
Supplementary material to GmGM: a fast Gaussian graphical model for multi-modal data

Anonymous Author(s)
Affiliation
Address
email

1	Contents	
2	1 Notation	1
3	2 Proofs	3
4	2.1 Permutations	3
5	2.2 Derivation of the probability density function	5
6	2.3 Gradient	6
7	2.4 Maximum Likelihood Estimate for the Eigenvectors	7
8	2.5 Maximum Likelihood Estimate for the Eigenvalues	8
9	3 Dependences	10
10	4 Experiments	10
11	4.1 Synthetic data	10
12	4.2 COIL video	10
13	4.3 EchoNet-Dynamic ECGs	11
14	4.4 Mouse embryo stem cell transcriptomics	11
15	4.5 10x Genomics flash frozen lymph node	12
16	4.6 LifeLines-DEEP metagenomics + metabolomics	12
17	5 Regularization	15
18	1 Notation	

19 In addition to the notation used in the main paper, we also introduce further notation to aid in the
20 proofs. For working with tensors, Kolda and Bader [8] proved to be an invaluable resource; we have
21 borrowed their notation in most cases. The only exception is that we have chosen to denote the
22 ℓ -mode matricization of a tensor \mathcal{T} as $\text{mat}_\ell[\mathcal{T}]$ rather than $\mathcal{T}_{(\ell)}$, to highlight its similarity to $\text{vec}[\mathcal{T}]$
23 and free up the subscript for other purposes.

24 To keep track of lengths of axes, we define the following notation:

- 25 • d_ℓ^γ is the length of axis ℓ
- 26 • $d_{>\ell}^\gamma$ is the product of lengths of all axes after ℓ
- 27 • $d_{<\ell}^\gamma$ is the product of lengths of all axes before ℓ
- 28 • $d_{\setminus\ell}^\gamma$ is the product of lengths of all axes except for ℓ
- 29 • d_\forall^γ is the product of lengths of all axes (i.e. the number of elements in \mathcal{D}^γ)
- 30 • $d_\forall = \sum_\gamma d_\forall^\gamma$ is the total number of elements across all datasets

31 In prior work, d_ℓ has been used to represent the lengths of axes but m_ℓ was used where we write $d_{\setminus\ell}$
 32 (such as in [5]). As prior work also used $\setminus\ell$ to represent leaving out the ℓ th axis in other contexts
 33 (such as in [6]), and the analogous definitions of $d_{>\ell}$ and $d_{<\ell}$ were convenient for use in proofs, we
 34 chose to introduce $d_{\setminus\ell}$ as the variable to represent leave-one-out length products. By representing all
 35 of these related concepts with similar symbols, we hope the maths will be easier to parse.

36 We will let \mathbf{I}_a be the $a \times a$ identity matrix, which allows a concise definition of the Kronecker sum:
 37 $\bigoplus_\ell \Psi_\ell = \sum_\ell \mathbf{I}_{d_{<\ell}} \otimes \Psi_\ell \otimes \mathbf{I}_{d_{>\ell}}$.

38 We make frequent use of the vectorization $\text{vec}[\mathbf{M}]$ of a matrix \mathbf{M} , and more generally of a tensor
 39 $\text{vec}[\mathcal{T}]$. We adopt the rows-first convention of vectorization, such that:

$$\text{vec} \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = [1 \quad 2 \quad 3 \quad 4] \quad (1)$$

40 While columns-first is more common, rows-first is more natural when we adopt the convention that
 41 rows are the first axis of tensor; this is the convention that matricization uses, and matricization
 42 is much more important for our work due to its role in defining the Gram matrices. Note that, for
 43 matrices, a rows-first vectorization of \mathbf{M} is equivalent to a columns-first vectorization of \mathbf{M}^T , so there
 44 is no fundamental difference between the two. For vectorizing a tensor, we proceed by stacking the
 45 rest of the axes in order, such that an element (i_1, \dots, i_K) in \mathcal{T} gets mapped to the element $\sum_\ell i_\ell d_{<\ell}$
 46 in $\text{vec}[\mathcal{T}]$.

47 We define the Gram matrices as $\mathbf{S}_\ell^\gamma = \text{mat}_\ell[\mathcal{D}^\gamma] \text{mat}_\ell[\mathcal{D}^\gamma]^T$. Typically we consider only the one-
 48 sample case but if you have multiple samples, indexed by a subscript i , then the Gram matrix becomes
 49 an average: $\mathbf{S}_\ell^\gamma = \frac{1}{n} \sum_i \text{mat}_\ell[\mathcal{D}_i^\gamma] \text{mat}_\ell[\mathcal{D}_i^\gamma]^T$.

50 An essential concept is that of the "stridewise-blockwise trace", defined as:

$$\text{tr}_b^a[\mathbf{M}] = [\text{tr}[\mathbf{M}(\mathbf{I}_a \otimes \mathbf{J}^{ij} \otimes \mathbf{I}_b)]]_{ij} \quad (2)$$

51 Where \mathbf{J}^{ij} is the matrix of zeros except at (i, j) where it has a 1. It is a generalization of the
 52 blockwise trace considered by Kalaitzis et al. [6], and is related to the $\text{proj}_{\mathcal{K}}$ operation defined by
 53 Greenewald, Zhou, and Hero III [5]. Specifically, $\text{proj}_{\mathcal{K}}[\mathbf{M}]$ is equivalent to $\bigoplus_\ell \text{tr}_{d_{>\ell}}^{d_{<\ell}}[\mathbf{M}]$ up to
 54 an additive diagonal factor (Lemma 33 from Greenewald, Zhou, and Hero III [5]). $\text{proj}_{\mathcal{K}}[\mathbf{M}]$ was
 55 defined to be the matrix that best approximates \mathbf{M} (in terms of the Frobenius norm) while being
 56 Kronecker-sum-decomposable. This matrix is not unique; the choice by Greenewald, Zhou, and
 57 Hero III [5] to include an additive factor was to enforce $\text{tr}[\text{proj}_{\mathcal{K}}[\mathbf{M}]] = 0$. We do not wish to
 58 enforce this constraint as it would be impossible to preserve in the multi-tensor case.

59 The parameter b of the stridewise-blockwise trace partitions the $m \times m$ matrix \mathbf{M} into a block matrix
 60 with $b \times b$ blocks of size $(\frac{m}{b} \times \frac{m}{b})$. The parameter a then partitions these blocks into a "strided"
 61 matrix with $a \times a$ strides containing $\frac{m}{ab} \times \frac{m}{ab}$ blocks. We take the trace of each stride, and the final
 62 matrix is the matrix of these traces. As this is conceptually complicated, we provide an example.

$$\text{tr}_2^2 \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \end{bmatrix} \quad (3)$$

$$= \text{tr}_2^2 \left[\begin{array}{c} \text{tr} \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \\ \text{tr} \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \\ \text{tr} \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \\ \text{tr} \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \\ \text{tr} \begin{bmatrix} 1 & 2 \\ 1 & 2 \end{bmatrix} \end{array} \quad \text{tr} \begin{bmatrix} 3 & 4 \\ 3 & 4 \end{bmatrix} \quad \text{tr} \begin{bmatrix} 5 & 6 \\ 5 & 6 \end{bmatrix} \quad \text{tr} \begin{bmatrix} 7 & 8 \\ 7 & 8 \end{bmatrix} \right] \quad (4)$$

$$= \text{tr}_2^2 \begin{bmatrix} 3 & 7 & 11 & 15 \\ 3 & 7 & 11 & 15 \\ 3 & 7 & 11 & 15 \\ 3 & 7 & 11 & 15 \end{bmatrix} \quad (5)$$

$$= \left[\begin{array}{c} \text{tr} \begin{bmatrix} 3 & 11 \\ 3 & 11 \end{bmatrix} \\ \text{tr} \begin{bmatrix} 3 & 11 \\ 3 & 11 \end{bmatrix} \\ \text{tr} \begin{bmatrix} 3 & 11 \\ 3 & 11 \end{bmatrix} \\ \text{tr} \begin{bmatrix} 3 & 11 \\ 3 & 11 \end{bmatrix} \end{array} \quad \text{tr} \begin{bmatrix} 7 & 15 \\ 7 & 15 \end{bmatrix} \right] \quad (6)$$

$$= \begin{bmatrix} 14 & 22 \\ 14 & 22 \end{bmatrix} \quad (7)$$

63 Notice the construction of the "strides" in Line 6 - the parameter of 2 told us to grab every second
64 element from each row and each column.

65 2 Proofs

66 We will assume that no dataset contains repeated axes (i.e. no single tensor has two axes represented
67 by the same graph), as this greatly affects the derived gradients. Shared axes - two tensors having
68 one or more axes in common - are allowed. The case of shared axes is, after all, the whole point of
69 developing this extension to prior work.

70 2.1 Permutations

71 Note that both $\text{vec}[\text{mat}_1[\mathcal{D}^\gamma]]$ and $\text{vec}[\text{mat}_\ell[\mathcal{D}^\gamma]]$ are row vectors containing the same ele-
72 ments, just in a different order. This means that there is a permutation matrix $\mathbf{P}_{\ell \rightarrow 1}$ such that
73 $\text{vec}[\text{mat}_1[\mathcal{D}^\gamma]]^T \mathbf{P}_{\ell \rightarrow 1} = \text{vec}[\text{mat}_\ell[\mathcal{D}^\gamma]]^T$.

74 **Lemma 1** (Rearrangement lemma). $\mathbf{P}_{\ell \rightarrow 1} (\mathbf{I}_{d_{<\ell}} \otimes \Psi_\ell \otimes \mathbf{I}_{d_{>\ell}}) \mathbf{P}_{\ell \rightarrow 1}^T = \Psi_\ell \otimes \mathbf{I}_{d_{\setminus \ell}}$

75 *Proof.* While vec , mat_ℓ and \otimes are defined as operations on matrices, for the purposes of permuta-
76 tions we can consider them as operations on indices. We can express them as follows:

$$\text{vec} : (i_1, \dots, i_K) \rightarrow \left(\sum_{\ell} i_{\ell} d_{<\ell} \right) \quad (8)$$

$$\text{mat}_{\ell} : (i_1, \dots, i_K) \rightarrow \left(i_{\ell}, \sum_{\ell' < \ell} i_{\ell'} d_{<\ell'} + \sum_{\ell' > \ell} i_{\ell'} \frac{d_{<\ell'}}{d_{\ell}} \right) \quad (9)$$

$$\otimes : ((i_1^1, i_1^2), \dots, (i_K^1, i_K^2)) \rightarrow \left(\sum_{\ell} i_{\ell}^1 d_{<\ell}, \sum_{\ell} i_{\ell}^2 d_{<\ell} \right) \quad (10)$$

77 We'll consider just the rows of \otimes , \otimes_{rows} - although the same argument applies with columns:

$$\otimes_{rows} : (i_1^1, \dots, i_K^1) \rightarrow \left(\sum_{\ell} i_{\ell}^1 d_{<\ell} \right) \quad (11)$$

78 Finally, we'll introduce the permutation operation $\sigma_{\ell \rightarrow 1}$ that will change the order of our Kronecker
79 product:

$$\sigma_{\ell \rightarrow 1} : ((i_1^1, i_1^2), \dots, (i_K^1, i_K^2)) \rightarrow (((i_{\ell}^1, i_{\ell}^2), (i_1^1, i_1^2), \dots, (i_{\ell-1}^1, i_{\ell-1}^2), (i_{\ell+1}^1, i_{\ell+1}^2), \dots, (i_K^1, i_K^2))) \quad (12)$$

80 And again without loss of generality we restrict ourself to $\sigma_{\ell \rightarrow 1}^{rows}$:

$$\sigma_{\ell \rightarrow 1}^{rows} : (i_1^1, \dots, i_K^1) \rightarrow (i_{\ell}^1, i_1^1, \dots, i_{\ell-1}^1, i_{\ell+1}^1, \dots, i_K^1) \quad (13)$$

81 After a Kronecker product our indices are in the form $\sum_{\ell} i_{\ell} d_{<\ell}$, and if we were to reorder it with
82 $\sigma_{\ell \rightarrow 1}$ they would be in the form $i_{\ell} + \sum_{\ell' < \ell} i_{\ell'} d_{<\ell'} d_{\ell} + \sum_{\ell' > \ell} i_{\ell'} d_{<\ell'}$. Likewise, if we had matricized
83 it we would have $\left(i_{\ell}, \sum_{\ell' < \ell} i_{\ell'} d_{<\ell'} + \sum_{\ell' > \ell} i_{\ell'} \frac{d_{<\ell'}}{d_{\ell}} \right)$, which is vectorized to $i_{\ell} + \sum_{\ell' < \ell} i_{\ell'} d_{<\ell'} d_{\ell} +$
84 $\sum_{\ell' > \ell} i_{\ell'} d_{<\ell'}$. These reorderings are the same, and hence the matrix that represents it is $\mathbf{P}_{\ell \rightarrow 1}$.

85 □

86 **2.2 Derivation of the probability density function**

87 Recall that the Kronecker-sum-structured normal distribution for a single tensor is defined as follows:

$$\text{vec} [\mathcal{D}^\gamma] \sim \mathcal{N} \left(\mathbf{0}, \left(\bigoplus_{\ell \in \gamma} \Psi_\ell \right)^{-1} \right) \iff \mathcal{D}^\gamma \sim \mathcal{N}_{KIS} \left(\{\Psi_\ell\}_{\ell \in \gamma} \right) \quad (14)$$

88 The log-likelihood for this distribution is given in [6] for the matrix case and [5] for the general tensor
 89 case. However, neither of these papers provide a derivation. As the full derivation will motivated the
 90 construction of lemmas useful for the proofs of Theorems 1 and 2, we will give it here. First, observe
 91 that the density function is that of a normal distribution.

$$p(\mathcal{D}^\gamma) = \frac{\sqrt{\left| \bigoplus_{\ell \in \gamma} \Psi_\ell \right|}}{(2\pi)^{\frac{d^\gamma}{2}}} e^{-\frac{1}{2} \text{vec}[\mathcal{D}^\gamma]^T (\bigoplus_{\ell} \Psi_\ell) \text{vec}[\mathcal{D}^\gamma]} \quad (15)$$

92 **Lemma 2** (\oplus -vec lemma). $\text{vec} [\mathcal{D}^\gamma]^T (\bigoplus_{\ell} \Psi_\ell) \text{vec} [\mathcal{D}^\gamma] = \sum_{\ell} \text{tr} [\mathbf{S}_\ell^\gamma \Psi_\ell]$

93 *Proof.* This proof relies on the following two properties of vec : $(\mathbf{A} \otimes \mathbf{B}) \text{vec} [\mathbf{C}] = \text{vec} [\mathbf{B}\mathbf{C}^T\mathbf{A}^T]$
 94 and $\text{tr} [\mathbf{A}^T\mathbf{B}] = \text{vec} [\mathbf{A}]^T \text{vec} [\mathbf{B}]$. The \mathbf{C} term picks up a transpose due to our use of the rows-first
 95 vectorization; when using columns-first notation the right hand side becomes $\text{vec} [\mathbf{B}\mathbf{C}\mathbf{A}^T]$.

$$\begin{aligned} \text{vec} [\mathcal{D}^\gamma] \left(\bigoplus_{\ell} \Psi_\ell \right) \text{vec} [\mathcal{D}^\gamma] &= \sum_{\ell} \text{vec} [\mathcal{D}^\gamma]^T (\mathbf{I}_{d_{<\ell}} \otimes \Psi_\ell \otimes \mathbf{I}_{d_{>\ell}}) \text{vec} [\mathcal{D}^\gamma] \quad (\text{Definition of } \bigoplus) \\ &= \sum_{\ell} \text{vec} [\text{mat}_1 [\mathcal{D}^\gamma]]^T (\mathbf{I}_{d_{<\ell}} \otimes \Psi_\ell \otimes \mathbf{I}_{d_{>\ell}}) \text{vec} [\text{mat}_1 [\mathcal{D}^\gamma]] \quad (16) \\ &= \sum_{\ell} \text{vec} [\text{mat}_\ell [\mathcal{D}^\gamma]]^T \mathbf{P}_{\ell \rightarrow 1}^T (\mathbf{I}_{d_{<\ell}} \otimes \Psi_\ell \otimes \mathbf{I}_{d_{>\ell}}) \mathbf{P}_{\ell \rightarrow 1} \text{vec} [\text{mat}_\ell [\mathcal{D}^\gamma]] \quad (17) \end{aligned}$$

$$= \sum_{\ell} \text{vec} [\text{mat}_\ell [\mathcal{D}^\gamma]]^T (\Psi_\ell \otimes \mathbf{I}_{d_{\setminus \ell}}) \text{vec} [\text{mat}_\ell [\mathcal{D}^\gamma]] \quad (\text{Rearrangement Lemma})$$

$$= \sum_{\ell} \text{vec} [\text{mat}_\ell [\mathcal{D}^\gamma]]^T \text{vec} [\text{mat}_\ell [\mathcal{D}^\gamma] \Psi_\ell^T] \quad (18)$$

$$= \sum_{\ell} \text{tr} [\mathbf{S}_\ell^\gamma \Psi_\ell] \quad (19)$$

96 □

97 With this lemma, the probability density function in the single-tensor case can be expressed in the
 98 form:

$$p(\mathcal{D}^\gamma) = \frac{\sqrt{\left| \bigoplus_{\ell \in \gamma} \Psi_\ell \right|}}{(2\pi)^{\frac{d^\gamma}{2}}} e^{-\frac{1}{2} \sum_{\ell} \text{tr} [\mathbf{S}_\ell^\gamma \Psi_\ell]} \quad (20)$$

99 Leading to the probability density function for the multi-tensor case as:

$$p(\{\mathcal{D}^\gamma\}) = \prod_{\gamma} \frac{\sqrt{|\bigoplus_{\ell \in \gamma} \Psi_{\ell}|}}{(2\pi)^{\frac{d_{\Psi}}{2}}} e^{-\frac{1}{2} \sum_{\ell} \text{tr}[\mathbf{S}_{\ell}^{\gamma} \Psi_{\ell}]} \quad (21)$$

$$= \frac{\prod_{\gamma} \sqrt{|\bigoplus_{\ell \in \gamma} \Psi_{\ell}|}}{(2\pi)^{\frac{d_{\Psi}}{2}}} e^{-\frac{1}{2} \sum_{\gamma} \sum_{\ell} \text{tr}[\mathbf{S}_{\ell}^{\gamma} \Psi_{\ell}]} \quad (22)$$

$$= \frac{\prod_{\gamma} \sqrt{|\bigoplus_{\ell \in \gamma} \Psi_{\ell}|}}{(2\pi)^{\frac{d_{\Psi}}{2}}} e^{-\frac{1}{2} \sum_{\ell} \text{tr}[\mathbf{S}_{\ell} \Psi_{\ell}]} \quad (23)$$

100 The negative log-likelihood is thus:

$$\text{NLL}(\{\mathcal{D}^\gamma\}) = \frac{d_{\Psi}}{2} \log(2\pi) + \frac{1}{2} \sum_{\ell} \text{tr}[\mathbf{S}_{\ell} \Psi_{\ell}] - \frac{1}{2} \sum_{\gamma} \log \left| \bigoplus_{\ell \in \gamma} \Psi_{\ell} \right| \quad (24)$$

101 2.3 Gradient

102 The derivation of the gradient of the negative log-likelihood is essentially the same as the derivation
 103 given by Kalaitzis et al. [6] for the original Bi-Graphical Lasso. Our derivation is complicated by the
 104 fact that we are considering general tensors rather than matrices. We'll let sym be the symmetricizing
 105 operator that must be applied as we are taking the derivative with respect to a symmetric matrix:
 106 $\text{sym}[\mathbf{M}] = \mathbf{K} \circ \mathbf{M}$, where \mathbf{K} is a matrix with 1s on the diagonal and 2s everywhere else. We'll also
 107 define \mathbf{J}^{ij} to be the matrix of zeros except for a 1 at position (i, j) .

$$\frac{d}{d\Psi_\ell} \text{NLL}(\{\mathcal{D}^\gamma\}) = \frac{1}{2} \text{sym}[\mathbf{S}_\ell] - \frac{1}{2} \sum_\gamma \frac{d}{d\Psi_\ell} \log \left| \bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right| \quad (25)$$

$$= \frac{1}{2} \text{sym}[\mathbf{S}_\ell] - \frac{1}{2} \sum_\gamma \text{tr} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \frac{d}{d\psi_\ell^{ij}} \bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right]_{ij} \quad (26)$$

$$= \frac{1}{2} \text{sym}[\mathbf{S}_\ell] - \frac{1}{2} \sum_\gamma \text{tr} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \left(\mathbf{I}_{d_{<\ell}} \otimes \frac{d}{d\psi_\ell^{ij}} \Psi_\ell \otimes \mathbf{I}_{d_{>\ell}} \right) \right]_{ij} \quad (27)$$

$$= \frac{1}{2} \text{sym}[\mathbf{S}_\ell] - \frac{1}{2} \sum_\gamma \text{tr} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \left(\mathbf{I}_{d_{<\ell}} \otimes (\mathbf{J}^{ij} + \mathbf{J}^{ji} - \delta_{ij} \mathbf{J}^{ij}) \otimes \mathbf{I}_{d_{>\ell}} \right) \right]_{ij} \quad (28)$$

$$= \frac{1}{2} \text{sym}[\mathbf{S}_\ell] - \frac{1}{2} \sum_\gamma \left[(2 - \delta_{ij}) \text{tr} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \left(\mathbf{I}_{d_{<\ell}} \otimes \mathbf{J}^{ij} \otimes \mathbf{I}_{d_{>\ell}} \right) \right] \right]_{ij} \quad (29)$$

$$= \frac{1}{2} \text{sym}[\mathbf{S}_\ell] - \frac{1}{2} \sum_\gamma (2\mathbf{J} - \mathbf{I}) \circ \text{tr} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \left(\mathbf{I}_{d_{<\ell}} \otimes \mathbf{J}^{ij} \otimes \mathbf{I}_{d_{>\ell}} \right) \right]_{ij} \quad (30)$$

$$= \frac{1}{2} \text{sym}[\mathbf{S}_\ell] - \frac{1}{2} \sum_\gamma \text{sym} \left[\text{tr}_{d_{>\ell}}^{d_{<\ell}} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \right] \right] \quad (31)$$

108 The MLE occurs when this gradient is zero, i.e. when the following equation is satisfied:

$$\mathbf{S}_\ell = \sum_\gamma \text{tr}_{d_{>\ell}}^{d_{<\ell}} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \right] \quad (32)$$

109 In other words, our effective Gram matrices are the best Kronecker-sum decomposition of the
 110 covariance matrix of the maximum likelihood estimate. Unfortunately, Kronecker-sum decomposition
 111 does not interact well with matrix inverses, so this does not directly yield an analytic solution. It does,
 112 however, yield a solution for the eigenvectors.

113 2.4 Maximum Likelihood Estimate for the Eigenvectors

114 We first produce two lemmas to aid in the derivation.

115 **Lemma 3** (Cyclic property of the stridewise-blockwise trace). *For any matrices \mathbf{M} , $\mathbf{A}_{a \times a}$, $\mathbf{B}_{b \times b}$,
 116 we have that $\text{tr}_b^a [(\mathbf{A} \otimes \mathbf{I} \otimes \mathbf{B}) \mathbf{M}] = \text{tr}_b^a [\mathbf{M} (\mathbf{A} \otimes \mathbf{I} \otimes \mathbf{B})]$*

117 *Proof.* This follows directly from the cyclic property of the (normal) trace operator and the definition
 118 of the stridewise-blockwise trace. \square

119 **Lemma 4** (Conjugacy extraction of the stridewise-blockwise trace). *For any matrices \mathbf{M} and \mathbf{V} , we
 120 have that $\text{tr}_b^a [(\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b) \mathbf{M} (\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b)^T] = \mathbf{V} \text{tr}_b^a [\mathbf{M}] \mathbf{V}^T$.*

Proof.

$$\text{tr}_b^a \left[(\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b) \mathbf{M} (\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b)^T \right] = \left[\text{tr} \left[(\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b) \mathbf{M} (\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b)^T (\mathbf{I}_a \otimes \mathbf{J}^{ij} \otimes \mathbf{I}_b) \right] \right]_{ij} \quad (\text{Definition of } \text{tr}_b^a)$$

121 Thanks to the Rearrangement Lemma, we can get this just in terms of the standard blockwise trace,
 122 for which there exists a convenient lemma from Dahl et al. [3] that does the heavy lifting for us.
 123 Unfortunately, this requires inserting permutation matrices into every nook and cranny.

$$= \left[\text{tr} \left[\mathbf{P} (\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b) \mathbf{P}^T \mathbf{P} \mathbf{M} \mathbf{P}^T \mathbf{P} (\mathbf{I}_a \otimes \mathbf{V} \otimes \mathbf{I}_b)^T \mathbf{P}^T \mathbf{P} (\mathbf{I}_a \otimes \mathbf{J}^{ij} \otimes \mathbf{I}_b) \mathbf{P}^T \right] \right]_{ij} \quad (33)$$

$$= \left[\text{tr} \left[(\mathbf{V} \otimes \mathbf{I}_{ab})^T \mathbf{P} \mathbf{M} \mathbf{P}^T (\mathbf{V} \otimes \mathbf{I}_{ab}) (\mathbf{J}^{ij} \otimes \mathbf{I}_{ab}) \right] \right]_{ij} \quad (34)$$

$$= \text{tr}_{ab} \left[(\mathbf{V} \otimes \mathbf{I}_{ab}) \mathbf{P} \mathbf{M} \mathbf{P}^T (\mathbf{V} \otimes \mathbf{I}_{ab})^T \right] \quad (\text{Definition of } \text{tr}_{ab})$$

$$= \mathbf{V} \text{tr}_{ab} \left[\mathbf{P} \mathbf{M} \mathbf{P}^T \right] \mathbf{V}^T \quad (\text{Lemma 2 of Dahl et al. [3]})$$

124 We then can see analogously that $\text{tr}_{ab} \left[\mathbf{P} \mathbf{M} \mathbf{P}^T \right] = \text{tr}_b^a \left[\mathbf{M} \right]$, completing the proof.

125

□

126 **Theorem 1.** Let $\mathbf{V}_\ell \mathbf{e}_\ell \mathbf{V}_\ell^T$ be the eigendecomposition of \mathbf{S}_ℓ . Then \mathbf{V}_ℓ are the eigenvectors of the
 127 maximum likelihood estimate of Ψ_ℓ .

Proof.

$$\mathbf{S}_\ell = \sum_\gamma \text{tr}_{d>\ell}^{d<\ell} \left[\left(\bigoplus_{\ell' \in \gamma} \Psi_{\ell'} \right)^{-1} \right] \quad (35)$$

$$= \sum_\gamma \text{tr}_{d>\ell}^{d<\ell} \left[\left(\bigoplus_{\ell' \in \gamma} \mathbf{V}_{\ell'} \Lambda_{\ell'} \mathbf{V}_{\ell'}^T \right)^{-1} \right] \quad (36)$$

$$= \sum_\gamma \text{tr}_{d>\ell}^{d<\ell} \left[\left(\bigotimes_{\ell'} \mathbf{V}_{\ell'} \right) \left(\bigoplus_{\ell' \in \gamma} \Lambda_{\ell'} \right)^{-1} \left(\bigotimes_{\ell'} \mathbf{V}_{\ell'} \right)^T \right] \quad (37)$$

$$= \sum_\gamma \text{tr}_{d>\ell}^{d<\ell} \left[\left(\mathbf{I}_{d<\ell} \otimes \mathbf{V}_\ell \otimes \mathbf{I}_{d>\ell} \right) \left(\bigoplus_{\ell' \in \gamma} \Lambda_{\ell'} \right)^{-1} \left(\mathbf{I}_{d<\ell} \otimes \mathbf{V}_\ell \otimes \mathbf{I}_{d>\ell} \right)^T \right] \quad (\text{Cyclic Property})$$

$$= \sum_\gamma \mathbf{V} \text{tr}_{d>\ell}^{d<\ell} \left[\left(\bigoplus_{\ell' \in \gamma} \Lambda_{\ell'} \right)^{-1} \right] \mathbf{V}^T \quad (\text{Conjugacy Extraction})$$

$$= \mathbf{V} \left(\sum_\gamma \text{tr}_{d>\ell}^{d<\ell} \left[\left(\bigoplus_{\ell' \in \gamma} \Lambda_{\ell'} \right)^{-1} \right] \right) \mathbf{V}^T \quad (38)$$

128 We conclude the proof by observing that the central matrix is diagonal, and thus the right hand side
 129 constitutes an eigendecomposition of \mathbf{S}_ℓ . Thus \mathbf{S}_ℓ and Ψ_ℓ share eigenvectors. □

130 2.5 Maximum Likelihood Estimate for the Eigenvalues

131 In the previous section, we derived the eigenvectors of the maximum likelihood estimate. While
 132 interesting (they correspond to the principal components of our data), we need the eigenvalues

133 to reconstruct Ψ_ℓ . Our strategy for this is to transform our data such that the precision matrices
 134 are diagonal, and estimate these diagonals. This transformation is stated in terms of the Tucker
 135 operator ($\llbracket \mathcal{D}^\gamma; \{\mathbf{V}_\ell^T\}_{\ell \in \gamma} \rrbracket$). In the case where \mathcal{D} is a matrix, we have that $\llbracket \mathbf{D}; \mathbf{V}_{rows}^T, \mathbf{V}_{cols}^T \rrbracket =$
 136 $\mathbf{V}_{rows} \mathbf{D} \mathbf{V}_{cols}^T$. While the definition of the Tucker operator can be given in terms of “n-mode prod-
 137 ucts”[8], we will only use the following property relating the Tucker operator to matricization Kolda
 138 [7]:

$$\begin{aligned} \mathcal{Y} &= \llbracket \mathcal{X}; \{\mathbf{M}_\ell\} \rrbracket \\ \implies \text{mat}_\ell[\mathcal{Y}] &= \mathbf{M}_\ell \text{mat}_\ell[\mathcal{X}] (\mathbf{M}_K \otimes \dots \otimes \mathbf{M}_{\ell+1} \otimes \mathbf{M}_{\ell-1} \otimes \dots \otimes \mathbf{M}_1)^T \quad (\text{Kolda [7]}) \end{aligned}$$

139 The Tucker operator is an important concept for our calculation of the eigenvalues, but it is only the
 140 existence of such an operator that is important for our work; we never need to calculate it.

141 **Lemma 5** (Eigendecompositions of the Kronecker-sum-structured normal distribution). *Suppose*
 142 $\{\mathcal{D}^\gamma\} \sim \mathcal{N}_{KS}(\{\Psi_\ell\})$. *Then $\{\llbracket \mathcal{D}^\gamma; \{\mathbf{V}_\ell^T\} \rrbracket\} \sim \mathcal{N}_{KS}(\{\Lambda_\ell\})$ and the effective Gram matrices of*
 143 *this distribution are given by the eigenvalues \mathbf{e}_ℓ of the effective Gram matrices \mathbf{S}_ℓ of the original*
 144 *distribution.*

145 *Proof.* We will prove this by showing that the probability density function is that of a Kronecker-
 146 sum-structured normal distribution with the given parameters.

147 In the first part of the proof, we will massage the probability density function into a convenient form -
 148 this does not depend on the Tucker decomposition, and holds for our original dataset as well.

$$p(\llbracket \mathcal{D}^\gamma; \{\mathbf{V}_\ell^T\}_{\ell \in \gamma} \rrbracket) = p(\{\mathcal{D}^\gamma\}) \quad (39)$$

$$= \frac{\prod_\gamma \sqrt{|\bigoplus_{\ell \in \mathcal{D}^\gamma} \Psi_\ell|}}{(2\pi)^{\frac{d_Y}{2}}} e^{-\frac{1}{2} \sum_\ell \text{tr}[\Psi_\ell \mathbf{S}_\ell]} \quad (40)$$

$$= \frac{\prod_\gamma \sqrt{|\bigoplus_{\ell \in \mathcal{D}^\gamma} \Lambda_\ell|}}{(2\pi)^{\frac{d_Y}{2}}} e^{-\frac{1}{2} \sum_\ell \text{tr}[\mathbf{V}_\ell \Lambda_\ell \mathbf{V}_\ell^T \mathbf{S}_\ell]} \quad (41)$$

$$= \frac{\prod_\gamma \sqrt{|\bigoplus_{\ell \in \mathcal{D}^\gamma} \Lambda_\ell|}}{(2\pi)^{\frac{d_Y}{2}}} e^{-\frac{1}{2} \sum_\ell \text{tr}[\Lambda_\ell \mathbf{V}_\ell^T \mathbf{S}_\ell \mathbf{V}_\ell]} \quad (42)$$

$$= \frac{\prod_\gamma \sqrt{|\bigoplus_{\ell \in \mathcal{D}^\gamma} \Lambda_\ell|}}{(2\pi)^{\frac{d_Y}{2}}} e^{-\frac{1}{2} \sum_\ell \text{tr}[\Lambda_\ell \mathbf{e}_\ell]} \quad (43)$$

149 To complete the proof, we must show that \mathbf{e}_ℓ are the effective Gram matrices for $\llbracket \mathcal{D}_j; \{\mathbf{V}_\ell^T\}_{\ell \in \mathcal{D}_j} \rrbracket$.

150 For brevity, let $\mathbf{V}_{\setminus \ell} = (\mathbf{V}_K \otimes \dots \otimes \mathbf{V}_{\ell+1} \otimes \mathbf{V}_{\ell-1} \otimes \dots \otimes \mathbf{V}_1)$.

$$\mathbf{e}_\ell = \mathbf{V}_\ell^T \mathbf{S}_\ell \mathbf{V}_\ell \quad (44)$$

$$= \sum_{\ell' \in \gamma} \frac{1}{n} \sum_i^n \mathbf{V}_\ell^T \text{mat}_\ell[\mathcal{D}_i^\gamma] \text{mat}_\ell[\mathcal{D}_i^\gamma]^T \mathbf{V}_\ell \quad (\text{Definition of } \mathbf{S}_\ell)$$

$$= \sum_{\ell' \in \gamma} \frac{1}{n} \sum_i^n \mathbf{V}_\ell^T \text{mat}_\ell[\mathcal{D}_i^\gamma] \mathbf{V}_{\setminus \ell}^T \mathbf{V}_{\setminus \ell} \text{mat}_\ell[\mathcal{D}_i^\gamma]^T \mathbf{V}_\ell \quad (45)$$

$$= \sum_{\ell' \in \gamma} \frac{1}{n} \sum_i^n \text{mat}_\ell \left[\llbracket \mathcal{D}_j; \{\mathbf{V}_\ell^T\}_{\ell \in \mathcal{D}_j} \rrbracket \right] \text{mat}_\ell \left[\llbracket \mathcal{D}_j; \{\mathbf{V}_\ell^T\}_{\ell \in \mathcal{D}_j} \rrbracket \right]^T \quad (46)$$

151 This completes the proof.

152

□

153 Since this transformed data is still normally distributed with Kronecker-sum structure, we can use the
154 previously derived gradient (Line 32):

$$\frac{d}{d\mathbf{\Lambda}_\ell} \text{NLL}(\{\mathcal{D}^\gamma\}) = \mathbf{e}_\ell - \sum_{\gamma} \text{tr}_{d_{>\ell}}^{d_{<\ell}} \left[\left(\bigoplus_{\ell' \in \gamma} \mathbf{\Lambda}_{\ell'} \right)^{-1} \right] \quad (47)$$

155 This yields Theorem 2:

156 **Theorem 2.** *Let $\{\mathbf{G}_\ell^\gamma\}$ be matrices such that the expression $\bigoplus_{\ell \in \gamma} \mathbf{G}_\ell^\gamma$ is the best Frobenius-norm
157 approximation of $\left(\bigoplus_{\ell \in \gamma} \mathbf{\Lambda}_\ell^t \right)^{-1}$. Then, for a learning rate μ_t , gradient descent can be performed
158 with the update equation $\mathbf{\Lambda}_\ell^{t+1} = \mathbf{\Lambda}_\ell^t - \mu_t \left[\mathbf{e}_\ell - \sum_{\gamma | \ell \in \gamma} \mathbf{G}_\ell^\gamma \right]$. As $\mathbf{\Psi}_\ell$ is positive definite, μ_t must
159 be chosen to prevent $\mathbf{\Lambda}_\ell^t$ from becoming negative.*

160 This is convenient because we have reduced our optimization task from one with $\sum_\ell d_\ell^2$ parameters
161 to one with $\sum_\ell d_\ell$ parameters.

162 3 Dependences

163 All tests and figures were generated on a Linux (Ubuntu 20.04) with an Intel Core i7 chip and 8GB of
164 RAM. Along with our code, we provide an environment file (environment.yml) that contains full
165 details of all the dependencies used. In our GitHub repository (<https://github.com/NeurIPS-GmGM-Paper/GmGM>), we give precise and simple instructions on how to create a conda environment with
166 the same setup as ours. Most of the packages used were specific to the experiments we ran. The
167 dependencies necessary for our algorithm were Python 3.9 and NumPy 1.23.5.
168

169 4 Experiments

170 4.1 Synthetic data

171 We generated random graphs by modelling each edge’s probability of existing as being drawn from
172 independent Bernoulli distributions. When estimating the runtimes, we ran all models five times and
173 averaged out the results. When creating precision-recall curves, we averaged the results of fifty runs
174 of the models. Due to space reasons, we omitted the precision-recall curves for the tensor-variate
175 case in our main paper, so we provide this here in Figure 1.

176 4.2 COIL video

177 We downloaded the processed COIL-20[10] dataset, and tested our model on it. We wanted to see
178 if our model could understand the structure of a video, which we expected to consist of two linear
179 graphs (for the rows and columns, i.e. each row is connected only to its neighbor rows) and a circular
180 graph (for the frames, because the video is of a 360° rotation). To generate these graphs, we ran
181 our algorithm on the duck video from the dataset, and then greedily kept the largest edge from each
182 vertex such that vertices in the final graph had at most two edges. If we shuffled our data (shuffle
183 rows, columns, and frames) and try to reconstruct it with these graphs, we get mixed results (Figure
184 2).

185 We can put a numeric value to the reconstruction, by measuring the percentage of the time that
186 our reconstructed edges connect two adjacent rows/columns/frames. We get an accuracy of 80%
187 for the rows, 91% for the columns, and 99% for the frames. This hints that it is quite good at
188 reconstructing frames of videos, but rows and columns are a harder task. This could be due to the
189 specific characteristics of this video, in which there are a lot of rows that spend most of their time
190 being mostly blue.

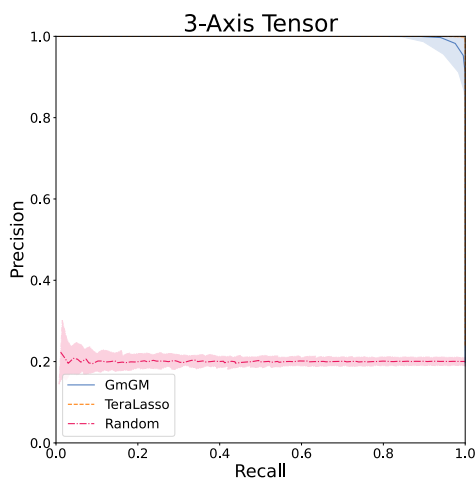


Figure 1: PR curves for the graphs generated from a 3-axis tensor. TeraLasso does almost perfectly; it can be hard to see as it is hugging the top right corner.

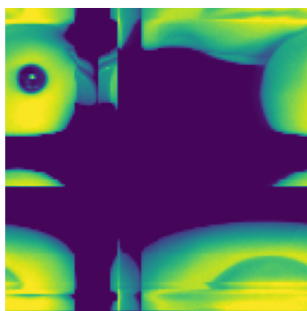


Figure 2: A reconstruction of the COIL-20 duck video after shuffling the rows, columns, and frames, using GmGM. While portions of the duck are well-reconstructed, it is clearly imperfect. Notably, the duck kisses itself.

191 4.3 EchoNet-Dynamic ECGs

192 We downloaded all of the EchoNet-Dynamic[11] data. This dataset did not have heartbeats labeled,
 193 so we picked a few videos at random and labeled them ourselves as a proof of concept. Specifically,
 194 we labeled every frame in which the mitral valve opened. Our goal was to see if the graphs produced
 195 by our algorithm could predict this opening. Table 1 contains the videos we picked, the labels we
 196 gave, and the labels we predicted.

197 Mitral valve predictions were done by taking GmGM’s output frames graph in precision matrix form,
 198 and measuring the mass along the diagonals. We treated this as a time series (since each diagonal
 199 corresponds to an increasing time offset from all frames simultaneously). We applied gaussian blur
 200 and then a continuous wavelet transform peak detection algorithm[4] to find which diagonals had
 201 the most mass (Figure 3). These represent the offsets corresponding with a heartbeat. Given the first
 202 mitral valve opening and these offsets, we predict the remaining openings.

203 4.4 Mouse embryo stem cell transcriptomics

204 We used the mouse embryo stem cell dataset E-MTAB-2805[2]. This dataset had already been labeled
 205 by what stage of the cell cycle each cell was in. The data was log-transformed, and we restricted the

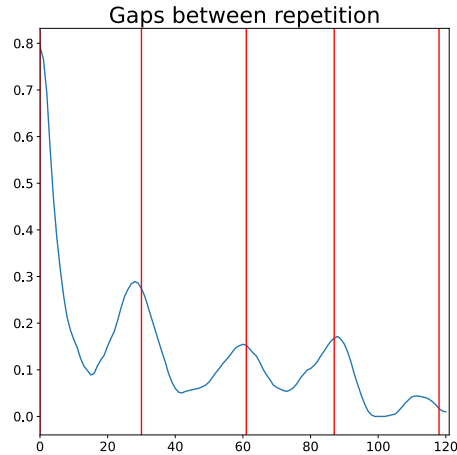


Figure 3: An example heartbeat offset detection, from EchoNet-Dynamic video 0XFE6E32991136338. The blue curve represents our Gaussian-blurred diagonal mass (if $x=10$, it represents the blurred mass of the 10th diagonal to the right of the main diagonal). The red lines represent the predicted peaks via a continuous wavelet transform peak detection algorithm. These represent offsets from the first mitral valve opening. For this video, the mitral valve opened on frame 17 and our first offset was on the 30th diagonal. Hence, we would predict the second mitral valve opening to occur at frame 47 (which, in this case, was correct).

206 gene set to the same mitosis-related genes used for Li et al. [9]’s analysis of this same dataset. We
 207 kept the top 100 edges in our output graphs for each vertex, and set the rest to zero.

208 4.5 10x Genomics flash frozen lymph node

209 For this experiment, we looked at a single-cell RNA-sequencing+ATAC-sequencing dataset from 10x
 210 Genomics[1]. We wanted to know whether clusters in UMAP-space make sense when viewed on
 211 GmGM’s predicted graphs, whether clusters on the graphs made sense in UMAP-space, and whether
 212 these clusters had any meaning. Before performing the experiment, we removed cells whose library
 213 size was three median absolute deviations from the median, and similarly removed genes and peaks
 214 if the the total amount of times they were expressed was three median absolute deviations from the
 215 median. In our output graphs, we kept the top 5 edges per vertex.

216 From Figures 4 and 5, we can see that the clusters indeed seem to make sense in both UMAP-space
 217 and on the GmGM graph, as they all form coherent regions in both spaces.

218 To validate that these clusters are meaningful, we performed a GO term enrichment analysis; the full
 219 results of this analysis are saved on our GitHub repository, but we summarize them here.

220 Clusters 3 and 7 are clearly distinct in both spaces, and this is reflected in their GO terms. Cluster
 221 3 was strongly associated with the CCKR signalling map and apoptosis, which none of the other
 222 clusters were. Cluster 7 was the most distinct, associated with the integrin signalling pathway, blood
 223 coagulation, and insulin. The other clusters all related to B and T cell-specific pathways. GmGM
 224 always grouped clusters 4 and 6 together, whereas UMAP would sometimes prefer to group cluster 6
 225 with the rest of the clusters (compare Figures 5a and 6).

226

227 4.6 LifeLines-DEEP metagenomics + metabolomics

228 We used the LifeLines-DEEP metagenomics and metabolomics datasets[13]. We did not do any pre-
 229 processing to the metabolomics, and we used the already pre-processed version of the metagenomics
 230 data from Prost, Gazut, and Bruls [12]. We kept only patients that appeared in both datasets, and

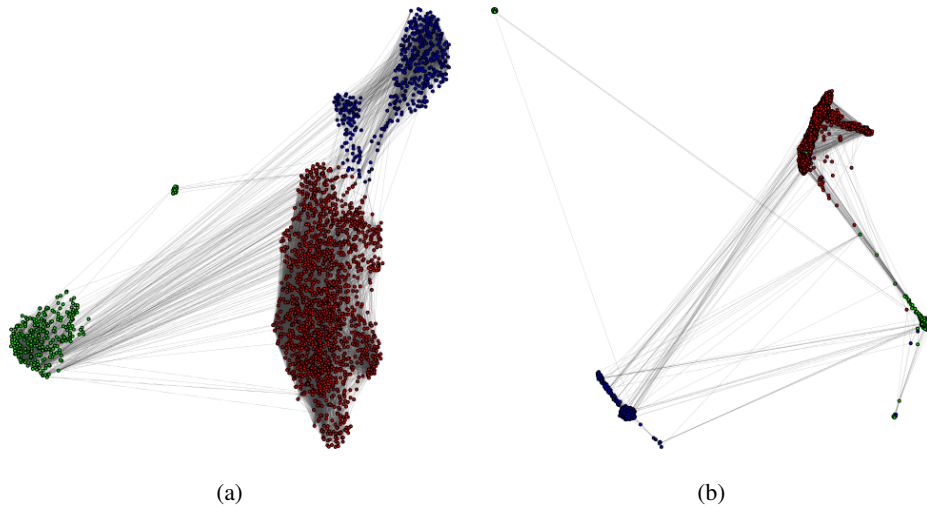


Figure 4: (a) UMAP of the cells in the 10x Genomics dataset. Colored by kmeans ($k=3$). (b) GmGM's predicted graph over those cells, colored using the same clusters as on UMAP and plotted using igraph without reference to the outputs of UMAP.

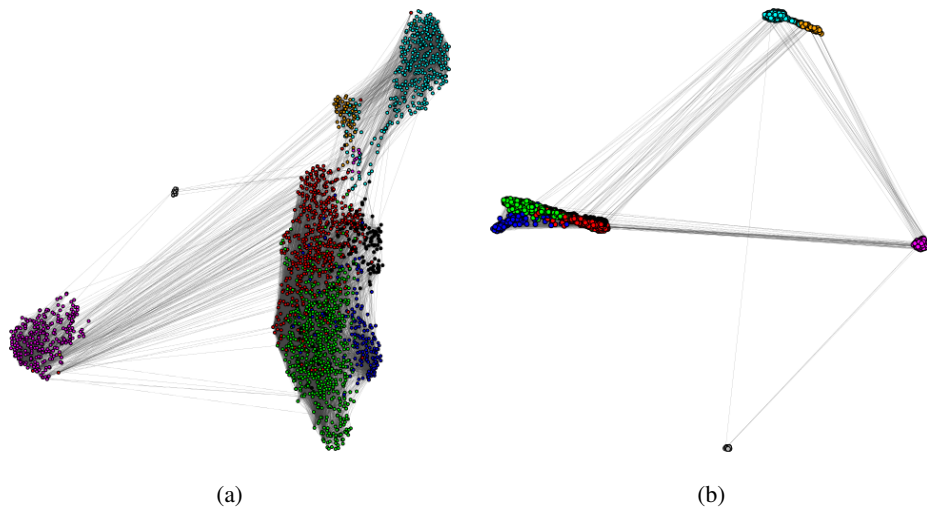


Figure 5: (a) UMAP of the cells in the 10x Genomics dataset. Colored using same clusters as GmGM. (b) GmGM's predicted graph over those cells, colored using Louvain clustering.


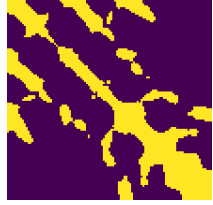


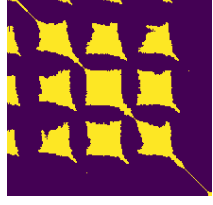
Video ID	Label	Predicted	Precision Matrix
0XFE6E32991136338	[17, 47, 77, 106]	[17, 47, 78, 104]	<p>Echocardiogram Frames</p> 
0XF072F7A9791B060	[24, 56, 100]	[24, 59, 90]	<p>Echocardiogram Frames</p> 
0XF70A3F712E03D87	[22, 66, 110]	[22, 67, 111]	<p>Echocardiogram Frames</p> 
0XF60BBEC9C303C98	[19, 67, 114, 162]	[19, 66, 115, 162]	<p>Echocardiogram Frames</p> 
0XF46CF63A2A1FA90	[25, 79, 134, 188]	[25, 80, 133, 184]	<p>Echocardiogram Frames</p> 

Table 1: Mitral valve labellings and precision matrices for the EchoNet-Dynamic dataset. The precision matrices, for the most part, seem to have clear off-diagonal structures, as expected, and the mitral valve prediction is generally quite good; it is only significantly off for the last opening in 0XF072F7A9791B060.

²³¹ log-transformed the data. We compared our model’s results to the model given by Prost, Gazut, and
²³² Brls [12] in the main paper.

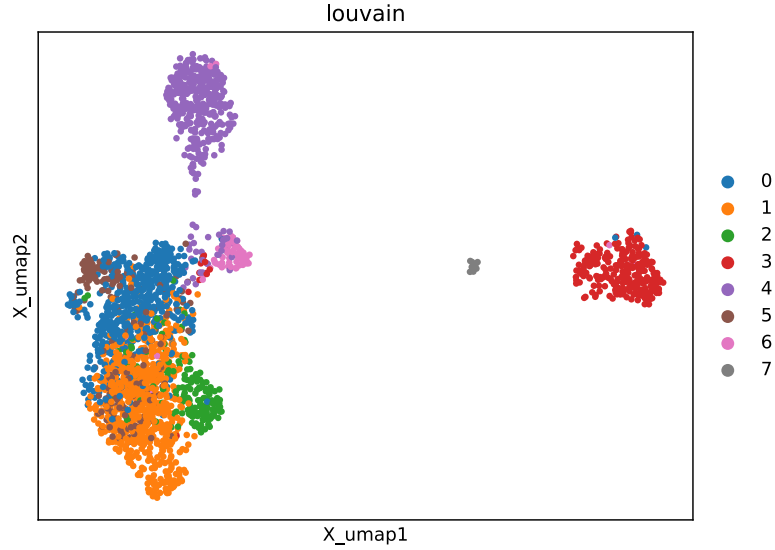


Figure 6: Another UMAP plot showing the same concept as Figure 5a, with the clusters labeled

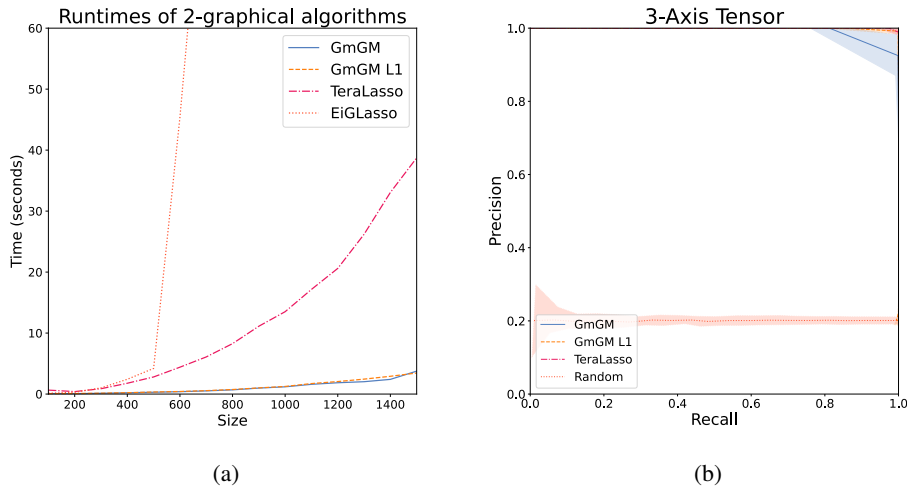


Figure 7: (a) Runtimes of our algorithm and prior work on matrix-variate data. Our regularized algorithm is denoted “GmGM L1”, and takes about the same time as the unregularized “GmGM”. (b) Precision-recall curves for tensor-variate data. TeraLasso and our regularized “GmGM L1” perform almost perfectly.

233 5 Regularization

234 As remarked in the main paper, our algorithm by default includes no regularization. This is because
 235 our algorithm leverages the fact that we have a closed-form expression for the eigenvectors of the
 236 maximum likelihood estimate to avoid costly eigendecompositions every iteration. We do not have a
 237 closed-form expression for the eigenvectors in the regularized case.

238 Nevertheless, we can add regularization to the eigenvalue estimation by performing an eigenrecom-
 239 position and regularizing that. Eigenrecomposition requires a matrix multiplication, which is quite
 240 costly compared to the cost of an unregularized iteration - both in practice, and asymptotically in
 241 the matrix-variate case (matrix multiplication is $O(\sum_{\ell} d_{\ell}^3)$ whereas an unregularized iteration is

242 $O(\prod_{\ell} d_{\ell})$). Thus, to regularize we first let our algorithm converge to the MLE before considering the
 243 penalty term. This allows us to avoid a major increase in runtime; our regularized algorithm runs in
 244 roughly the same time as the unregularized one (Figure 7a).

245 It is important to note that this estimator is slightly different than the standard Lasso estimator, as the
 246 standard estimator would minimize $\|\Psi_{\ell}\|_1$ and our estimator minimizes $\|\hat{\mathbf{V}}_{\ell}\mathbf{\Lambda}_{\ell}\hat{\mathbf{V}}_{\ell}^T\|_1$ (where only
 247 the eigenvalues $\mathbf{\Lambda}_{\ell}$ are free to vary). It can be derived as follows:

$$\frac{\partial}{\partial\lambda_i}\|\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T\|_1 = \frac{\partial}{\partial\lambda_i}\left\|\sum_j\lambda_jv_{ja}v_{bj}\right\|_1 \quad (48)$$

$$= \left[\frac{\partial}{\partial\lambda_i}\left|\sum_j\lambda_jv_{ja}v_{bj}\right|\right]_{ab} \quad (49)$$

$$= \left[\frac{\partial}{\partial\lambda_i}\text{sign}\left[\sum_j\lambda_jv_{ja}v_{bj}\right]v_{ia}v_{bi}\right]_{ab} \quad (50)$$

$$= [\text{sign}[\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T]_{ab}v_{ia}v_{bi}]_{ab} \quad (51)$$

$$= \mathbf{v}_i^T\text{sign}[\mathbf{V}\mathbf{\Lambda}\mathbf{V}^T]\mathbf{v}_i \quad (52)$$

248 Despite this difference, it performs comparably to prior work. We show in Figure 7b the precision-
 249 recall curves for the 3-axis case, and observe that it performs almost perfectly. This is notable as it
 250 was the case that the unregularized algorithm performed worse than prior work.

251 References

- 252 [1] 10x Genomics. *Flash-Frozen Lymph Node with B Cell Lymphoma (14k sorted nuclei)*. en. May
 253 2021. URL: [https://www.10xgenomics.com/resources/datasets/fresh-frozen-](https://www.10xgenomics.com/resources/datasets/fresh-frozen-lymph-node-with-b-cell-lymphoma-14-k-sorted-nuclei-1-standard-2-0-0)
 254 [lymph-node-with-b-cell-lymphoma-14-k-sorted-nuclei-1-standard-2-0-0](https://www.10xgenomics.com/resources/datasets/fresh-frozen-lymph-node-with-b-cell-lymphoma-14-k-sorted-nuclei-1-standard-2-0-0)
 255 (visited on 05/10/2023).
- 256 [2] Florian Buettner et al. “Computational analysis of cell-to-cell heterogeneity in single-cell
 257 RNA-sequencing data reveals hidden subpopulations of cells”. en. In: *Nature Biotechnology*
 258 33.2 (Feb. 2015). Number: 2 Publisher: Nature Publishing Group, pp. 155–160. ISSN: 1546-
 259 1696. DOI: 10.1038/nbt.3102. URL: <https://www.nature.com/articles/nbt.3102>
 260 (visited on 05/10/2023).
- 261 [3] Andy Dahl et al. *Network inference in matrix-variate Gaussian models with non-independent*
 262 *noise*. arXiv:1312.1622 [stat]. Dec. 2013. DOI: 10.48550/arXiv.1312.1622. URL: [http:](http://arxiv.org/abs/1312.1622)
 263 [://arxiv.org/abs/1312.1622](http://arxiv.org/abs/1312.1622) (visited on 02/27/2023).
- 264 [4] Pan Du, Warren A. Kibbe, and Simon M. Lin. “Improved peak detection in mass spectrum by
 265 incorporating continuous wavelet transform-based pattern matching”. In: *Bioinformatics* 22.17
 266 (Sept. 2006), pp. 2059–2065. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/bt1355.
 267 URL: <https://doi.org/10.1093/bioinformatics/bt1355> (visited on 05/11/2023).
- 268 [5] Kristjan Greenewald, Shuheng Zhou, and Alfred Hero III. *Tensor Graphical Lasso (TeraLasso)*.
 269 arXiv:1705.03983 [stat]. Sept. 2019. URL: <http://arxiv.org/abs/1705.03983> (visited
 270 on 02/24/2023).
- 271 [6] Alfredo Kalaitzis et al. “The Bigraphical Lasso”. en. In: *Proceedings of the 30th International*
 272 *Conference on Machine Learning*. ISSN: 1938-7228. PMLR, May 2013, pp. 1229–1237. URL:
 273 <https://proceedings.mlr.press/v28/kalaitzis13.html> (visited on 02/24/2023).
- 274 [7] Tamara Kolda. *Multilinear operators for higher-order decompositions*. en. Tech. rep.
 275 SAND2006-2081, 923081. Apr. 2006, SAND2006–2081, 923081. DOI: 10.2172/923081.
 276 URL: <https://www.osti.gov/servlets/purl/923081/> (visited on 05/05/2023).
- 277 [8] Tamara G. Kolda and Brett W. Bader. “Tensor Decompositions and Applications”. en. In:
 278 *SIAM Review* 51.3 (Aug. 2009), pp. 455–500. ISSN: 0036-1445, 1095-7200. DOI: 10.1137/
 279 07070111X. URL: <http://epubs.siam.org/doi/10.1137/07070111X> (visited on
 280 02/26/2023).

- 281 [9] Sijia Li et al. *Scalable Bigraphical Lasso: Two-way Sparse Network Inference for Count*
282 *Data*. arXiv:2203.07912 [cs, stat]. Mar. 2022. URL: <http://arxiv.org/abs/2203.07912>
283 (visited on 02/24/2023).
- 284 [10] Sameer A Nene, Shree K Nayar, and Hiroshi Murase. “Columbia Object Image Library
285 (COIL-20)”. en. In: ().
- 286 [11] David Ouyang et al. “Video-based AI for beat-to-beat assessment of cardiac function”. en. In:
287 *Nature* 580.7802 (Apr. 2020). Number: 7802 Publisher: Nature Publishing Group, pp. 252–256.
288 ISSN: 1476-4687. DOI: 10.1038/s41586-020-2145-8. URL: [https://www.nature.com/](https://www.nature.com/articles/s41586-020-2145-8)
289 [articles/s41586-020-2145-8](https://www.nature.com/articles/s41586-020-2145-8) (visited on 05/10/2023).
- 290 [12] Vincent Prost, Stéphane Gazut, and Thomas Bröls. “A zero inflated log-normal model for
291 inference of sparse microbial association networks”. en. In: *PLOS Computational Biology*
292 17.6 (June 2021). Publisher: Public Library of Science, e1009089. ISSN: 1553-7358. DOI: 10.
293 1371/journal.pcbi.1009089. URL: [https://journals.plos.org/ploscompbiol/](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1009089)
294 [article?id=10.1371/journal.pcbi.1009089](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1009089) (visited on 02/24/2023).
- 295 [13] Eetje F. Tigchelaar et al. “Cohort profile: LifeLines DEEP, a prospective, general population
296 cohort study in the northern Netherlands: study design and baseline characteristics”. en. In:
297 *BMJ Open* 5.8 (Aug. 2015). Publisher: British Medical Journal Publishing Group Section: Epi-
298 demiology, e006772. ISSN: 2044-6055, 2044-6055. DOI: 10.1136/bmjopen-2014-006772.
299 URL: <https://bmjopen.bmj.com/content/5/8/e006772> (visited on 04/30/2023).