

Bayesian Semi-nonnegative Tri-matrix Factorization to Identify Associations between Pathways and sub-types in cancer data

Sunho Park, Nabhonil Kar and Tae Hyun Hwang

Quantitative Health Sciences Cleveland Clinic,

9500 Euclid Ave. Cleveland, OH 44195

{parks, karn2, hwangt}@ccf.org

https://www.lerner.ccf.org/qhs/

This supplementary material provides more details about our proposed method. We first summarize the probability distributions used in our main paper in the following table.

| Distribution | PDF | mean | variance | note |
|-------------------------------|---|--|---|---|
| Bernoulli($z \rho$) | $\rho^z(1-\rho)^{(1-z)}$ | ρ | $\rho(1-\rho)$ | $z \in \{0, 1\}, \rho \in [0, 1]$ |
| Gamma($\tau a, b$) | $\frac{1}{\Gamma(a)}b^a\tau^{a-1}e^{-b\tau}$ | $\frac{a}{b}$ | $\frac{a}{b^2}$ | $\tau > 0, a > 0, b > 0$ |
| Exponential($s \lambda$) | $\lambda e^{-\lambda s}$ | λ^{-1} | λ^{-2} | $s \in [0, \infty]$ |
| $\mathcal{N}(x \mu, \sigma)$ | $\frac{1}{\sqrt{2\pi}\sigma}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ | μ | σ | $x \in \mathbb{R}$ |
| $\mathcal{TN}(x \mu, \sigma)$ | $\frac{\mathcal{N}(x \mu, \sigma)}{1-\Phi(-\frac{\mu}{\sqrt{\sigma}})}$ | $\mu + \sqrt{\sigma}h_1(-\frac{\mu}{\sqrt{\sigma}})$ | $\sqrt{\sigma}\left[1 - h_2(-\frac{\mu}{\sqrt{\sigma}})\right]$ | $x \in \mathbb{R}_+, h_1(x) = \frac{\mathcal{N}(x 0,1)}{1-\Phi(x)},$ $h_2(x) = h_1(x)[h_1(x) - x]$ |

1. Model Summary

The observation matrix is decomposed into the sub-matrices in the following way:

$$\mathbf{X} \approx \mathbf{U}\mathbf{S}\mathbf{V}^\top = \mathbf{U}\mathbf{S}(\mathbf{Z} \circ \mathbf{V})^\top, \quad (1)$$

where \circ stands for an element-wise multiplication operator. Denoting all the latent variables by $\Theta \triangleq \{\mathbf{S}, \mathbf{V}, \mathbf{Z}, \mathbf{G}\}$, the joint probability distributions of the model is given as follows

$$p(\mathbf{X}, \Theta) = p(\mathbf{X}|\mathbf{S}, \mathbf{Z}, \mathbf{V}, \tau)p(\tau)p(\mathbf{S})p(\mathbf{V}, \mathbf{Z}|\mathbf{G})p(\mathbf{G}). \quad (2)$$

where

$$p(\mathbf{X}|\mathbf{U}, \mathbf{S}, \mathbf{Z}, \mathbf{V}, \gamma) = \prod_{(i,j) \in \Omega} p(X_{ij}|\mathbf{u}_i^\top \mathbf{S}(\mathbf{z}_j \circ \mathbf{v}_j), \tau) \quad (3)$$

$$p(\mathbf{S}) = \prod_{k=1}^K \prod_{r=1}^R p(S_{kr}), \quad (4)$$

$$p(\mathbf{V}, \mathbf{Z}|\mathbf{G}) = \prod_{r=1}^R \prod_{j=1}^D p(V_{jr}Z_{jr}|G_{jr}), \quad (5)$$

$$p(\mathbf{G}) = \prod_{r=1}^R p(\vec{g}_r|\mathbf{m}_r, \mathbf{L}), \quad (6)$$

where $\mathbf{L} = \mathbf{I} - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$ is a normalized Laplacian matrix and \mathbf{A} is an adjacency matrix driven from a protein-protein interaction network ($A_{ij} = 1$ if $i \neq j$ and there is a connection between the genes i and j on the network, and otherwise $A_{ij} = 0$). Note that, the mean vector \mathbf{m}_r is set according to the membership information encoded in the pathways Z^0 : $m_{jr} = \xi_+$ if $Z_{jr}^0 = 1$, otherwise $m_{jr} = \xi_-$, where $\xi_+ > 0$ and $\xi_- < 0$ (in our all experiments, we use $\xi_+ = 3$ and $\xi_- = -5$). The form of each probability distribution in (2) is given as follows.

$$p(X_{ij} | \mathbf{u}_i^\top \mathbf{S}(\mathbf{z}_j \circ \mathbf{v}_j), \gamma) = \mathcal{N}(X_{ij} | \mathbf{u}_i^\top \mathbf{S}(\mathbf{z}_j \circ \mathbf{v}_j), \gamma) \quad (7)$$

$$p(\gamma) = \text{Gamma}(\gamma | \alpha_a^0, \alpha_b^0), \quad (8)$$

$$p(S_{kr}) = \text{Exponential}(S_{kr} | \lambda_{kr}^{S0}), \quad (9)$$

$$p(V_{jr}, Z_{jr} | G_{jr}) = \mathcal{N}(Z_{jr} V_{jr} | 0, \sigma_{jr}^{V0}) (\rho_{jr}(G_{jr}))^{Z_{jr}} (1 - \rho_{jr}(G_{jr}))^{(1-Z_{jr})}. \quad (10)$$

2. Variational Inference

The posterior distributions over the latent variables are approximately computed in the framework of variational inference. The variational distributions that approximate the true posterior distributions over the latent variables are assumed to be factorized as follows

$$q(\Theta) = q(\gamma) \left(\prod_{k=1}^K \prod_{r=1}^R q(S_{kr}) \right) \left(\prod_{j=1}^D \prod_{r=1}^R q(V_{jr}, Z_{jr}) q(G_{jr}) \right), \quad (11)$$

where

$$q(\gamma) = \text{Gamma}(\gamma | \alpha_a, \alpha_b), \quad (12)$$

$$q(S_{kr}) = \mathcal{TN}(S_{kr} | \mu_{kr}^S, \sigma_{kr}^S), \quad (13)$$

$$q(V_{jr}, Z_{jr}) = \mathcal{N}(V_{jr} | Z_{jr} \mu_{jr}^V, Z_{jr} \sigma_{jr}^V + (1 - Z_{jr}) \sigma_{jr}^{V0}) \hat{\rho}_{jr}^{Z_{jr}} (1 - \hat{\rho}_{jr})^{(1-Z_{jr})}, \quad (14)$$

$$q(G_{jr}) = \mathcal{N}(G_{jr} | \mu_{jr}^g, \sigma_{jr}^g). \quad (15)$$

The variational distributions can be computed by maximizing the lower bound with respect to (w.r.t) the variational distributions. Denoting a set of all the latent variables by $\Theta = \{\gamma, \mathbf{S}, \mathbf{Z}, \mathbf{V}, \mathbf{G}\}$, we can show that the log-likelihood can be decomposed as follows

$$\log p(\mathbf{X}) = \mathcal{L}(q) + \text{KL}(q||p) \quad (16)$$

where

$$\mathcal{L}(q) = \int q(\Theta) \log \frac{p(\mathbf{X}, \Theta)}{q(\Theta)} d\Theta, \quad (17)$$

$$\text{KL}(q||p) = - \int q(\Theta) \log \frac{p(\Theta | \mathbf{X})}{q(\Theta)} d\Theta. \quad (18)$$

where $\text{KL}(q||p)$ is Kullback-Leibler (KL) divergence between the variational distribution and the true posterior distribution and is always nonnegative ($\text{KL}(q||p) = 0$ if and only if $q=p$). Thus, we can easily see that the log-likelihood is lower-bound by the variational lower bound $\mathcal{L}(q)$ and thus the variational distributions can be updated by maximizing $\mathcal{L}(q)$ w.r.t their

parameters. Note that, the variational bound of our model is expressed as follows

$$\begin{aligned}
\mathcal{L}(q) &= \mathbb{E}_{q(\Theta)} \left[\log \frac{p(\mathbf{X}, \Theta)}{q(\Theta)} \right] \\
&= \mathbb{E}_{q(\Theta)} \left[\log p(\mathbf{X} | \mathbf{S}, \mathbf{Z}, \mathbf{V}, \gamma) p(\gamma) p(\mathbf{S}) p(\mathbf{V}, \mathbf{Z} | \mathbf{G}) p(\mathbf{G}) \right] \\
&\quad - \mathbb{E}_{q(\Theta)} \left[\log q(\tau) q(\mathbf{S}) q(\mathbf{V}, \mathbf{Z}) q(\mathbf{G}) \right].
\end{aligned} \tag{19}$$

The variational distributions $q(\gamma)$, $\{q(S_{kr})\}$ and $\{q(V_{jr}, Z_{jr})\}$ can be updated in closed form. Letting Θ_l be the variable we want to update at each turn and $\Theta^{\setminus l}$ be the remaining variables, the optimal solution of $q(\Theta_l)$ can be given by the stationary condition for the factor $q(\Theta_l)$ in the maximization problem, i.e., maximize $_{q(\Theta_l)} \mathcal{L}(q)$:

$$\log q(\Theta_l) \propto \mathbb{E}_{q(\Theta^{\setminus l})} [\log(p(\mathbf{X}, \Theta))]. \tag{20}$$

On the other hand, for $\{q(G_{jr})\}$, their means and variances can be updated by any iterative gradient-based optimization methods, e.g., limited-memory BFGS used in our experiments. We provide detailed derivations of each update in the following subsections.

2.1. Update of the variational distributions over γ and \mathbf{S}

The variational distribution over the precision γ can be updated as follows

$$q(\tau) = \text{Gamma}(\tau | \hat{\alpha}_a, \hat{\alpha}_b) \tag{21}$$

where

$$\hat{\alpha}_a = \alpha_a^0 + \frac{|\Omega|}{2}, \tag{22}$$

$$\hat{\alpha}_b = \alpha_b^0 + \frac{1}{2} \sum_{(i,j) \in \Omega} \mathbb{E}_q \left[\left(X_{ij} - \mathbf{u}_i^\top \mathbf{S}(\mathbf{z}_j \circ \mathbf{v}_j) \right)^2 \right], \tag{23}$$

where Ω is a set of indices of the observations and

$$\begin{aligned}
&\mathbb{E}_{q(\Theta)} \left[\left(X_{ij} - \mathbf{u}_i^\top \mathbf{S}(\mathbf{z}_j \circ \mathbf{v}_j) \right)^2 \right] \\
&= \left(X_{ij} - \sum_{k=1}^K \sum_{r=1}^R \hat{U}_{ik} \langle S_{kr} \rangle \langle Z_{jr} V_{jr} \rangle \right)^2 + \sum_{k=1}^K \sum_{r=1}^R \left[\hat{U}_{ik}^2 \langle S_{kr}^2 \rangle \langle Z_{jr} V_{jr}^2 \rangle - \hat{U}_{ik}^2 \langle S_{kr} \rangle^2 \langle Z_{jr} V_{jr} \rangle^2 \right] \\
&\quad + \sum_{k=1}^K \sum_{r=1}^R \sum_{k' \neq k}^K \left[\hat{U}_{ik} \langle S_{kr} \rangle \left(\langle Z_{jr} V_{jr}^2 \rangle - \langle Z_{jr} V_{jr} \rangle^2 \right) \hat{U}_{ik'} \langle S_{k'r} \rangle \right],
\end{aligned} \tag{24}$$

where $\langle Z_{jr} V_{jr} \rangle = \rho_{jr} \mu_{jr}^V$ and $\langle Z_{jr}^2 V_{jr}^2 \rangle = \langle Z_{jr} V_{jr}^2 \rangle = \rho_{jr} (\sigma_{jr}^V + (\mu_{jr}^V)^2)$.

The variable \mathbf{S} can be updated as follows

$$q(S_{kr}) = \mathcal{TN}(S_{kr} | \mu_{kr}^S, \sigma_{ij}^S), \tag{25}$$

where

$$\sigma_{kr}^S = \left(\langle \gamma \rangle \sum_{(i,j) \in \Omega} \hat{U}_{ik}^2 \langle Z_{jr} V_{jr}^2 \rangle \right)^{-1}, \quad (26)$$

$$\begin{aligned} \mu_{kr}^S = \sigma_{kr}^S & \left[-\lambda_{kr}^S + \langle \gamma \rangle \sum_{(i,j) \in \Omega} \left(\left(X_{ij} - \sum_{(k',r') \neq (k,r)} \hat{U}_{ik'} \langle S_{k'r'} \rangle \langle Z_{jr'} V_{jr'} \rangle \right) \hat{U}_{ik} \langle Z_{jr} V_{jr} \rangle \right. \right. \\ & \left. \left. - \hat{U}_{ik} \left(\langle Z_{jr} V_{jr}^2 \rangle - \langle Z_{jr} V_{jr} \rangle^2 \right) \sum_{k' \neq k} \hat{U}_{ik'} \langle S_{k'r'} \rangle \right) \right]. \end{aligned} \quad (27)$$

2.2. Update of the variational distributions over Z and V

Each pair of elements, $\{Z_{jr}, V_{jr}\}$, can be updated by the inference method in.¹ From the stationary condition for $q(Z_{jr}, V_{jr})$ when maximizing the variational bound \mathcal{L} in (19), we have

$$q(V_{jr}, Z_{jr}) = \frac{1}{\mathcal{Z}} \exp \left\{ \langle \log p(X|\Theta) \rangle \mathcal{N}(Z_{jr} V_{jr} | 0, \sigma_{jr}^{V0}) \langle \Phi(G_{jr}) \rangle^{Z_{jr}} \langle (1 - \Phi(G_{jr})) \rangle^{(1-Z_{jr})} \right\}, \quad (28)$$

where \mathcal{Z} is a normalization constant. We can see that $q(V_{jr}, Z_{jr})$ can be factorized as

$$q(V_{jr}, Z_{jr}) = q(V_{jr} | Z_{jr}) q(Z_{jr}). \quad (29)$$

The marginal probability distribution over the binary variable Z_{jr} can be calculated as follows

$$q(Z_{jr} = 1) = \rho_{jr} = \frac{1}{1 + \exp \{-\xi_{jr}\}}, \quad (30)$$

where

$$\begin{aligned} \xi_{jr} &= \log q(Z_{jr} = 1) - \log q(Z_{jr} = 0) \\ &= \langle \log \Phi(G_{jr}) \rangle - \langle \log (1 - \Phi(G_{jr})) \rangle - \frac{1}{2} \log \sigma_{jr}^{V0} + \frac{1}{2} \frac{(\mu_{jr}^V)^2}{\sigma_{jr}^V} + \frac{1}{2} \log \sigma_{jr}^V. \end{aligned} \quad (31)$$

where the expectations in the second equality are approximated using Jensen's inequality:

$$\langle \log \Phi(G_{jr}) \rangle \approx \log \Phi \left(\frac{\mu_{jr}^g}{\sqrt{1 + \sigma_{jr}^g}} \right), \quad (32)$$

$$\langle \log (1 - \Phi(G_{jr})) \rangle = \langle \log (\Phi(-G_{jr})) \rangle \approx \log \Phi \left(\frac{-\mu_{jr}^g}{\sqrt{1 + \sigma_{jr}^g}} \right) \quad (33)$$

The conditional variational distribution of V_{jr} given Z_{jr} is given by

$$q(V_{jr} | Z_{jr} = 0) = \mathcal{N}(V_{jr} | 0, \sigma_{jr}^{V0}), \quad (34)$$

$$q(V_{jr} | Z_{jr} = 1) = \mathcal{N}(V_{jr} | \mu_{jr}^V, \sigma_{jr}^V), \quad (35)$$

where

$$\sigma_{jr}^V = \left[(\sigma_{jr}^{V0})^{-1} + \langle \tau \rangle \sum_{i \in \Omega_j} \left(\left(\sum_{k=1}^K \hat{U}_{ik} \langle S_{kr} \rangle \right)^2 + \sum_{k=1}^K \hat{U}_{ik}^2 \left(\langle S_{kr}^2 \rangle - \langle S_{kr} \rangle^2 \right) \right) \right]^{-1}, \quad (36)$$

$$\mu_{jr}^V = \sigma_{jr}^V \left[\frac{\mu_{jr}^{V0}}{\sigma_{jr}^{V0}} + \langle \tau \rangle \sum_{i \in \Omega_j} \left(\left(X_{ij} - \sum_{k=1}^K \sum_{r' \neq r} \hat{U}_{ik} \langle S_{kr'} \rangle \langle Z_{jr'} V_{jr'} \rangle \right) \sum_{k=1}^K \hat{U}_{ik} \langle S_{kr} \rangle \right) \right]. \quad (37)$$

As a summary, the joint probability distribution is simply rewritten as follow

$$q(V_{jr}, Z_{jr}) = \mathcal{N}\left(V_{jr}|Z_{jr}\mu_{jr}^V, Z_{jr}\sigma_{jr}^V + (1 - Z_{jr})\sigma_{jr}^{V0}\right)\rho_{jr}^{Z_{jr}}(1 - \rho_{jr})^{1-Z_{jr}}. \quad (38)$$

2.3. Update of the variational distributions over \mathbf{G}

The optimization problem (19) can be reduced as follows

$$\text{maximize}_{q(\mathbf{G})} \mathcal{L}_g, \quad (39)$$

where \mathcal{L}_g is a function including only terms which are related to the variable \mathbf{G} :

$$\mathcal{L}_g = \mathbb{E}_{q(\mathbf{Z})q(\mathbf{G})} \left[\log p(\mathbf{Z}|\mathbf{G})p(\mathbf{G}) \right] - \mathbb{E}_{q(\mathbf{G})} \left[\log q(\mathbf{G}) \right]. \quad (40)$$

The first term of \mathcal{L}_g in eq. (39) can be calculated as follows

$$\begin{aligned} & \mathbb{E}_{q(\mathbf{Z})q(\mathbf{G})} \left[\log p(\mathbf{Z}|\mathbf{G}) \right] \\ &= \sum_{j,r} \langle Z_{jr} \rangle \langle \log \Phi(G_{jr}) \rangle + \langle (1 - Z_{jr}) \rangle \langle \log(1 - \Phi(G_{jr})) \rangle \\ &\approx \sum_{j,r} \rho_{jr} \log \Phi\left(\frac{\mu_{jr}^g}{\sqrt{1 + \sigma_{jr}^g}}\right) + (1 - \rho_{jr}) \log \Phi\left(\frac{-\mu_{jr}^g}{\sqrt{1 + \sigma_{jr}^g}}\right), \end{aligned} \quad (41)$$

where we have used the same techniques (using Jensen's inequality) as in the previous subsection. We then calculate the third term, a sum of entropy terms of R Gaussian distributions:

$$-\mathbb{E}_{q(\mathbf{G})} \left[\log q(\mathbf{G}) \right] = \sum_{r=1}^R H(q(\vec{g}_r)) = \frac{1}{2} \sum_{r=1}^R \log \left(\prod_{j=1}^N \sigma_{jr}^g \right) + c \quad (42)$$

where c is a constant, which is independent of the variable \mathbf{G} . The second term is a cross entropy between two Gaussian distributions, $p(\mathbf{G})$ and $q(\mathbf{G})$, calculated as follows:

$$\begin{aligned} \mathbb{E}_{q(\mathbf{G})} \left[\log p(\mathbf{G}) \right] &= \sum_{r=1}^R -H(q(\vec{g}_r)) - \text{KL}(q(\vec{g}_r)|p(\vec{g}_r)) \\ &= -\frac{1}{2} \sum_{r=1}^R \left(\left((\boldsymbol{\mu}_r^g - \mathbf{m}_r)^\top \mathbf{L}^{-1} (\boldsymbol{\mu}_r^g - \mathbf{m}_r) \right) + \left(\sum_{j=1}^D \sigma_{jr}^g [\mathbf{L}^{-1}]_{jj} \right) \right). \end{aligned} \quad (43)$$

The gradient of \mathcal{L}_g w.r.t. the parameters $\{\mu_{jr}^g, \sigma_{jr}^g\}$ also can be easily calculated. We update these parameters using limited-memory BFGS in our experiments

3. Experimental settings

We here explain how to initialize our factorization model. We need to initialize the parameters of the prior distributions (referred to as prior hyperparameters) and the parameters of the variational distributions (referred to as variational parameters). Regarding to the variational parameters, e.g., the mean and variance for the case where the variational distribution is a Gaussian distribution, note that the variational distributions are updated cyclically, i.e., for each step, we update one variational distributions, fixing the others. Thus, we need to initialize some (not all) variational parameters as well as the prior hyperparameters.

The prior hyperparameters are set as follows.

- (For the noise precision γ) $\alpha_a^0 = \alpha_b^0 = 0.1$
- (For each association S_{kr}) $\lambda_{kr}^{S0} = 10$ for all k, r
- (For each element in the centroid, V_{jr}) $\mu_{jr}^{V0} = 0$ and $\sigma_{jr}^{V0} = 1$
- (For the prior mean vectors of the GPs, \mathbf{m}_r) $\xi_+ = 5$ and $\xi_- = -5$

Based on our experiences on synthetic simulation and multiple gene expression datasets, the factorization result of the method is not sensitive to the most parameters' initial settings in general. However, we should note that ξ_+ and ξ_- represent the prior belief on the initial pathway membership information Z^0 . If we set $\xi_+ = \xi_- = 0$, the prior probability of the on-off binary variable $Z_{jr} = 1$ is 0.5 regardless of whether the r th pathway includes the j th gene or not, i.e., $Z_{jr}^0 = 1$ or $Z_{jr}^0 = 0$. The more extreme values ξ_+ and ξ_- have, i.e. $\xi_+ \gg 0$ and $\xi_- \ll 0$, the stronger the prior belief we place on the initial pathway information \mathbf{Z}^0 . The setting we use in our experiments, i.e., $\xi_+ = 5$ and $\xi_- = -5$, usually gives satisfactory factorization results. However, users can adjust them according to their prior belief on pathway information.

The variational parameters are set as follows.

- (For S_{kr}) $\mu_{kr}^S \sim \text{Uniform}([0, 1])$, for all k, r
- (For \mathbf{V}) each column is set to the centroid from K-means run on the input data \mathbf{X} if $R < N$, and to the randomly chosen sample (row) from the input matrix \mathbf{X} otherwise.
- (For G_{jr}) $\mu_{jr}^g = m_{jr}$ and $\sigma_{jr}^g = \exp(-\zeta)$, where $\zeta \sim \mathcal{N}(0, 0.1^2)$

4. Bayesian semi-nonnegative- vs Point estimate non-negative factorization

As explained in the main paper, many types of genomic data are given in a form of real-valued matrix after relevant normalization or transformation steps. However, the NMF formulation do not allow negative values in the observation matrix. One of the standard ways to handle negative values for the NMF formulation is to fold the original matrix by columns:² every column (gene) will be represented in two new columns in a new matrix, one of which contains only the positive values (up-regulations) and the other column only the magnitudes of the negative values (down-regulations). However, the folding approach would cause computational difficulties if it was applied to our tri-matrix factorization problem. In addition to the doubled number of columns, the gene-gene interaction network will requires 2^3 times larger computational complexity to calculate the inverse matrix (please refer eq. (43)). We instead follow the concept of the semi-nonnegative matrix proposed in,³ where the matrix containing centroids (bases vectors) is allowed to have negative values in the matrix but the encoding matrix is still restricted to be nonnegative. Furthermore, we formulate the semi-nonnegative (tri-)matrix factorization in the framework of Bayesian learning. We provide two specific examples that show the superiority of our method over the nonnegative tri-matrix factorization (NTriPath⁴ implemented with the folding approach to deal with negative values).

4.1. Baseline example

We begin by presenting a simple example that contains a basic structure in the observation matrix and other inputs. We discuss the results of this straightforward settings before examining the two cases in which our proposed Bayesian framework provides clear advantages over

the deterministic approach. A detailed overview of data generation is now presented for our first example.

Inspired by the biological setting of the gene expression application in the main paper, we use the same terminology here as in the main script (i.e. subgroups, genes, pathways) in discussing the problem formulation and results of our simulated experiments. As in the main application, these experiments attempt to decompose patterns of upregulated and downregulated genes within different patient subgroup. In this preliminary example, as well as in subsequent experiments, the observation matrix $\mathbf{X} \in \mathbb{R}^{200 \times 800}$ consists of 4 subgroups (each containing 50 samples) with some defined pattern among the 800 genes (which are grouped into sets of 100). Within each subgroup, a set of 100 genes can represent upregulation, downregulation or background noise. For upregulated and downregulated genes, samples are drawn from a Gaussian $\mathcal{N}(1, 2)$ or Gaussian $\mathcal{N}(-1, 2)$, respectively. For background noise, samples are drawn from a Gaussian $\mathcal{N}(0, 1)$. A simple block structure was determined with each subsample containing 2-3 "selected" (either upregulated or downregulated) gene sets (see Figure 1a). Subgroups are encoded in $\mathbf{U} \in \mathbb{R}_+^{200 \times 4}$ using simple 1-of- K encoding ($U_{ij} \in \{0, 1\}$ and $\sum_j U_{ij} = 1$) (see Fig. 1b). Gene-pathway prior knowledge is encoded in the matrix $\mathbf{Z}^0 \in \mathbb{R}_+^{800 \times 4}$ (expressed as $\mathbf{V}_+^0 \in \mathbb{R}^{800 \times 4}$ in the deterministic settings) and is initialized to contain similar structure to the observation \mathbf{X} . That is, we allow \mathbf{Z}^0 to contain 4 pathways that reflect the same pattern of selected genes within the subgroups of \mathbf{X} by setting $\mathbf{Z}_{ij}^0 = 1$ for all the genes i that are upregulated or downregulated within the subgroup that we have designated to pathway j (see Fig. 1c). The gene-gene interaction network $\mathbf{A} \in \mathbb{R}^{800 \times 800}$ is initialized with approximately 10% sparsity and contains random symmetric connections (with no self connections; see Fig. 1d). Our motivation for this simple design is to allow our method to arrive at a simple and predictable solution for the model's learned factors, namely, the subgroup-pathway association matrix $\mathbf{S} \in \mathbb{R}_+^{4 \times 4}$, the real-valued pathway-gene association matrix $\mathbf{V} \in \mathbb{R}^{800 \times 4}$ and the updated pathway-gene binary membership matrix $\mathbf{Z} \in \mathbb{R}_+^{800 \times 4}$.

4.2. Robustness against noise

We compare the performance of our Bayesian factorization method and NTriPath in the case where an input matrix is contaminated by background noise with different noise levels. Our objective here is to show whether each method is robust against noise. We assume that the observation matrix $\widetilde{\mathbf{X}}$ is generated by adding white Gaussian noises to the data input matrix \mathbf{X} defined in Section 4.1, i.e., $\widetilde{\mathbf{X}} = \mathbf{X} + \widetilde{\mathbf{E}}$, where $\widetilde{E}_{ij} \sim \mathcal{N}(0, \gamma_n^{-1})$ and the noise variances γ_n^{-1} increases from 1^2 to 10^2 . We train both methods on the noisy observation matrix $\widetilde{\mathbf{X}}$ and test how each method correctly identifies the associations between the sub-types and the pathways.

Before conducting the experiments, we first explain how we select the hyperparameters for NTriPath whose model formulation is given as follows

$$\text{minimize}_{\mathbf{S} \geq 0, \mathbf{V} \geq 0} \frac{1}{2} f(\mathbf{S}, \mathbf{V}), \quad (44)$$

where the objective function is define as:

$$f(\mathbf{S}, \mathbf{V}) = \|\mathbf{X} - \mathbf{U}\mathbf{S}\mathbf{V}^\top\|_F^2 + \lambda_S \|\mathbf{S}\|_1^2 + \lambda_V \|\mathbf{V}\|_1^2 + \lambda_Z \|\mathbf{V} - \mathbf{Z}^0\|_F^2 + \lambda_{V_L} \text{tr}\{\mathbf{V}^\top \mathbf{L} \mathbf{V}\}. \quad (45)$$

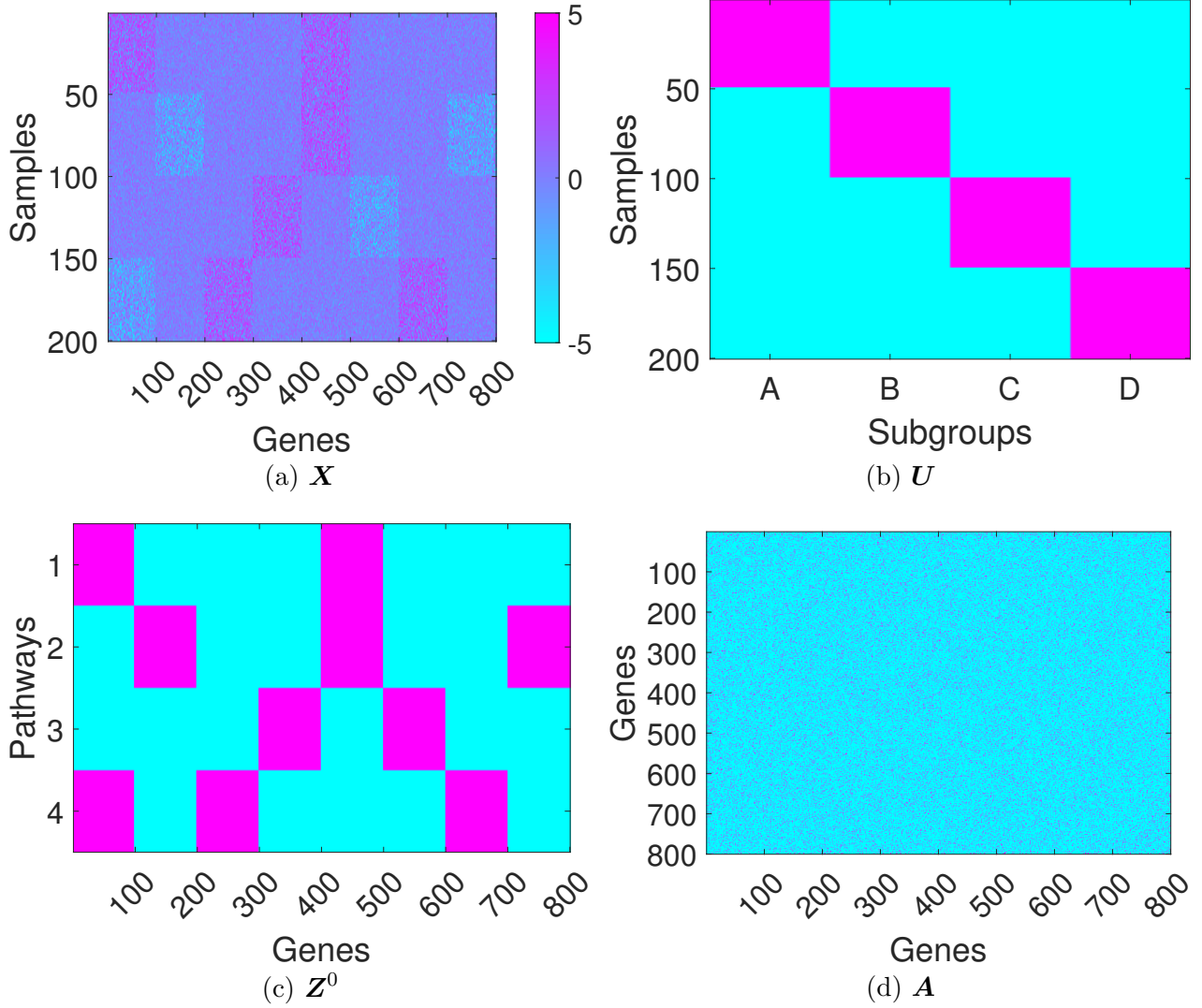


Fig. 1. Inputs

NTriPath involves 4 regularization constants which should be specified by user. Although we know the ground truth associations between the sub-types and the pathways for this simple dataset, identifying associations from input data is clearly an unsupervised learning problem since true associations between sub-types and pathways are not known in general. Thus, it is not clear how to tune the regularization parameters for the given input data. In addition, it requires high computational burden for large scale datasets when the best combination of the hyperparameters is searched among a set of candidate values, (the search space will be 4D grid space) by cross validation (CV). For simplicity, we fix $\lambda_{V_L} = \lambda_Z = 1$ as in our previous work. We tune only λ_S and λ_V which are related to the sparseness of the metrics. We select the best combination from 2D grid space (each grid space is defined as $\{0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1\}$) by comparing the least reconstruction error. On the other hand, our method is able to automatically tune the model complexity by integrating over the latent variables.

We report the performance of both methods in Figure 4.2. Since we know the ground truth

associations for this dataset, we can calculate the accuracy of each method based on how many associations each method correctly predicts. We repeat each experiment 20 times at each noise variance. As we can see in the figure, our method shows overall stable performance.

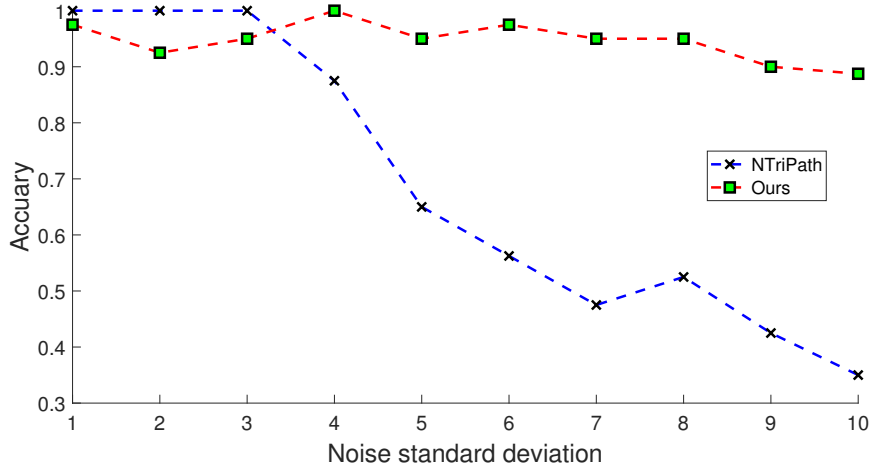


Fig. 2. The prediction performance of our method compared to NTriPath.

5. TCGA gastric cancer and metastatic gastric cancer immunotherapy clinical-trial datasets: additional information

We here include the list of the selected pathways from both data sets in the experimental results section in the main text. Please see Table 2 for the TCGA gastric cancer data and Table 3 for the metastatic gastric cancer immunotherapy clinical-trial data.

References

1. M. K. Titsias and M. Lázaro-Gredilla, Spike and slab variational inference for multi-task and multiple kernel learning, in *Advances in Neural Information Processing Systems (NIPS)*, eds. J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira and K. Q. Weinberger 2011 pp. 2339–2347.
2. P. Kim and B. Tidor, B. subsystem identification through dimensionality reduction of large-scale gene expression dat, *Genome Research* **13**, p. 17061718 (2003).
3. C. Ding, T. Li and M. I. Jordan, *Convex and Semi-Nonnegative Matrix Factorizations*, Tech. Rep. 60428, Lawrence Berkeley National Lab (2006).
4. S. Park, S.-J. Kim, D. Yu, S. Pea-Llopis, J. Gao, J. S. Park, B. Chen, J. Norris, X. Wang, M. Chen, M. Kim, J. Yong, Z. Wardak, K. Choe, M. Story, T. Starr, J.-H. Cheong and T. H. Hwang, An integrative somatic mutation analysis to identify pathways linked with survival outcomes across 19 cancer types, *Bioinformatics* **32**, 1643 (2016).

Table 2. Summary of the top 3-ranked pathways associated with the molecular sub-types obtained from the TCGA gastric cancer dataset.

| sub-types | rank | #members | member genes |
|-----------|------|----------|--|
| CIN | 1 | 12 | ADAMTS4,CELA1,CTRB1,DERL1,DERL2,DERL3,KLK5,MFI2,MMP26,PRSS1,PRSS3,SERPINA1 |
| | 2 | 14 | COL2A1,COL3A1,COL9A1,COL9A2,COL9A3,COMP,FN1,MAG,MAP1B,MBP,NGFR,PLP1,PRNP,RTN4R |
| | 3 | 13 | BARD1,BRCA1,CSTF1,CSTF2,CSTF3,FEZ1,HTATSF1,IKBKAP,MED21,PIN1,POLR2A,RBBP8,SUPT5H |
| EBV | 1 | 12 | ADAMTS4,CELA1,CTRB1,DERL1,DERL2,DERL3,KLK5,MFI2,MMP26,PRSS1,PRSS3,SERPINA1 |
| | 2 | 13 | C3,F2,F2RL3,FCER2,HP,ICAM2,ICAM4,ITGAM,ITGAX,ITGB2,JAM2,JAM3,TJP1 |
| | 3 | 14 | CD44,EED,FN1,ICAM4,ITGA4,ITGAE,ITGB1,ITGB7,LGALS8,MADCAM1,PXN,TLN1,VCAM1,VCAN |
| GS | 1 | 10 | CALM1,CPE,GCG,GLP1R,GPRASP1,GRM5,MEP1A,MEP1B,OPRM1,VIPR1 |
| | 2 | 1 | GNAQ |
| | 3 | 2 | CD200,CD200R1 |
| MSI | 1 | 3 | CCDC67,CCDC85B,EIF3E |
| | 2 | 1 | GNAQ |
| | 3 | 3 | NUP155,NUPL2,ZFYVE9 |

Table 3. Summary of the top 3-ranked pathways associated with the treatment response obtained from the metastatic gastric cancer immunotherapy clinical-trial dataset.

| sub-types | rank | #members | member genes |
|---------------|------|----------|--|
| responder | 1 | 14 | ATN1,ECM1,ELN,FBLN1,FBLN2,FBN1,FBN2,FN1,HSPG2,ITGB1,LTBP1,MFAP2,PRELP,VCAN |
| | 2 | 11 | CCL19,CCL21,CCL5,CCR3,CXCL11,CXCL13,CXCL9,DPP4,IGFBP7,PF4,VCAN |
| | 3 | 10 | CCL11,CCL5,CCR3,CPAMD8,CXCL11,CXCL13,CXCL9,DPP4,FAP,PF4 |
| non-responder | 1 | 3 | SLC1A4,SLC1A5,TBC1D17 |
| | 2 | 3 | CCDC85B,KRTAP4-12,LMO2 |
| | 3 | 3 | NMU,NMUR1,NMUR2 |