $\nabla_V f_V(X_i | \Lambda_i^{(s)})$ is readily known in closed form, optimization using (2) has both numerical and theoretical issues. For the numerical issue, we need to observe that the likelihood $f_V(X_i | \Lambda_i^{(s)})$ evaluations in the denominator need to be performed in log space to avoid underflows and usually require good starting points for *V*, which are particularly challenging when the data dimension *p* is large.

From a theoretical perspective, the likelihood along its gradient evaluation depends heavily on the asymptotic behavior of sample size S. It has been shown in (Lee, 1995) that the estimator will be asymptotically biased if the MC sample size S does not grow

faster than data sample size \sqrt{n} . In fact, we verified with numerical studies that the gradient estimation can potentially require a larger MC sample size $S \gg \sqrt{n}$ to stabilize $\nabla_V f_\theta(X_i | \Lambda_i^{(s)})$ in (2). Consequently, when optimizing the likelihood via SML, there is a trade-off between computation efficiency and estimation bias. For modern applications of large data dimensions, the sample size *S* required to control the MC variance can be quite large, thus preventing the practical applications of such methods.

1.1.4. Real world applications

Although PCA has been traditionally used for many real world applications (Li et al., 2024), in the past decades, the generalization of the deterministic PCA factorization to the exponential family has enabled a series of benchmark models across different fields. For example, the non-negative matrix factorization (NMF; Lee and Choi, 1999) in computer vision generalized the data distributional assumption to Poisson. The non-Gaussian state space model (Kitagawa, 1987) in time series generalized the data distributional assumption to non-Gaussian using non-parametric estimation; the Skip-gram model (Levy and Goldberg, 2014; Mikolov et al., 2013) in natural language processing generalized the data assumption to multinomial. Perhaps most relevant to statistics factor model research, (Wedel and Kamakura, 2001) and (Wu and Zhang, 2003) generalized the data likelihood to exponential family distribution while allowing the random specification of a latent factor Λ .

Perhaps due to the infeasibility of the SML estimation described in the previous subsection, the lack of a practical estimation method has limited the application of the random factorization to only Bernoulli factor models with an identity link function, i.e, the random dot product model (RDPM; Hoff et al., 2002; Young and Scheinerman, 2007). Despite its restrictive identity link assumption, the RDPM has established its popularity on its empirical evidence from network analysis. After addressing the SML estimation problem, we also feel that empirical evidence of such a generalized model has still been missing from the literature.

1.2. Organization of the paper

The paper is organized as follows: in Section 2 we introduce a more general *exponential* factor model that addresses the shortcomings listed in 1.1.1 and 1.1.2; next, in Section 3, we discuss our main contributions—a collection of efficient and robust optimization strategies for inference, tackling the points in 1.1.3; Section 4 demonstrates the effectiveness of our factorization with simulated examples and applications on benchmark data from various fields; finally, Section 5 concludes with a summary of the innovations and directions for future work.

2. Exponential Factor Models

2.1. Guaranteeing factor identifiability

We start by addressing the model issues raised in Section 1. Given our concern with computational efficiency, our first issue is

non-identifiability; as discussed in 1.1.1, we need to constraint the factors to avoid lack of identifiability due to rotations. From now on we adopt the following standardization of the factors:

- (i) $\Lambda_i \stackrel{\text{iid}}{\sim} N(0, I_q)$ for $i \in [n]$ as usual, with the distribution of the rows of Λ being invariant to orthogonal transformations;
- (ii) V has scaled pairwise orthogonal columns, that is, V = UD with U ∈ S_{p,q}(ℝ), a p-frame in the Stiefel manifold of order q, and D = Diag_{j∈[q]}{d_j} with d₁ ≥ ··· ≥ d_q > 0, so that V^TV = D². We denote this space for V as S̃_{p,q}(ℝ).

This setup makes the factorization model identifiable since for any arbitrary $T \in O(q)$, $V^* = VT^{\top}$ can only belong to $\widetilde{S}_{p,q}(\mathbb{R})$ if T^{\top} commutes with a diagonal matrix, that is, if $T \in O(1)^q$, and so V is unique (up to column sign changes, as in the SVD). In practice, given any pair of factors $\widehat{\Lambda}$ and \widehat{V} we just need to find the singular value decomposition of $\widehat{\Lambda}\widehat{V}^{\top} = \Lambda DU^{\top}$ to identify Λ and V = UD.

2.2. Generalizing the normal likelihood

Next, we relax the convenient but often unrealistic Gaussian assumptions in the likelihood and settle with a more general mean and variance specification in the spirit of *quasi-likelihood*. We assume that $X_i | \Lambda_i \sim F(V\Lambda_i, \Phi_i)$ where F belongs to the exponential family with link function g and variance function \mathcal{V} , that is,

$$\mathbb{E}(X_i|\Lambda_i) \doteq \mu_i = g^{-1}(\eta_i), \text{ with } \eta_i = V\Lambda_i + \eta_0, \text{ and}$$

$$\operatorname{Var}(X_i|\Lambda_i) = \Phi_i \mathbb{V}(\mu_i),$$
(3)

where $\mathbb{V}(\mu_i) = \text{Diag}\{\mathcal{V}(\mu_i)\}$ is the diagonal variance and $\eta_0 \in \mathbb{R}^p$ is the latent center of the factor model. This way, we can more naturally represent data X belonging to fields other than real numbers; common cases are binary data with F being Bernoulli or binomial (with weights) and count data with F being Poisson or negative binomial. In particular, the negative binomial distribution offers enhancements over the Poisson distribution by effectively accommodating the over-dispersion characteristic often observed in count data; see, e.g., (Xia, 2020) for a detailed treatment of negative binomial distributions as a compound Poisson type. To accommodate entry-wise weights, as motivated in Section 1.1.2, we set $\Phi_i = \Phi W_i^{-1}$ where $W_i = \text{Diag}_{j=1,...,p}\{w_{ij}\}$ are the known weights, that is, $\Phi_i = \text{Diag}_{j=1,...,p}\{\phi_j/w_{ij}\}$. This setup implies $\mathbb{E}(X_{ij}|\Lambda_i) = \mu_{ij}$ and $\text{Var}(X_{ij}|\Lambda_i) = \phi_j \mathcal{V}(\mu_{ij})/w_{ij}$. From these two moment conditions, we can adopt the extended quasilikelihood (Nelder and Pregibon, 1987) to define:

$$\log f_{V,\eta_0,\Phi}(X_i|\Lambda_i) = -\sum_{j=1}^{p} \frac{w_{ij}}{\phi_j} \int_{\mu_{ij}}^{X_{ij}} \frac{X_{ij} - t}{\mathcal{V}(t)} dt - \frac{1}{2} \log\left(2\pi \frac{\phi_j \mathcal{V}(X_{ij})}{w_{ij}}\right).$$
(4)

We call this the exponential factor model (EFM).

The MLE estimate of $\theta = (V, \eta_0, \Phi)$ then requires access to the marginal density for each observation $i \in [n]$,

$$\log f_{\theta}(X_i) = \int_{\Lambda_i} \log f_{\theta}(X_i | \Lambda_i) f(\Lambda_i) d\Lambda_i$$
(5)

where f is the density of the standard multivariate normal density of order q:

$$\widehat{\theta} = \operatorname*{argmax}_{V \in \widetilde{\mathcal{S}}_{p,q}(\mathbb{R}), \eta_0 \in \mathbb{R}^p, \phi_1, \dots, \phi_p > 0} \sum_{i=1}^n \log f_{V, \Phi}(X_i).$$
(6)

2.3. Modeling covariance

Such a generalized factor framework can be used to efficiently estimate the covariance structure of high dimensional data. Specifically, applying the total variance formula we can derive the covariance of a new observation $X \mid \lambda \sim F(V\lambda, \Phi)$ given $\lambda \sim N(0, I_q)$,

$$\operatorname{Cov}_{\theta}(X) = \mathbb{E}_{\lambda} \left(\operatorname{Var}_{\theta}(X|\lambda) \right) + \operatorname{Var}_{\lambda} \left(\mathbb{E}_{\theta}(X|\lambda) \right).$$
(7)

Here, the first term is a diagonal matrix but the second term requires an outer product that induces correlations among the entries of X,

$$\mathbb{E}_{\lambda} \big(\operatorname{Var}_{\theta}(X|\lambda) \big) = \operatorname{Diag}_{j \in [p]} \Big\{ \phi_{j} \mathbb{E}_{\lambda} \big[\mathcal{V} \circ g^{-1} \big((V\lambda)_{j} \big) \big] \Big\},$$

$$\operatorname{Var}_{\lambda} \big(\mathbb{E}_{\theta}(X|\lambda) \big) = \mathbb{E}_{\lambda} \Big[\big(g^{-1}(V\lambda) - \mu_{\lambda} \big) \big(g^{-1}(V\lambda) - \mu_{\lambda} \big)^{\mathsf{T}} \Big],$$
(8)

where $\mu_{\lambda} = \mathbb{E}_{\lambda}(g^{-1}(V\lambda))$.

In the special case of Gaussian distribution, we have $g(\mu) = \mu$ and $\mathcal{V}(\mu) = 1$ and so, as expected, $\operatorname{Cov}_{\theta}(X) = \Phi + VV^{\top}$. Using this result, Fan et al. (2008) demonstrated the efficiency of the plug-in covariance estimator via estimation of $V, \phi, \operatorname{Cov}(\Lambda)$. Similarly, we demonstrate that the plug-in efficiency of (V, Φ) in high dimensional covariance estimation is preserved under our generalized quasi-factor setup using Eq (8).

3. Approximate but Efficient and Robust Optimization

As we mentioned in Section 1.1.3, directly optimizing Eq (6) through simulated gradients has both numerical and theoretical issues. Next, we demonstrate how we can conduct maximum like-lihood estimation on the factor matrix V and conditional variances Φ through some approximate but efficient and robust algorithms.

3.1. EM optimization with small q

Similar to the Gaussian factor model, one common estimation method for such a latent space model should be Expectation Maximization (EM). In our setup, the EM can be formulated as follow:

```
• E-Step:
```

Given parameter at iteration step t: $\theta^{(t)}$, we compute:

$$\begin{split} \mathbb{E}_{\Lambda_{i}|X_{i}}(\log f_{\theta^{(t)}}(X_{i},\Lambda_{i})) &= \int_{\Lambda_{i}}\log f_{\theta^{(t)}}(X_{i},\Lambda_{i})f_{\theta^{(t)}}(\Lambda_{i}|X_{i})d\Lambda_{i}\\ &\approx \int_{\Lambda_{i}}\log f_{\theta^{(t)}}(X_{i},\Lambda_{i})\widetilde{f}_{\theta^{(t)}}(\Lambda_{i}|X_{i})d\Lambda_{i}\\ &= \mathbb{E}_{\widetilde{\Lambda}_{i}|X_{i}}(\log f_{\theta^{(t)}}(X_{i},\Lambda_{i})) \end{split}$$

where $\tilde{f}_{\theta^{(t)}}(\Lambda_i | X_i)$ is the Gaussian approximated posterior density:

$$f_{\theta^{(t)}}(\Lambda_i | X_i) = f_N(\Lambda_i, H_{\widehat{\Lambda}_i})$$

with $\widehat{\Lambda}_i$ and $H_{\widehat{\Lambda}_i}$ easily obtained from penalized GLM problem:

$$\widehat{\Lambda}_{i} = \underset{\Lambda_{i}}{\operatorname{argmax}} \log f_{\theta^{(t)}}(X_{i}|\Lambda_{i}) + \log f(\Lambda_{i})$$

$$H_{\widehat{\Lambda}_{i}} = \nabla_{\Lambda_{i}}^{(2)} \Big(\log f_{\theta^{(t)}}(X_{i}|\Lambda_{i}) + \log f(\Lambda_{i})\Big) (\widehat{\Lambda}_{i})$$
(10)

Generally denote response $y \in \mathbb{R}^n$, covariate $x \in \mathbb{R}^{n \times p}$, coefficient $\beta \in \mathbb{R}^p$, prior $\mu_0 \in \mathbb{R}^p$, $\Sigma_0 \in \mathbb{R}^{p \times p}$, dispersion, $\phi \in \mathbb{R}$, and weight $w \in \mathbb{R}^n$, the solution to the Bayesian/penalized GLM of Eq (10) can be formulated as:

$$\underset{\beta^{(t+1)}}{\operatorname{argmin}} \sum_{i=1}^{n} w_i \log f_{\beta^{(t)}}(y_i | x_i) + \log f_N(\mu_0, \Sigma_0)$$
(11)

The solution can be obtained by solving $\beta^{(t+1)}$ from the following equation:

$$(x^{\top}Sx + \Sigma_0^{-1})\beta^{(t+1)} = x^{\top}Sz + \mu_0\Sigma_0^{-1}$$
(12)

where:

-
$$S \in \mathbb{R}^{n \times n}$$
 is a diagonal matrix with $S_{ii} = \frac{w_i}{(g'(\mu_i))^2 \phi \mathcal{V}(\mu_i)}$
- $z \in \mathbb{R}^n$ is working response with $z = x^\top \beta^{(t)} + (y - \mu)g'(\mu)$

In the observation that $\Lambda_i \in \mathbb{R}^q$, the expectation could be numerically integrated when the latent dimension q is small. Using the standard Gaussian quadrature method with number of nodes being denoted as m, we can write the Eq (9) as:

$$\mathbb{E}_{\Lambda_i|X_i}(\log f_{\theta^{(t)}}(X_i,\Lambda_i)) = \sum_{l=1}^m \log f_{\theta^{(t)}}(X_i,\Lambda_{il})\widetilde{f}_{\theta^{(t)}}(\Lambda_{il}|X_i)$$
(13)

where $\{\Lambda_{il}\}_{l=1}^{m}$ are the quadrature node of evaluation partitioning on the domain of Λ_i : dom (Λ_i) . To effectively find those quadrature nodes, we can locate the $\{\Lambda_{il}\}_{l=1}^{m}$ such that the integrand function $\{\tilde{f}_{\theta^{(t)}}(\Lambda_{il})\}_{l=1}^{m}$ is non-zero. Those points $\{\Lambda_{il}\}_{l=1}^{m}$ can be defined according to the contour of the density $\tilde{f}_{\theta^{(t)}}(\Lambda_i) = f_N(\hat{\Lambda}_i, H_{\hat{\Lambda}_i})$, which is sufficiently characterized by an exploration of $\hat{\Lambda}$ and $H_{\hat{\Lambda}_i}$ using the Gaussian property.

Remark This Gaussian approximation idea is similar to (Rue et al., 2009). That is, we firstly locate the center $\widehat{\Lambda}_i$ of function $\log f_{\theta^{(t)}}(\Lambda_i|X_i)$, then we explore the *m* points using the hessian $H_{\widehat{\Lambda}_i}$ of function $\log f_{\theta^{(t)}}(\Lambda_i|X_i)$. Those exploration will be exact if $\log(X_i, \Lambda_i)$ is Gaussian and can not be bad under the observation that $\log f_{\theta^{(t)}}(\Lambda_i|X_i) = \log f_{\theta^{(t)}}(\Lambda_i|X_i) + \log f(\Lambda_i)$ and the assumption that $f(\Lambda_i)$ is a Gaussian density.

• M-Step:

We can then solve the M-step by finding the parameters of $\theta^{(t+1)}$ that maximizing $\mathbb{E}_{\widetilde{\Lambda}_t|X_i}(\log f_{\theta^{(t)}}(X_i, \Lambda_i))$:

$$\theta^{(t+1)} \approx \underset{\theta=(V,\Phi,\eta_0)}{\operatorname{argmin}} - \sum_{i=1}^{n} \mathbb{E}_{\widetilde{\Lambda}_i \mid X_i}(\log f_{\theta^{(t)}}(X_i,\Lambda_i))$$
(14)

For this, observe that after the exploration from the E step, we finalize the quadrature nodes $\{\Lambda_{il}\}_{l=1}^{m}$ and can thus treat $\tilde{f}_{\theta^{(t)}}(\Lambda_{il}|X_i)$ as quasi-GLM weight to solve weighted GLM problem below:

$$\theta^{(t+1)} = \underset{\theta=(V,\eta_0)}{\operatorname{argmin}} - \sum_{i=1}^{n} \sum_{l=1}^{m} \log f_{\theta^{(t)}}(X_i, \Lambda_{il}) \widetilde{f}_{\theta^{(t)}}(\Lambda_{il}|X_i)$$
(15)

The update of $\{\phi_j\}_{j=1}^p$ is simplified to be the Pearson residual:

$$\widehat{\phi}_j = \frac{1}{n} \sum_{i=1}^n \frac{(X_{ij} - \widehat{\mu}_{ij})^2}{\mathcal{V}(\widehat{\mu}_{ij})}, j \in [p]$$
(16)

However, both the the complexity and the error of this optimization is proportional to the latent dimension q. In some of the existing literature (Fan et al., 2020), it is often assumed that $q \ll p$ so that the factorization needs to be applied. For other cases that we do need a larger q, we explore two alternative optimization methods utilizing Stochastic Gradient Descent (SGD). The algorithm for this EM optimization is summarized in Algorithm 1.

Algorithm 1: Numerical EM (recommended for small q) Data: $X \in \mathbb{F}^{n \times p}$ with $\mathbb{F} = \mathbb{N}_+, \mathbb{R}_+ \cdots$ Input factorization rank q, number of Gaussian nodes m, and maximum iteration T Initialization Initialize (V_0, η_0) using centered DMF or SVD, Initialize Φ_j using Eq (16) for t=0 : T do compute location $\mu = \widehat{\Lambda}_i$ and scale $\Sigma = H_{\widehat{\Lambda}_i}$ of Gaussian integrator using Eq (10) compute the integrator $\{\widetilde{f}_{\theta^{(t)}}(\Lambda_{il}|X_i)\}_{l=1}^m$ using Eq (13) solve V_t, η_t via weighted GLM problem in Eq (15) solve dispersion Φ_t via pearson residual using Eq (16). identify $\widetilde{V}_T, \widetilde{\eta}_T$ according to Section 2.1 Result: MLE estimator of $\widehat{V} = \widetilde{V}_T, \widehat{\eta}_0 = \widetilde{\eta}_T, \widehat{\Phi} = \Phi_T$

3.2. SGD optimization with large q

Denote $X \stackrel{d}{=} Y$ as random variable X is equivalent to random variable Y in distribution. Under some regularity conditions to allow the exchange of differentiation and integration, the gradient of likelihood can be written as:

$$\nabla_{\theta} \left[\sum_{i=1}^{n} \log(f_{\theta}(X_{i})) \right] \\
= \sum_{i=1}^{n} \frac{\nabla_{\theta}(\int_{\Lambda_{i}} f_{\theta}(X_{i}, \Lambda_{i}) d\Lambda_{i})}{\int_{\Lambda_{i}} f_{\theta}(X_{i}, \Lambda_{i}) d\Lambda_{i}} \\
= \sum_{i=1}^{n} \int_{\Lambda_{i}} \nabla_{\theta} \left[f_{\theta}(X_{i}, \Lambda_{i}) \right] \frac{1}{f_{\theta}(X_{i})} d\Lambda_{i}$$

$$= \sum_{i=1}^{n} \int_{\Lambda_{i}} \nabla_{\theta} \left[\log f_{\theta}(X_{i}, \Lambda_{i}) \right] \frac{f_{\theta}(X_{i}, \Lambda_{i})}{f_{\theta}(X_{i})} d\Lambda_{i}$$

$$= \sum_{i=1}^{n} \mathbb{E}_{\widetilde{\Lambda}_{i}} \left[\nabla_{\theta} \left[\log f_{\theta}(X_{i}, \widetilde{\Lambda}_{i}) \right] \right], \widetilde{\Lambda}_{i} \stackrel{d}{=} \Lambda_{i} |X_{i}|$$

If we can evaluate Eq (17) efficiently and accurately, we could then update the EFM parameters through gradient descent with step size α :

$$\theta_{t+1} = \theta_t - \alpha \nabla_\theta \left[-\sum_{i=1}^n \log(f_V(X_i)) \right]$$
(18)

Ignoring for now the expectation $\mathbb{E}_{\widetilde{\Lambda}_i}(\dot{J}, \nabla_{\theta}[f_{\theta}(X_i, \Lambda_i)]$ are in fact available in closed form given a specification of our model in Eq (3). We provided a summary of those gradient below:

• Gradient for V_i :

$$-\nabla_{V_j} \log f_{\theta}(X_i, \Lambda_i) = \frac{w_{ij}}{\phi_j} \nabla_{V_j} \left(\int_{\mu_{ij}}^{X_{ij}} \frac{1}{\mathcal{V}(t)} (X_{ij} - t) dt \right)$$
$$= -\frac{w_{ij}}{\phi_j} \frac{1}{\mathcal{V}(\mu_{ij})} (X_{ij} - \mu_{ij}) \frac{\partial \mu_{ij}}{\partial \eta_{ij}} \frac{\partial \eta_{ij}}{\partial V_j}$$
$$= -\frac{w_{ij}}{\phi_j} \frac{1}{\mathcal{V}(\mu_{ij})} (X_{ij} - \mu_{ij}) \frac{1}{g'(\mu_{ij})} \Lambda_i$$
(19)

• Gradient for ϕ_i :

Calculating the gradient for ϕ will require a function form of $c(X_{ij}, \Phi_{ij})$, which doesn't entitle easy closed form. But as argued in (Nelder and Pregibon, 1987), little is lost if we focus on the unnormalized version of of quasi-likelihood since the normalization value $c(X_{ij}, \Phi_{ij})$ contains almost no information about main parameter *V*. Ignoring $c(X_{ij}, \Phi_{ij})$, we repeat the derivation

for parameter $\phi_j \in \mathbb{R}_+$:

$$\begin{aligned} \nabla_{\phi_j}(-\log f_{\theta}(X_i, \Lambda_i)) &= -\frac{w_{ij}}{\phi_j^2} \left(\int_{\mu_{ij}}^{X_{ij}} \frac{1}{\mathcal{V}(t)} (X_i - t) dt \right) + \frac{1}{2\phi_j} \\ &= -\frac{w_{ij}}{2\phi_j^2} Q(X_{ij}; \mu_{ij}) + \frac{1}{2\phi_j}, \forall j \in [p] \\ &= \frac{1}{2\phi_j} \left(-\frac{w_{ij}}{\phi_j} Q(X_{ij}; \mu_{ij}) + 1 \right), \forall j \in [p] \\ &\approx \frac{1}{2\phi_j} \left(-\frac{w_{ij} (X_{ij} - \mu_{ij})^2}{\phi_j \mathcal{V}(\mu_{ij})} + 1 \right), \forall j \in [p] \end{aligned}$$
where $Q(X_{ij}; \mu_{ij}) = -2 \int_{\mu_{ij}}^{\mu_{ij}} \frac{1}{2\mathcal{V}(t)} (X_i - t) dt = Q(X_{ij}; X_{ij}) - Q$

where $Q(X_{ij}; \mu_{ij}) = -2 \int_{X_{ij}}^{\mu_{ij}} \frac{1}{V(t)} (X_i - t) dt = Q(X_{ij}; X_{ij}) - Q(X_{ij}; \mu_{ij})$ is the quasi-deviance function.

• Gradient for η_0 : Ignoring $c(X_{ij}, \Phi_{ij})$, we repeat the derivation for $\eta_{0j} \in \mathbb{R}$:

$$\nabla_{\eta_{0j}}(\log f_{\theta}(X_{i},\Lambda_{i})) = \frac{w_{ij}}{\phi_{j}} \nabla_{\eta_{0j}} \left(\int_{\mu_{ij}}^{X_{ij}} \frac{1}{\mathcal{V}(t)} (X_{i}-t) dt \right)$$
$$= -\frac{w_{ij}}{\phi_{j}} \frac{1}{\mathcal{V}(\mu_{ij})} (X_{ij} - \mu_{ij}) \frac{\partial \mu_{ij}}{\partial \eta_{ij}} \frac{\partial \eta_{ij}}{\partial \eta_{0j}}$$
$$= \frac{1}{\Phi_{jj}} \frac{1}{\mathcal{V}(\mu_{ij})g'(\mu_{ij})} (\mu_{ij} - X_{ij}), \forall j \in [p]$$
(21)

In our later optimization, we additionally need the Hessian of the log $f_{\theta}(X_i, \Lambda_i)$.

• Hessian for V_i and Λ_i :

The Hessian for V_j and Λ_i are symmetric with respect to each other with a simple notation change, below we use H_{V_j} as an example:

$$\mathbb{E}_{X_{i}}\left[\nabla_{V_{j}}^{(2)}(-\log f_{\theta}(X_{i},\Lambda_{i}))\right] = \\\mathbb{E}_{X_{i}}\left[\left(\nabla_{V_{j}}\left(G(\mu_{i})^{\top}[S^{\frac{1}{2}}(\mu_{i})\Phi S^{\frac{1}{2}}(\mu_{i})]^{-1}\right)(\mu_{i}-X_{i})+\right. \\\left.G(\mu_{i})^{\top}[S^{\frac{1}{2}}(\mu_{i})\Phi S^{\frac{1}{2}}(\mu_{i})]^{-1}\nabla_{V_{j}}(\mu_{i}-X_{i})\right)\Lambda_{i}\right] (22) \\= \sum_{j=1}^{p}G(\mu_{ij})^{2}[S^{\frac{1}{2}}(\mu_{ij})\Phi_{jj}S^{\frac{1}{2}}(\mu_{ij})]^{-1}\frac{\partial\mu_{ij}}{\partial\eta_{ij}}\frac{\partial\eta_{ij}}{\partial V_{j}}\Lambda_{i} \\= \Lambda_{i}\left(\operatorname{Diag}(G(\mu_{i})^{2})[S^{\frac{1}{2}}(\mu_{i})\Phi S^{\frac{1}{2}}(\mu_{i})]^{-1}\right)\Lambda_{i}^{\top}$$

where the first components is 0 because $\mathbb{E}(\mu_i - X_i) = \mathbf{0}_p$.

The random sampling for the gradient of every observation $i \in [n]$ in Eq (17) is however expensive. We here additionally leverage modern stochastic optimization to randomly sample partial data to approximate the optimization gradient. The optimization is termed Stochastic Gradient Descent (SGD) due to its interpretation as a stochastic approximation of the actual gradient function. Since the method effectively reduces the sample size by applying a subsampling on the original dataset, the SGD can be used to accelerate all optimization algorithms with explicit gradient formulation. We here use the SML optimization as an example to illustrate the implementation.

Taking last equality from Eq (17) and applying the law of large number, we can compute the gradient using stochastic sample of size *B* and *S*:

$$\sum_{i=1}^{n} \nabla_{\theta} \left(\log(f_{\theta}(X_{i})) \right) \approx n \mathbb{E}_{X} \left[\nabla_{\theta} \left(\log(f_{\theta}(X)) \right) \right]$$

$$= n \mathbb{E}_{X^{(B)}} \left[\nabla_{\theta} \left(\log(f_{\theta}(X_{b})) \right) \right], X^{(B)} \stackrel{d}{=} X$$

$$\approx \frac{n}{B} \sum_{b=1}^{B} \nabla_{\theta} \left(\log(f_{\theta}(X_{b}^{(B)})) \right)$$

$$= \frac{n}{B} \sum_{b=1}^{B} \mathbb{E}_{\widetilde{\Lambda}_{b}} \left[\nabla_{\theta} \left[\log f_{\theta}(X_{b}^{(B)}, \widetilde{\Lambda}_{b}) \right] \right], \widetilde{\Lambda}_{b} \stackrel{d}{=} N(\mathbf{0}_{q}, I_{q})$$

$$\approx \frac{n}{B} \sum_{b=1}^{B} \sum_{s=1}^{S} \frac{1}{S} \nabla_{\theta} \left[\log f_{\theta}(X_{b}^{(B)}, \widetilde{\Lambda}_{b,s}^{(S)}) \right], \widetilde{\Lambda}_{b,s} \stackrel{d}{=} N(\mathbf{0}_{q}, I_{q})$$

The batch size *B* is chosen to be smaller than sample size *n*, which scales the original complexity with a factor of $\frac{B}{n}$ per iteration. To maintain the relationship $X^{(B)} \stackrel{d}{=} X$, one can sample with replacement from the original data *X*. In addition, since this stochastic sampled gradient is proposed to maximize the true likelihood instead of the simulated likelihood, the sample size *S* required to compute the optimization gradients does not need to grow in an order of the actual sample size (e.g. $S = n^{1/2}$ as required for SML(Lee, 1995)).

As it is compared to second-order optimization such as Newton's method, step size selection is of crucial importance for stochastic gradient descent. A large step size will make the algorithm oscillate while a small step size hardly improves our likelihood function. The theoretical analysis states that we should choose the step size according to the conditional number of the parameter hessian matrix (Bertsekas, 1999). When the hessian matrix is not available, recent researchers have refined the step size selection by utilizing the momentum and the scale of the parameters. Specifically, AdaGrad (Duchi et al., 2011) scales the gradient update with the gradient's second moment while the RMSProp (Ieleman and Hinton, 2012) employed exponential decay to smooth out the gradient direction. More recently, combining both RMSprop and AdaGrad, Adam (Kingma and Ba, 2015) method has gained its well-deserved attention in modern stochastic optimization research. Here we adopt the Adam method with the parameter recommended in the original Adam paper ($\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$). To avoid the oscillation around the minimum, we also employed a decay learning rate with $\gamma_t = \frac{\alpha}{1+0.5t}$. When it is necessary, the tuning on the hyperparameter α can be conducted by randomly sampling α on a log grid.

To tackle specifically the expectation $\mathbb{E}_{\widetilde{\Lambda}_i}(\cdot)$ with $\widetilde{\Lambda}_i \stackrel{d}{=} \Lambda_i | X_i$, we

propose two optimization algorithms in the following subsections.

3.2.1. Computing the gradient using laplacian approximation

The evaluation of the gradient is thus equivalently an evaluation on the posterior moment of general function $g(\cdot)$. To such integration, laplacian approximation (Tierney and Kadane, 1986) has been frequently studied. Denote

•
$$g(\Lambda_i) : \mathbb{R}^q \to \mathbb{R}^q, g(\Lambda_i) = \nabla_{\theta} [\log f_{\theta}(X_i, \Lambda_i)]$$

• $h(\Lambda_i) : \mathbb{R}^q \to \mathbb{R}, h(\Lambda_i) = -\frac{1}{p} \sum_{j=1}^p \log f_{\theta}(X_{ij}|\Lambda_i) + f(\Lambda_i)$

It is easy to verify that $h(\Lambda_i)$ is a constant order function of p as $p \to \infty$ and that $g(\Lambda_i)$ does not growth with p. Interestingly, even if $g(\Lambda_i)$ is permitted to grow with p, provided that the growth rate is bounded by $O(e^{g_0p^{1-\delta}})$, a valid approximation can still be derived with an error estimate of $O(e^{p^{1-\delta}})$. The precise approximation in this general scenario can be found in Lemma 5 of (Xia and Zhang, 2023). The evaluation follows from the derivation below with $\Lambda_i = \Lambda_i | X_i$ being the posterior:

$$\sum_{i=1}^{n} \mathbb{E}_{\widetilde{\Lambda}_{i}}[g(\widetilde{\Lambda}_{i})] = \int_{\mathbb{R}^{q}} g(\Lambda_{i}) f(\Lambda_{i}|X_{i}) d\Lambda_{i}$$
$$= \frac{\int_{\mathbb{R}^{q}} g(\Lambda_{i}) \exp(-ph(\Lambda_{i})) d\Lambda_{i}}{\int_{\mathbb{R}^{q}} \exp(-ph(\Lambda_{i})) d\Lambda_{i}}$$
(24)

To further simplify the notation, we denote $H_{\Lambda_i} = \nabla_{\Lambda_i}^{(2)} h(\widehat{\Lambda}_i)$, and $U_{\Lambda_i} = \nabla_{\Lambda_i} h(\widehat{\Lambda}_i)$. If we choose $\widehat{\Lambda}_i$ to maximize $-h(\Lambda_i)$, we will have the gradient $U_{\Lambda_i} = \nabla_{\Lambda_i} h(\widehat{\Lambda}_i) = \mathbf{0}_q$ and the numerator can be simplified with the leading term (Tierney et al., 1989):

$$\int_{\mathbb{R}^{q}} g(\Lambda_{i}) \exp(-ph(\Lambda_{i})) d\Lambda_{i}$$

= $g(\widehat{\Lambda}_{i}) \exp(-ph(\widehat{\Lambda}_{i})) \int_{\mathbb{R}^{q}} \exp(-\frac{p}{2}(\Lambda_{i} - \widehat{\Lambda}_{i})^{\mathsf{T}} H_{\Lambda_{i}}(\Lambda_{i} - \widehat{\Lambda}_{i}))$
= $g(\widehat{\Lambda}_{i}) \exp(-ph(\widehat{\Lambda}_{i}))(2\pi/p)^{q/2} |\Sigma(\widehat{\Lambda}_{i})|^{-1/2} (1 + O(1/p))$
(25)

As a corollary of this O(1/p) approximation result, the denominator of $\frac{\int_{\mathbb{R}^{q}} g(\Lambda_{i}) \exp(-ph(\Lambda_{i}))d\Lambda_{i}}{\int_{\mathbb{R}^{q}} \exp(-ph(\Lambda_{i}))d\Lambda_{i}}$ is a special case of the numerator with $g(\Lambda_{i}) = 1$, a more accurate approximation with a relative order $O(1/p^{2})$ can be obtained by applying approximation in Eq (25) to both the numerator and denominator. Such an approximation however would require additionally either of the following components:

- The higher order of derivative $\{\nabla_V^{(k)}g(\widehat{\Lambda}_i)\}_{k=1}^2$ and $\nabla_{\Lambda_i}^{(3)}g(\widehat{\Lambda}_i)$.
- The maximization solution of $\widehat{\Lambda}_i = \operatorname{argmin}_{\Lambda_i} \frac{1}{p} \log g(\Lambda_i) + h(\Lambda_i).$

The approximation with the first evaluation is named as Standard Form while the one with the second evaluation is named as the Fully Exponential Form (Tierney et al., 1989).

However, Eq (25) still requires an evaluation of $h(\widehat{\Lambda}_i) = -\frac{1}{p} \sum_{j=1}^{p} \log(f_{\theta}(X_{ij}|\Lambda_i))$, which will asymptotically approach 0 with bad initialization of *V* as $p \to \infty$. Fortunately, for an order of O(1/p) approximation, a joint approximation of the denominator and numerator integral with the Standard Form (Tierney et al., 1989) suggests $g(\widehat{\Lambda}_i)$ would equivalently provide O(1/p) approximation by simply plugging $\widehat{\Lambda}_i = \operatorname{argmin}_{\Lambda_i} h(\Lambda_i)$ into function $h(\cdot)$:

$$\sum_{i=1}^{n} \mathbb{E}_{\widetilde{\Lambda}_{i}}[g(\widetilde{\Lambda}_{i}))] = \sum_{i=1}^{n} g(\widehat{\Lambda}_{i}) (1 + O(1/p))$$
(26)

To facilitate later reference, we name this optimization as Laplacian optimization, whose gradient evaluation is accurate for large pwith a relative error rate of O(1/p). Although our setup assumes a Gaussian latent variable, the optimization can be generally applied to non-Gaussian latent prior with a simple modification of $f(\Lambda_i)$ in

$$h(\Lambda_i) = \frac{1}{p} \sum_{j=1}^{p} \log f_{\theta}(X_{ij}|\Lambda_i) + f(\Lambda_i)$$

The optimization can be readily accommodated as weighted MAP solution of Bayesian GLM regression for the given prior of Λ_i with density $f(\Lambda_i)$:

$$\widehat{\Lambda}_{i} = \underset{\Lambda_{i}}{\operatorname{argmax}} \frac{1}{p} \sum_{j=1}^{p} \log f_{\theta}(X_{ij} | \Lambda_{i}) + f(\Lambda_{i})$$
(27)

3.2.2. Computing the gradient using posterior sampling

When *p* is of moderate dimension, the evaluation of Eq (2) still would require a good starting point, yet the gradient evaluation using Eq (25) is inaccurate. In this case, we observe from Eq (17) that if the posterior distribution of $\Lambda_i | X_i$ can be easily simulated without numerical issue, then computing the Eq (17) using stochastic sampling will be handy. Our second optimization idea thus comes from approximating the posterior distribution instead of approximating the likelihood gradient.

If we conduct Taylor expansion of the $f_{\theta}(X_i|\Lambda_i)$ to the second order around the stationary point of $\widehat{\Lambda}_i$, we can obtain the following data likelihood approximation:

$$\log(f_{\theta}(X_{i}|\Lambda_{i})) \approx \log(f_{\theta}(X_{i}|\widehat{\Lambda}_{i})) + (\Lambda_{i} - \widehat{\Lambda}_{i})^{\top} U_{\Lambda_{i}}(\widehat{\Lambda}_{i}) + \frac{1}{2} (\Lambda_{i} - \widehat{\Lambda}_{i})^{\top} H_{\Lambda_{i}}(\widehat{\Lambda}_{i}) (\Lambda_{i} - \widehat{\Lambda}_{i}) \qquad (28)$$
$$\propto (\Lambda_{i} - \widehat{\Lambda}_{i})^{\top} H_{\Lambda_{i}}(\widehat{\Lambda}_{i}) (\Lambda_{i} - \widehat{\Lambda}_{i})$$

The required stationary point of $\widehat{\Lambda}_i$ can be obtained by solving the following normal equation for *n* rows of *X* with index i = 1, ..., n:

$$V^{\top}D_{i}.V\lambda_{i}^{(t+1)} = V^{\top}D_{i}.(V\lambda_{i}^{(t)} + D_{i}^{-1}G_{i}.) = V^{\top}D_{i}.Z_{i}^{(t)}, \quad (29)$$

where λ_i denotes the *i*-th row of matrix Λ . ϕ is the dispersion parameter from exponential family. $Z^{(t)}, D_{i}^{-1}, G_{i}$ are defined as:.

•
$$S_{ij} = \frac{w_{ij}}{\phi_j} \frac{g^{-1'}(\eta_{ij})^2}{V(\mu_{ij})}$$
 and $G_{ij} = \frac{g^{-1'}(\eta_{ij})}{V(\mu_{ij})} \frac{w_{ij}}{\phi_j} (X_{ij} - \mu_{ij})$

- $D_{i.} \stackrel{d}{=} \text{Diag}\{S_{i.}^{(t)}\}$ with $S_{i.}$ denotes the *i*-th row of *S*. Similarly, $\text{Diag}\{S_{\cdot i}^{(t)}\}$ with $S_{\cdot j}$ denotes the *j*-th column of *S*.
- $Z^{(t)}$ is the working response:

$$Z_{ij}^{(t)} = \eta_{ij}^{(t)} + \frac{G_{ij}^{(t)}}{S_{ij}^{(t)}} = \eta_{ij}^{(t)} + \frac{X_{ij} - \mu_{ij}^{(t)}}{g^{-1'}(\eta_{ij}^{(t)})}.$$
 (30)

After which, the posterior distribution of the latent variable Λ can be approximated with the Gaussian formula:

$$f(\Lambda_i|X_i) \propto f(X_i|\Lambda_i) f(\Lambda_i) \approx f_N(\widehat{\Lambda}_i, H_{\Lambda_i}^{-1}(\widehat{\Lambda}_i)) f_N(0, I_q)$$

= $f_N(\mu_i, \Sigma_i)$ (31)

where $f_N(\mu, \Sigma)$ is the multivariate normal density with mean μ and variance Σ . We have from the Gaussian integration formula that $\mu_i = [I_q + H_{\Lambda_i}^{-1}(\widehat{\Lambda}_i)]^{-1}\widehat{\Lambda}_i, \Sigma_i = [I_q + H_{\Lambda_i}(\widehat{\Lambda}_i)]^{-1}$. Under this closed-form expression of the posterior, the complex-

Under this closed-form expression of the posterior, the complexity of evaluating the gradient Eq (23) becomes as small as sampling from a multivariate normal distribution with parameter μ_i , Σ_i . As for the quality of this approximation, we appeal to a similar statement in (Rue et al., 2009) that we are implicitly assuming the shape of the posterior is determined solely by the prior and the likelihood only contributes to the location and scale. The approximation can be inaccurate when the prior is non-Gaussian but this is not the case for the factor model where $f(\Lambda_i) \sim N(0, I_q)$.

For the convenience of later reference, this second optimization method is named as Posterior Sampling optimization. This posterior sampling optimization is in fact closely related to MCEM algorithm (Dempster et al., 1977), whose convergence result is proven to be superior compared to SML(Jank and Booth, 2003). To observe this connection, notice that we can introduce a new probability measure $Q(\Lambda_i)$ to reformulate the optimization problem in Eq (6) as:

$$\widehat{V} = \operatorname*{argmax}_{\theta} \sum_{i=1}^{n} \log \int_{\Lambda_{i}} f_{\theta}(X_{i}|\Lambda_{i}) f(\Lambda_{i}) d\Lambda_{i}$$

$$= \operatorname*{argmax}_{\theta} \sum_{i=1}^{n} \log \int_{\Lambda_{i}} \frac{f_{\theta}(X_{i}|\Lambda_{i}) f(\Lambda_{i})}{Q(\Lambda_{i})} Q(\Lambda_{i}) d\Lambda_{i} \qquad (32)$$

$$= \operatorname*{argmax}_{\theta} \sum_{i=1}^{n} \log \left(\mathbb{E}_{Q(\Lambda_{i})} \left[\frac{f_{\theta}(X_{i},\Lambda_{i})}{Q(\Lambda_{i})} \right] \right)$$

Now if we apply Jensen's inequality to switch the order of expectations and log operation:

$$f_{\theta}(X) = \sum_{i=1}^{n} \log \mathbb{E}_{Q(\Lambda_i)} \left[\frac{f_{\theta}(X_i, \Lambda_i)}{Q(\Lambda_i)} \right]$$

$$\geq \sum_{i=1}^{n} \mathbb{E}_{Q(\Lambda_i)} \left[\log(\frac{f_{\theta}(X_i, \Lambda_i)}{Q(\Lambda_i)}) \right]$$
(33)

The derived inequality becomes equality if $\log(\frac{f_{\theta}(X_i, \Lambda_i)}{Q(\Lambda_i)})$ is a constant, which can only be achieved by choosing the posterior $Q(\Lambda_i) \propto f_{\theta}(X_i, \Lambda_i) = f_{\theta}(X_i|\Lambda_i) f(\Lambda_i)$.

Then the "M step" of the EM algorithm optimizes:

$$\underset{\theta}{\operatorname{argmax}} \sum_{i=1}^{n} \mathbb{E}_{\mathcal{Q}(\Lambda_i)} [\log(\frac{f_{\theta}(X_i, \Lambda_i)}{\mathcal{Q}(\Lambda_i)})]$$

whose gradient, under similar regularity conditions to switch the order of derivative and integration, can be obtained in the following form:

$$\sum_{i=1}^{n} \int_{\Lambda_{i}} \nabla_{V} [\log f_{\theta}(X_{i}, \Lambda_{i})] Q(X_{i}) d\Lambda_{i}$$

Hence, our proposed stochastic gradient descent is equivalent to an EM algorithm iteratively maximizing the marginalized likelihood. However, as indicated in (Caffo et al., 2005), this MCEM using simulated gradient for optimization often requires an adaptive change of the sample size S_t to converge successfully. Recent research (Jank, 2006) circumvent the choice of adaptive S_t by averaging the past iterations. The average is oftentimes weighted with an emphasis on the recent iterations, which is in fact equivalent to the step size selection of Adam optimization (Kingma and Ba, 2015).

We summarize those two SGD algorithms in Algorithm 2

Remark Although the iteration ends after *T* passes the data, similar early stopping criteria (Yao et al., 2007) can be adopted if the main objective of the model is to make future predictions. If one hopes to focus on the interpretability of the factorized components, one can stop the algorithm with small V_{t+1} updates.

4. Examples and Results

In this section, we first demonstrate the result of simulation experiments, which validates the effectiveness and superiority of our SGD estimation compared to the SML estimation. Then with benchmark dataset in computer vision and network analysis, we compare our EFM factorization result against other commonly applied factorizations such as Non-negative Matrix Factorization (**NMF**), t-distributed stochastic neighbor embedding (t-SNE) and deviance matrix factorization (**DMF**, Wang and Carvalho, 2023). The factorization ranks q on those empirical datasets are determined according to a rank determination proposition in (Wang and Carvalho, 2023).

4.1. Simulated data

To validate the effectiveness of the optimization algorithm, we applied our optimization to some simulated datasets. We design small, simulated datasets where the marginalized likelihood can be evaluated using SML with a large sample size S. To avoid potential confusion, we denote S as the sample size used to evaluate the gradient and denote R as the sample size used to evaluate the

Algorithm 2: Adam SGD (recommended for large q) **Data:** $X \in \mathbb{F}^{n \times p}$ with $\mathbb{F} = \mathbb{N}_+, \mathbb{R}_+ \cdots$ **Notations** Z standard normal, $\Sigma^{-1/2}$ cholesky decomposition of Σ , \circ element-wise product. Stochastic sample of size *B* and *S*: $X^{(B)}$ and $\Lambda^{(S)}$. **Input** Batch size *B*, Sample size *S*, learning rate α , factorization rank q and maximum iteration T**Initialization** Initialize V_0 , η_0 using centered DMF or SVD, Φ_0 through pearson residual; set Adam param: $V_{d\theta} = \mathbf{0}, S_{d\theta} = \mathbf{0}, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$ **for** *t***=**0 : *T* **do** sample with replacement batch $X^{(B)}$ from data X if if p is large then compute $\widehat{\Lambda}_i$ as the MAP solution defined in Eq (27) compute the gradient $\widetilde{\nabla}_{\theta} = \sum_{i=1}^{n} g(\widehat{\Lambda}_i)$ as Eq (26) else compute Λ_i according to Eq (29) compute $\vec{\mu}(\widehat{\Lambda}_i), \Sigma(\widehat{\Lambda}_i)$ according to Eq (31) draw S samples of $\Lambda^{(S)} = \vec{\mu} + (\Sigma)^{1/2} Z^{(S)}$ compute $\widetilde{\nabla}_{\theta} = \sum_{i=1}^{n} \nabla_{\theta} (\log(f_{\theta}(X_i)))$ via Eq (23) update $V_{d\theta} = \beta_1 V_{d\theta} + (1 - \beta_1) \widetilde{\nabla}_V //$ Adam momentum update $\sum_{d\theta} = \beta_2 S_{d\theta} + (1 - \beta_2) (\widetilde{\nabla}_V \circ \widetilde{\nabla}_V) //$ Adam scale obtain $\tilde{V}_{d\theta} = V_{d\theta} / (1 - \beta_1) //$ bias correct 1st moment obtain $\widetilde{S}_{d\theta} = S_{d\theta} / (1 - \beta_2) //$ bias correct 2nd moment update $V_{t+1} = V_t - \frac{\alpha}{1+0.5t} \frac{V_{d\theta}}{\sqrt{S_{d\theta}+\epsilon}}$ decompose(SVD) $V'_{t+1} = UDS^{\top}$ with $d_1 \ge \cdots \ge d_q > 0$ identify $\tilde{V}_T, \tilde{\eta}_T$ according to Section 2.1 **Result:** MLE estimator of $\widehat{V} = V_T$, $\widehat{\eta}_0 = \eta_T$, $\widehat{\Phi} = \Phi_T$, posterior approximation of the *i*-th latent variable parameter $\vec{\mu}(\Lambda_i), \Sigma(\Lambda_i)$

marginal likelihood. The notation is further clarified according to the marginal likelihood evaluation below:

$$\mathcal{L}(V) = \sum_{i=1}^{n} \log f(X_{i}|V)$$

$$\approx \sum_{i=1}^{n} \log(\sum_{r=1}^{R} \frac{1}{R} f(X_{i}|\Lambda_{i}^{(R)}V^{\top}))$$

$$= \sum_{i=1}^{n} \log(\frac{1}{R} \sum_{r=1}^{R} \exp(\log f(X_{i}|\Lambda_{i}^{(R)}V^{\top})))$$

$$= \sum_{i=1}^{n} \operatorname{LogSumExp}(\{\log f(X_{i}|\Lambda_{i}^{(R)}V^{\top})\}) - n \log(R)$$
(34)

Notice that the evaluation of Eq (34) is not required for the implementation of our Algorithm 2. The evaluation is only introduced to compare the effectiveness and efficiency of those optimization to decrease the integrated negative log-likelihood, which can only be obtained using SML with large *R* for non-Gaussian data likelihood.

To compare the quality of our optimization algorithm with the SML solution, we experimented with simulated data from Negative Binomial ($\phi = 20$), Binomial and Poisson. The parameters are chosen as $n = 500, B = 128, q = 2, p = 10, \alpha = 0.5$ with their canonical link function. Note that the sample size *n* is chosen to be small for an efficient evaluation of the loss function using Monte Carlo. Based upon the true generating parameter V^* , we empirically verify that an accurate evaluation of the likelihood using Eq (34) would require R = 1,500, which is three times larger than the sample size n = 500. This observation is consistent to existing literature yet interesting to practitioners since the common literature is concentrated on the discussion of asymptotic efficiency with a lower bound of $R > \sqrt{n}$ (Lee, 1995). In reality, the Monte Carlo samples required for accurate gradient or likelihood evaluation obviously depend on the variance of the gradient and likelihood of the "specific" dataset, which has no upper bound. We here adopted R = 1,500 to accurately monitor the loss decrease per unit of time with the same initialization point. The evaluation time using sample size R is later subtracted from the optimization time for a fair comparison.

To investigate the dimensionality and variance effect on different optimization algorithms, we first conducted two experiments with p = 5 and p = 10 and then conducted two additional experiments with large p = 512. Notice that the dimension size p = 5, p = 10 are designed to accommodate the numerical stability of SML optimization, which has evaluation issues with large p as mentioned in Section 1.1.3. For each of the optimization, we fix the initialization and the random seed to fairly compare the optimization paths.

We abbreviated LAPL for Laplacian Approximation optimization, PS for posterior sampling optimization, and SML for simulated maximum likelihood optimization. To also investigate the sampling requirement, we varied sample size *S* from $S = \{50, 300, 500\}$. The optimization paths with p = 5 are plotted below: As we can see



Figure 1: EFM Optimization Comparison, p =5

from Figure 1, the Posterior Sampling (PS) optimization is not very sensitive to the choice of sample size *S* while the SML optimization solution varies greatly according to a different choice of sample size with larger *S* leads to faster decrement. The Laplacian approximation decreases the loss function at the slowest speed due to the approximation error of O(1/p) in gradient evaluation

(Eq (24)). The result indicates that the EM and PS optimization should be preferred on small data dimension p due to its efficiency, less sensitivity of sample size S, and numerical stability.

In theory, the LAPL optimization should become more accurate with a relative error with respect to the data dimensionality O(1/p). To observe the effectiveness of LAPL with moderate dimensionality, we continued the same experiments with p = 10. We also adopted S = 500 for both SML and PS optimization to compare the optimization efficiency. The optimization paths are again recorded with respect to the Adam optimization steps: As we can see from



Figure 2: EFM Optimization Comparison, p =10

Figure 2, the posterior sampling optimization is still the most efficient among the three optimizations with more loss descended per unit of time. However, as it reaches the convergence region, the PS optimization contains higher variance as it is compared to the LAPL. This observation indicate that the potential superiority of the Laplacian optimization when the gradient evaluation contains large variance. The SML optimization undoubtedly decreases the negative likelihood at the slowest rate and should not be ever considered for practical applications.

To effectively compare the optimization performance on largedimensional dataset, we experimented with p = 512, under which scenario, the SML completely lost its power due to the numerical stability issue. In fact, we observe that the SML consistently increase the loss with respect to the optimization steps. We thus compare only the PS of different sample sizes S and LAPL optimization with their optimization paths prorated across the optimization time: As we can see from Figure 3, despite the potentially larger Monte Carlo error in gradient evaluation with a smaller sample size S = 50, the PS optimization converges as it is compared to the LAPL and EM optimization. Such a behaviour is very desired as it also indicates low computational budget required for large dimension optimization. However, the PS optimization would potentially require a large S when the variance of the gradient evaluation increases. To validate this argument, we increased the magnitude of $V^* = U^*D^*$ in simulation through a multiplication c > 1 on its diagonal elements. This multiplication will enlarge the magnitude of gradient according to the function relationship in Table 1.

We then continued the data simulation process with large dimension p = 512 to compare the performance between PS and LAPL:

Table 1: Distributions and their respective gradient and hessian functions.

Distribution F	Negative Log-likelihood $l(X)$	Gradient $-\nabla_{V_j} l(X)$	Hessian $-\nabla_{\Lambda_i}^{(2)} l(X)$
Gaussian(identity)	$\frac{1}{2\sigma^2}\sum_{j=1}^{p}(X_j - \Lambda V_j)^{\top}(X_j - \Lambda V_j)$	$-\frac{1}{\sigma^2}\Lambda^{\top}(X_j - \Lambda V_i^{\top})$	$\frac{1}{\sigma^2}V^{\top}V$
Poisson(log)	$\sum_{i=1}^{p} 1_{n}^{\top} \exp(\Lambda V_{i}^{\top}) - X_{i}^{\top} (\Lambda V_{i}^{\top})$	$-X_{i}^{\top}\Lambda + \exp(\Lambda V_{i}^{\top})^{\top}\Lambda$	$V^{T} \mathrm{Diag}[\exp(\Lambda_i V^{T})]V$
Gamma(log)	$-\phi[\sum_{i=1}^{p} X_{i}^{T} \Lambda V_{i}^{T} + \log(-\Lambda V_{i}^{T})]$	$(-1/\Lambda V_i^{T})\Lambda - X_i^{T}\Lambda$	$V^{T}\mathrm{Diag}^{2}[\phi/\Lambda_{i}V^{T}]V$
Binomial(logit)	$\sum_{j=1}^{p} -(w_j \circ X_j^{T})\Lambda V_j^{T} +$	$[w_j/(1 + \exp(-\Lambda V_j^{T}))]^{T}\Lambda -$	$V^{\top}(w_i \circ \text{Diag}[\exp(\Lambda_i V^{\top})$
	$w_{i}^{\dagger} \log(1_{n} + \exp(\Lambda V_{i}^{\dagger}))$	$(w_j \circ X_j)^{\top} \Lambda$	$/(1 + \exp(\Lambda_i V^{T} a))])V$
Negative Binomial(α)	$\sum_{j=1}^{p} -X_{j}^{\top} [\Lambda V_{j}^{\top} -$	$-X_j^{\top}\Lambda +$	V^{T} Diag[exp($\Lambda_i V^{T}$)+
	$\alpha 1_n^\top \log(1_q - \exp(\Lambda V_j^\top))]$	$\alpha \exp(\Lambda V_j^{T})^{T} / (1_n^{T} - \exp(\Lambda V_j^{T})) \Lambda$	$\alpha \exp(2\Lambda_i V^{\top})]V$



Figure 3: EFM Optimization Comparison, p =512 and small V



Figure 4: EFM Optimization Comparison, p =512 and large V

As we can see from Figure 4, when the data demonstrates high variance in the sampled gradients, the PS optimization deteriorates by indicating a higher requirement for the sample size *S*. As a conclusion, the LAPL optimization can should be preferred considering the scenarios of non-Gaussian prior, an even larger dimensionality, and potential high variance in gradient evaluation using posterior sampling. Due to its efficiency demonstrated in the simulation studies, we adopted the Posterior Sampling Optimization for our later

empirical studies with careful monitoring on the loss decrement.

4.2. Covariance modeling

One of the major application of Gaussian factor model is its efficiency in covariance estimation for high dimensional data(Fan et al., 2008). Using similar setup in (Fan et al., 2008), we simulate three quasi-factor data with n = 756, p = [66, 116, ..., 466] and four families (qausi-possion, negative-binomial, binomial, poisson). For comparison, we used the same prior configuration for both factors Λ and projection matrix V according to the setup in Table 1 of (Fan et al., 2008). That is, for each p and family:

• we firstly simulate for $i \in [n], j \in [p]$

$$\begin{split} \Lambda_{i} &\sim N \begin{bmatrix} (0.023558\\ 0.012989\\ 0.020714 \end{bmatrix} \begin{pmatrix} 1.2507 & 0 & 0\\ 0 & 0.31564 & 0\\ 0 & 0 & 0.19303 \end{pmatrix} \end{bmatrix} \\ V_{j} &\sim N \begin{bmatrix} (0.78282\\ 0.51803\\ 0.41003 \end{pmatrix} \begin{pmatrix} 0.029145 & 0.023873 & 0.010184\\ 0.023873 & 0.053951 & -0.006967\\ 0.010184 & 0.006967 & 0.086856 \end{pmatrix} \end{bmatrix} \\ \Phi_{j} &\sim \begin{cases} Gamma(\alpha = 4.0713, \beta = 0.1623) & \text{for quasi-possion} \\ 1 & & \text{for others} \end{cases} \\ w_{ij} &\sim \begin{cases} Poisson(20) & \text{for binomial} \\ 1 & & \text{for others} \end{cases} \end{split}$$

- condition on simulated (Λ, V, Φ) , we generate $X \in \mathbb{F}^{n \times p}$ using the four quasi-family with quasi-density f satisfying Eq (3).
- based upon generated data X, quasi-family defined by density f and prior of Λ in (35), we estimate $\theta = (V, \Phi, \Lambda | X)$ by solving Eq (6).
- we compute covariance via
 - (a) naive sample covariance $\widehat{\Sigma}_{sam} = \frac{XX^{\top}}{n-1} \frac{XI^{\top}X^{\top}}{n(n-1)}$.
 - (b) total covariance $\widehat{\Sigma}$ by plugging the estimated θ into Eq (7).



Figure 5: Estimation Error on Covariance

- (c) true covariance Σ by plugging the actual V, Φ and Λ prior into Eq (7).
- we compute the error of covariance estimation via:
 - (a) Frobenius norm $\|\Sigma \widehat{\Sigma}\|_F$

 - (b) Entropy loss $\operatorname{Tr}(\widehat{\Sigma}\Sigma^{-1}) \log |\widehat{\Sigma}\Sigma^{-1}| p$ (c) Normalized loss $\frac{1}{\sqrt{p}} \|\Sigma^{-1/2}(\widehat{\Sigma} \Sigma)\Sigma^{-1/2}\|_F$
- we repeat the above process k times.

With k = 5, we have obtained the estimation error of covariance accordingly:

As we observe from Figure 5, we have obtained similar Non-Gaussian covariance estimation error as in it was shown previously (Fan et al., 2008) for the Gaussian case. Judging from the 12 normalized and 12 entropy distance, our EM optimization estimates the covariance matrix of high dimensional data in a much more accurately manner as it is compared to the naive estimation using covariance formula.

4.3. Computer vision data

To illustrate the advantages that our EFM can provide more representative factorization, we also conducted experiments on computer vision dataset. Perhaps one of the most popular computer vision dataset is the MNIST dataset, which contains 70,000 handwriting pictures labeled from 0 to 9. However, the modern machine learning researchers have evidenced that the classification task of the MNIST dataset might be too simple in the sense that an appropriately tuned classical machine learning algorithm can easily achieve 97% accuracy¹. With also 70,000 pictures of 10 classes, the Fashion-MNIST dataset (Xiao et al., 2017) is proposed to replace the original dataset by constructing a more complicated classification problem. We here examine our EFM factorization when applied to the Fashion-MNIST dataset and compare our factorized components against DMF, NMF, and t-SNE.

To determine the rank, we adopted the rank determination proposition in (Wang and Carvalho, 2023). The resulting eigenvalue plot in Figure 6 indicates a potential ranks three or seven for the factorization. We adopt rank three for visualization convenience and the simulated likelihood demonstrates convergence result after 15 epochs with $B = 256, S = 50, \alpha = 0.5$. To illustrate the superior capability for model generalization, we used only the first 2,000 samples of the 70,000 data as our training dataset to estimate V. After obtaining this V, we conduct penalized GLM regression based

¹https://paperswithcode.com/sota/image-classification-on-mnist



Figure 6: Simulated Likelihood and Eigenvalue Gap for MNIST

upon another 2,000 sampled testing set X to estimate $\widehat{\Lambda}_i | X$. Those $\widehat{\Lambda}_i | X$ can be considered as the out of sample latent estimation based upon EFM estimated \widehat{V} . Due to the factorization algorithm setup, the t-SNE and NMF result are based up on the 2,000 training dataset and the DMF and EFM results are obtained on the separate 2,000 testing set. We summarize the factorized result in Figure 7.

As it is indicated by the Fashion-MNIST factorization result, both our EFM and the t-SNE methods indicates great separability on the 10 classes with some mistakes on pullover, shirt and coat, which are actually similar classes when we look at the image representation. The NMF performs the worst by having utilizing only two dimensions to separate the 10 different classes. Without the stochastic optimization and regularization, the DMF performs no better than our EFM on the testing samples with the potential reason of overfitting.

To illustrate that our EFM method also provide reasonable uncertainty quantification, we also conducted experiments on the ORL face dataset, which contains 40 subjects with pictures taken under 10 different conditions. Those pictures are in 64x64 pixel dimension, which provides high resolution pixel for restoration. To add uncertainty to the image restoration, we cropped part of the face image by setting the pixels to 0. For example, in Figure 8, we can see that the mouths of the person is covered with white background pixel. Similar to the eigen-face decomposition, We adopt q = 41for factorization since we expect to find 40 individual face eigenvectors and one "average face". We fit EFM with negative binomial and estimated $\hat{\phi} = 10.8861$ by using the moment estimator. For comparison, we also conducted eigen-face restoration with rank 40 after centering the "average face". The result is summarized in Figure 9:

Our EFM not only restore the faces much more accurately compared to eigen-face, but also could quantify the uncertainty in the image restoration. With laplacian approximation, we compute the MAP estimator of $\widehat{\Lambda}_i$ and $\widehat{\Sigma}_i$ and simulate the latent variable $\Lambda_i^{(s)}$ accordingly. The simulated $\Lambda_i^{(s)}$ can then be combined with \widehat{V} to provide simulated human faces. As it is shown in Figure 10, the uncertainty due to the crop of the image is correctly identified after centering the simulated faces, which demonstrates different mouth characteristics.

4.4. Network analysis

Another interesting application of our weighted EFM is on social network analysis. One of the popular interest in this field is to summarize the large adjacency matrix using lower-dimensional representations. Such lower-dimensional representations are named nodes embedding and can be further applied for many statistical inference (e.g. community detection, link prediction). Recently, emerging interests has been directed to embedding inference based upon multiple networks. Those multiple networks are formally observed as multiple interaction graphs, $\{\mathcal{G}^{(1)}, \ldots, \mathcal{G}^{(k)}\}$, which consists of multiple edge relationships, $\{\mathcal{E}^{(1)}, \ldots, \mathcal{E}^{(k)}\}$, for the same sets of vertices *V*. This emerging field of research is named as multi-layer or multiplex network analysis (Kivelä et al., 2014). Denoting the number of vertices as n = |V|, the multiplex network inference starts by transforming those multiple graphs into adjacency matrices $\{A^{(1)}, \ldots, A^{(k)}\}$ of same the dimension of size $n \times n$ with the diagram provided below:

Factorization and inference jointly on those constructed $\{A^{(1)}, \ldots, A^{(k)}\}$ have been shown to provide better nodes embedding with potential application of community detection and link prediction (Wang et al., 2019; Jones and Rubin-Delanchy, 2020). One method to enable the joint inference on adjacency matrices is to effectively combine $\{A^{(1)}, \ldots, A^{(k)}\}$ into a single adjacency matrix \widetilde{A} . Chapter two of (Draves, 2022) provided a decent introduction to various aggregation techniques. To briefly summarize some of the relevant aggregation techniques used in this section, there are

• Average Adjacency Spectral Embedding(AASE) (Tang et al., 2018) which simply average the adjacency matrix through

$$\widetilde{A}_{ij} = \frac{1}{k} \sum_{l=1}^{k} A_{ij}^{(l)}$$
(36)

• Unfolded Adjacency Spectral Embedding (UASE) that concatenate the adjacency matrices columns wisely

$$\widetilde{A} = [A^{(1)}, A^{(2)}, \dots, A^{(k)}] \in \mathbb{Z}^{n \times nk}$$
(37)

• Omnibus Embedding (OE) by constructing pair-wisely the following adjacency matrix:

$$\widetilde{A} = \begin{bmatrix} A^{(1)}, & \frac{1}{2}(A^{(1)} + A^{(2)}), & \dots, & \frac{1}{2}(A^{(1)} + A^{(k)}) \\ \frac{1}{2}(A^{(1)} + A^{(2)}), & A^{(2)}, & \dots, & \frac{1}{2}(A^{(2)} + A^{(k)}) \\ \vdots, & \vdots, & \vdots, & \vdots \\ \frac{1}{2}(A^{(k)} + A^{(1)}), & \frac{1}{2}(A^{(k)} + A^{(2)}), & \dots, & A^{(k)} \end{bmatrix}$$
(38)

With definition $D = \text{Diag}(d_1, \ldots, d_q), d_1 \ge \ldots \ge d_q$ and $S_{n,q}$ as the space of $n \times q$ matrix that has orthogonal columns, we can conduct SVD for the asymmetric \widetilde{A} from UASE and eigen-



Figure 7: Fashion-MNIST Factorization result(L to R shows results for EFM, DMF, NMF and t-SNE)



Figure 8: Cropped ORL Face

decomposition for the symmetric \widetilde{A} from ASE, OE:

$$\widetilde{A}^{(ASE)} = UDU, U \in S_{n,q}$$

$$\widetilde{A}^{(UASE)} = UDV^{\top}, U \in S_{nk,q}, V \in S_{n,q}$$

$$\widetilde{A}^{(OE)} = UDU, U \in S_{nk,q}$$
(39)

The nodes embedding Λ can then be defined as $\Lambda = UD^{1/2}$ with dimension $n \times q$.

However, one obvious shortcomings of such an aggregation is that the factorization implicitly assumes an equal contribution from each of the layers. In reality, we know that each layer of the network is at least different according to different level of sparsity. Treating equally the interaction in a dense graph and the interaction in a sparse graph is problematic by over-emphasizing the interactions on the dense graphs. Additionally for the temporal network that consists of the edge relationship of same vertices across different time-step, it intuitively makes more sense if we apply higher weights to the more recent adjacency matrices compared to an equal aggregation on those edge relationship since those recent adjacency matrices are more powerful in the prediction of future interaction.

Naturally one immediate improvement to an equal aggregation of those individual networks is to apply different weights through



Figure 9: Restored Face: Left(EFM), Right(Eigen-Face)

the factor inference. The EFM provides solution to this aggregation technique by allowing entry-wise/layer-wise weight to the aggregated interaction. With this flexibility on heuristic weight specification, we demonstrate that we can factorize to obtain improved embedding results for multiplex network analysis.

We thus explore our EFM inference on the AUCS dataset ² (Dickison et al., 2016). The dataset records interactions among its 61 employees (n = 61) at the Department of Computer Science at Aarhus University. As we confirmed with the author, among those 61 employees, there are 55 employees with labels from one of the eight research groups. There are 6 employees that do not belong to any of the eight research groups. The research group labels of those 55 employees can thus be treated as known community structure to validate the effectiveness of community inference.

For the whole community, interaction are recorded according to five different online and offline relationships, whose adjacency

²https://manliodedomenico.com/data.php



Figure 10: Simulated ORL Face after Centering

matrices can be represented with the following diagram:



As we observe from the interaction network data, the co-author network is definitely more sparse as it is compared to other layer of the network. By following the heuristic introduced in the beginning of this section, we propose to weight each interaction of different layer of the network according to the sparsity of the layer. Specifically, with $\lambda_{max}^{(k)}$ denoted as the largest eigen-value of adjacency matrix $A^{(k)}$, we

- weight each of the interaction according to $1/\lambda_{max}^{(k)}$ to ensure each matrix has its largest eigen-value equal to 1.
- weight all the diagonal terms with value 0 since the zero interaction of a node with respect to itself does not necessarily contain any information

- weight all the remaining zero terms as the minimal value of the non zero terms in the weight matrix since the interaction is usually sparse with bias toward 0.
- assign value 1 to the interaction/adjacency matrix A_{ij} to value 1 as long as there exists an interaction between node *i* and node *j*.

The more formal mathematical definition is provided below:

$$W_{ij} = \begin{cases} \sum_{l=1}^{k} \frac{1}{\lambda_{max}^{(l)}} A_{ij}^{(l)}, & \text{if } A_{ij}^{k} = 1\\ 0, & \forall i = j\\ \min(\{W_{ij}, W_{ij} \neq 0\}), & \forall A_{ij} = 0 \end{cases}$$
(40)
$$A_{ij} = \begin{cases} 1, & \text{if } \exists \ k \ \text{s.t } A_{ij}^{k} = 1\\ 0, & \text{if } \forall \ k \ \text{s.t } A_{ij}^{k} = 0 \end{cases}$$

With weight *W* and adjacency matrix *A* definition in Eq (40), we then apply binomial EFM with logit link to obtain three dimensional nodes embedding. We can then visualize those factorized embedding according the separability of the known research group community labels. For a comparison to the {AASE, UASE, OM} embedding techniques, SVD or eigen-decomposition are also conducted on their corresponding aggregated adjacency matrix \tilde{A} to obtain their corresponding nodes embedding defined in Eq (39). For the Omnibus Embedding, we choose the first *n* vectors since it has dimension *nk*. The visualization comparison on those factorized embedding $\Lambda = UD^{1/2}$ is provided below:

As we can see from Figure 11 that the weighted EFM separated more research groups as it is compared to a naive SVD on any of the embedded graphs. The classification result is also consistent with the existing literature (Magnani et al., 2021) who have claimed that there are five major research groups identified by publisher of the dataset.

5. Conclusion

We propose a EFM with an efficient optimization algorithm. The model assumption is justified from orientational statistics and thus provide more representative factorized result for many interesting application. The optimization algorithm improves the simulated likelihood estimation (SML) by eliminating the asymptotic estimation bias with moderate simulation sample size *S*. Additionally, utilizing the SGD optimization, our EFM generalizes better than alternative factorization models such as DMF. Both the simulation studies and empirical studies provide compelling evidence to those advantages.

References

J. Kevin Ford, Robert C. MacCallum, and Marianne Tait. The application of exploratory factor analysis in applied psychology: A critical review and analysis. *Personnel Psychology*, 1986.



Figure 11: AUCS embedding visualization comparison

- Simon J D Prince, James H Elder, Jonathan Warrell, and Fatima M Felisberti. Tied factor analysis for face recognition across large pose differences. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2008.
- Eugene F Fama and Kenneth R. French. A five-factor asset pricing model. *Journal of Financial Economics*, 2015.
- Tianchen Xu, Ryan T., and Demmer Gen Li. Zero-inflated poisson factor model with application to microbiome read counts. *Biometrics*, 2021.
- Alexander Basilevsky. Statistical factor analysis and related methods. *Wiley Series in Probability and Mathematical Statistics*, 1994.
- Michel Wedel and Wagner A. Kamakura. Factor analysis with (mixed) observed and latent variables in the exponential family. *Psychometrika*, 66:513–30, 2001.
- Michel Wedel, Ulf Bo Ckenholt, and Wagner A. Kamakurac. Factor models for multivariate count data. *Journal of Multivariate Analysis*, 2003.
- David J Bartholomew, Martin Knott, and Irini Moustaki. *Latent* variable models and factor analysis: A unified approach. John Wiley & Sons, 2011.
- Michael E. Tipping and Christopher M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, Series B*, 1998.
- Theodore Wilbur Anderson. The use of factor analysis in the statistical analysis of multiple time series. *Psychometrika*, 1963.
- Liang Wang and Luis Carvalho. Deviance matrix factorization. *Electronic Journal of Statistics*, 17(2):3762–3810, 2023.
- Jianqing Fan, Yingying Fan, and Jinchi Lv. High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics*, 147(1):186–197, 2008.

- Michael Borenstein, Larry V Hedges, Julian PT Higgins, and Hannah R Rothstein. A basic introduction to fixed-effect and randomeffects models for meta-analysis. *Research synthesis methods*, 1 (2):97–111, 2010.
- Chris K Carter and Robert Kohn. Markov chain monte carlo in conditionally gaussian state space models. *Biometrika*, 83(3): 589–601, 1996.
- Alexander Shapiro. Identifiability of factor analysis: Some results and open problems. *Linear Algebra and its Applications*, 70:1–7, 1985.
- Henry F Kaiser. The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23(3):187–200, 1958.
- Rémi Gribonval and Karin Schnass. Dictionary identification—sparse matrix-factorization via *l*_1-minimization. *IEEE Transactions on Information Theory*, 56(7):3523–3539, 2010.
- Daniel D. Lee and Seung Hoe Choi. Learning the parts of objects by nonnegative matrix factorization. *Nature*, 401, 1999.
- Zhao Li, Xindong Wu, and Hong Peng. Nonnegative matrix factorization on orthogonal subspace. *Pattern Recognition Letters*, 31 (9):905–911, 2010.
- Yunqian Ma and Yun Fu. *Manifold learning theory and applications*, volume 434. CRC press Boca Raton, 2012.
- Richard Socher Jeffrey Pennington and Christopher Manning. Glove: Global vectors for word representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- Mahdi M Kalayeh, Haroon Idrees, and Mubarak Shah. Nmf-knn: Image annotation using weighted multi-view non-negative matrix factorization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 184–191, 2014.
- Mark A Davenport, Yaniv Plan, Ewout Van Den Berg, and Mary Wootters. 1-bit matrix completion. *Information and Inference:* A Journal of the IMA, 3(3):189–223, 2014.

- Clemens Reimann, Peter Filzmoser, and Robert G Garrett. Factor analysis applied to regional geochemical data: problems and possibilities. *Applied geochemistry*, 17(3):185–206, 2002.
- Prem Gopalan, Jake M. Hofman, and David M. Blei. Scalable recommendation with hierarchical poisson factorization. *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, 2015.
- Havard Rue, Sara Martino, and Nicolas Chopin. Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. *Journal of the Royal Statistical Society, Series B*, 2009.
- Ningning Wu and Jing Zhang. Factor analysis based anomaly detection. *IEEE Systems, Man and Cybernetics SocietyInformation Assurance Workshop*, 2003.
- Lung-fei Lee. Asymptotic bias in simulated maximum likelihood estimation of discrete choice models. *Econometric Theory*, 11: 437–83, 1995.
- Zhenglin Li, Haibei Zhu, Houze Liu, Jintong Song, and Qishuo Cheng. Comprehensive evaluation of mal-api-2019 dataset by machine learning in malware detection. *arXiv preprint arXiv:2403.02232*, 2024.
- Genshiro Kitagawa. Non-gaussian state—space modeling of nonstationary time series. *Journal of the American statistical association*, 82(400):1032–1041, 1987.
- Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. *Advances in neural information processing systems*, 27, 2014.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- Peter D Hoff, Adrian E Raftery, and Mark S Handcock. Latent space approaches to social network analysis. *Journal of the american Statistical association*, 97(460):1090–1098, 2002.
- Stephen J Young and Edward R Scheinerman. Random dot product graph models for social networks. In Algorithms and Models for the Web-Graph: 5th International Workshop, WAW 2007, San Diego, CA, USA, December 11-12, 2007. Proceedings 5, pages 138–149. Springer, 2007.
- Weixuan Xia. The average of a negative-binomial lévy process and a class of lerch distributions. *Communications in Statistics-Theory and Methods*, 49(4):1008–1024, 2020.
- John A Nelder and Daryl Pregibon. An extended quasi-likelihood function. *Biometrika*, 74(2):221–232, 1987.
- Jianqing Fan, Jianhua Guo, and Shurong Zheng. Estimating number of factors by adjusted eigenvalues thresholding. *Journal of the American Statistical Association*, 2020.

- Dimitri P. Bertsekas. Nonlinear programming: 2nd edition. *Athena Scientific*, 1999.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 2011.
- T. Ieleman and G Hinton. Neural networks for machine learning. *Technical report*, 2012.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- Luke Tierney and Joseph B Kadane. Accurate approximations for posterior moments and marginal densities. *Journal of the american statistical association*, 81(393):82–86, 1986.
- Weixuan Xia and Yuyang Zhang. On the absolute-value integral of a brownian motion with drift: Exact and asymptotic formulae. *arXiv preprint arXiv:2312.04172*, 2023.
- Luke Tierney, Robert E Kass, and Joseph B Kadane. Fully exponential laplace approximations to expectations and variances of nonpositive functions. *Journal of the american statistical association*, 84(407):710–716, 1989.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39:1–22, 1977.
- Wolfgang Jank and James Booth. Efficiency of monte carlo em and simulated maximum likelihood in two-stage hierarchical models. *Journal of Computational and Graphical Statistics*, 12:214–29, 2003.
- Brian S Caffo, Wolfgang Jank, and Galin L Jones. Ascent-based monte carlo expectation–maximization. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2): 235–251, 2005.
- Wolfgang Jank. Implementing and diagnosing the stochastic approximation em algorithm. *Journal of Computational and Graphical Statistics*, 15:803–29, 2006.
- Yuan Yao, Lorenzo Rosasco, and Andrea Caponnetto. On early stopping in gradient descent learning. *Constructive Approximation*, 26(2):289–315, 2007.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- Mikko Kivelä, Alex Arenas, Marc Barthelemy, James P Gleeson, Yamir Moreno, and Mason A Porter. Multilayer networks. *Journal of complex networks*, 2(3):203–271, 2014.

- Shangsi Wang, Jesús Arroyo, Joshua T Vogelstein, and Carey E Priebe. Joint embedding of graphs. *IEEE transactions on pattern analysis and machine intelligence*, 43(4):1324–1336, 2019.
- Andrew Jones and Patrick Rubin-Delanchy. The multilayer random dot product graph. *arXiv preprint arXiv:2007.10455*, 2020.
- Benjamin Draves. *Joint spectral embeddings of random dot product graphs*. PhD thesis, Boston University, 2022.
- Runze Tang, Michael Ketcha, Alexandra Badea, Evan D Calabrese, Daniel S Margulies, Joshua T Vogelstein, Carey E Priebe, and Daniel L Sussman. Connectome smoothing via low-rank approximations. *IEEE transactions on medical imaging*, 38(6): 1446–1456, 2018.
- Mark E Dickison, Matteo Magnani, and Luca Rossi. *Multilayer social networks*. Cambridge University Press, 2016.
- Matteo Magnani, Luca Rossi, and Davide Vega. Analysis of multiplex social networks with r. *Journal of Statistical Software*, 98: 1–30, 2021.