

## A Proofs

### A.1 Proof of theorem [1](#)

$\mathbf{W} \in \mathcal{S}_W$  is the weight matrix of a graph with  $R$  connected components  $\{C_1, \dots, C_R\}$  partitioning  $[n]$ . Since  $k$  is upper bounded by a constant, there exists  $M_+ > 1$  that upper bounds  $k$ . Let  $\mathcal{T}$  be the adjacency matrix of a spanning forest of  $\mathbf{W}$ , since each edge of  $\mathbf{W}$  is bounded by  $n$ , one has:

$$\begin{aligned} \int f_k(\mathbf{X}, \mathbf{W}) \lambda_{\mathcal{S}_C}(d\mathbf{X}) &= \int \prod_{(i,j) \in [n]^2} k(\mathbf{X}_i - \mathbf{X}_j)^{W_{ij}} \lambda_{\mathcal{S}_C}(d\mathbf{X}) \\ &\leq M_+^{n^3} \int \prod_{(i,j) \in [n]^2} k(\mathbf{X}_i - \mathbf{X}_j)^{\mathcal{T}_{ij}} \lambda_{\mathcal{S}_C}(d\mathbf{X}) \\ &\leq M_+^{n^3} \prod_{r \in [R]} \int \prod_{(i,j) \in C_r^2} k(\mathbf{X}_i - \mathbf{X}_j)^{\mathcal{T}_{ij}} \lambda_{\mathcal{S}_C}(d\mathbf{X}). \end{aligned} \quad (10)$$

Let  $r \in [R]$ . The spanning tree corresponding to the  $r^{\text{th}}$  connected component called  $\mathcal{T}^r$  has exactly  $n_r - 1$  edges. There exists a leaf node  $\ell \in [n]$  of  $\mathcal{T}^r$  and let  $\tilde{\ell}$  be the node linked to it. Consider a bijective map  $\sigma: C_r \setminus \{\ell\} \rightarrow [n_r - 1]$  such that  $\sigma(\tilde{\ell}) = 1$  and for  $(i, j) \in (C_r \setminus \{\ell\})^2$ ,  $\sigma(i) \leq \sigma(j)$  implies that node  $i$  has a shorter path on  $\overline{\mathcal{T}^r}$ <sup>1</sup> to  $\ell$  than node  $j$ . There exists a bijective map  $e: [2 : n_r - 1] \rightarrow [n_r - 2]$  such that for  $i \in [2 : n_r - 1]$ ,  $\overline{\mathcal{T}^r}_{\sigma^{-1}(i), \sigma^{-1}(e(i))} > 0$  and node  $\sigma^{-1}(e(i))$  has a shorter path on  $\overline{\mathcal{T}^r}$  to node  $\ell$  than node  $\sigma^{-1}(i)$ . Recall that since  $\mathbf{X} \in \mathcal{S}_C$  one has:  $\sum_{i \in C_r} \mathbf{X}_i = 0$  hence  $\mathbf{X}_\ell = -\sum_{i \neq \ell} \mathbf{X}_i$ . Let us now consider the linear map  $\phi^r$  such that:

$$\forall i \in [n_r - 1], \quad \phi^r(\mathbf{X}_i) = \begin{cases} \mathbf{X}_{\sigma^{-1}(i)} + \sum_{j \in [n_r - 1]} \mathbf{X}_{\sigma^{-1}(j)} & \text{if } i = 1 \\ \mathbf{X}_{\sigma^{-1}(i)} - \mathbf{X}_{\sigma^{-1}(e(i))} & \text{otherwise.} \end{cases}$$

We now show that the change of variable  $\phi^r$  is a  $\mathcal{C}^1$  diffeomorphism by proving that its Jacobian has full rank. Ordering the columns with the map  $\sigma$ , the latter takes the form:

$$\mathbf{J}_{\phi^r} = \begin{pmatrix} 2 & 1 & 1 & \dots & 1 \\ & 1 & 0 & \dots & 0 \\ & & \ddots & \ddots & \vdots \\ & \mathbf{A} & & \ddots & 0 \\ & & & & 1 \end{pmatrix}$$

where  $\mathbf{A}$  is a strictly lower triangular matrix such that for all  $i \in [2 : n_r - 1]$ ,  $A_{ie(i)} = -1$  and for all  $t \neq e(i)$ ,  $A_{it} = 0$ . The above can be factorized as:

$$\mathbf{J}_{\phi^r} = \begin{pmatrix} \alpha_{n_r-1} & \alpha_{n_r-2} & \dots & \alpha_2 & \alpha_1 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & \dots & \dots & 0 & 1 \end{pmatrix}^{-1} \begin{pmatrix} 1 & 0 & \dots & \dots & 0 \\ & 1 & \ddots & & \vdots \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & \vdots \\ & \mathbf{A} & & \ddots & 0 \\ & & & & 1 \end{pmatrix}$$

where  $\alpha_1 = -1$  and for  $\ell > 1$ ,  $\alpha_\ell = \sum_{j < \ell} \alpha_j \mathbb{1}_{e(n_r-j)=n_r-\ell} - 1$ . With this in place, for  $i \in [n_r - 1]$ ,  $\alpha_i \neq 0$  in particular  $\alpha_{n_r-1} \neq 0$  therefore  $|\mathbf{J}_{\phi^r}| \neq 0$  and  $\phi^r$  is a  $\mathcal{C}^1$  diffeomorphism. This change of variable yields:

$$\begin{aligned} \int \prod_{(i,j) \in C_r^2} k(\mathbf{X}_i - \mathbf{X}_j)^{\mathcal{T}_{ij}} \lambda_{\mathcal{S}_C}(d\mathbf{X}) &= \int \bigotimes_{i \in [n_r - 1]} k(\mathbf{Y}_i) |\mathbf{J}_{\phi^r}(\mathbf{Y})|^{-1} \lambda_{\mathbb{R}^p}(d\mathbf{Y}) \\ &= |\mathbf{J}_{\phi^r}|^{-1} \prod_{i \in [n_r - 1]} \int k(\mathbf{Y}_i) \lambda_{\mathbb{R}^p}(d\mathbf{Y}_i) \end{aligned}$$

<sup>1</sup>Symmetrized version *i.e.*  $\overline{\mathcal{T}^r} = \mathcal{T}^r + (\mathcal{T}^r)^\top$ .

using the Fubini Tonelli theorem. The result follows from  $\lambda_{\mathbb{R}^p}$ -integrability of  $k$  and upper bound [10](#).

## A.2 Proof of proposition [1](#)

Let  $\mathcal{P} \in \{B, D, E\}$ ,  $k$  be a valid kernel (assumptions of theorem [1](#)) with  $\mathbf{K}_X = (k(\mathbf{X}_i - \mathbf{X}_j))_{(i,j) \in [n]^2}$  and  $\boldsymbol{\pi} \in \mathbb{R}_+^{n \times n}$ . Let  $\mathbf{W} \sim \mathbb{P}_{\mathcal{P},k}^\varepsilon(\cdot; \boldsymbol{\pi}, 1)$ . Inversion of conditional with Bayes rule gives:

$$\forall \mathbf{W} \in \mathcal{S}_W, \quad \mathbb{P}(\mathbf{W}|\mathbf{X}) \propto \mathcal{C}_k^\varepsilon(\mathbf{W})^{-1} f^\varepsilon(\mathbf{X}, \mathbf{W}) f_k(\mathbf{X}, \mathbf{W}) \mathbb{P}_{\mathcal{P},k}^\varepsilon(\mathbf{W}; \boldsymbol{\pi}, 1) \quad (11)$$

where the prior reads:

$$\mathbb{P}_{\mathcal{P},k}^\varepsilon(\mathbf{W}; \boldsymbol{\pi}, 1) \propto \mathcal{C}_k^\varepsilon(\mathbf{W}) \Omega_{\mathcal{P}}(\mathbf{W}) \prod_{(i,j) \in [n]^2} \pi_{ij}^{W_{ij}}. \quad (12)$$

Hence the joint normalizing constant simplifies such that:

$$\forall \mathbf{W} \in \mathcal{S}_W, \quad \mathbb{P}(\mathbf{W}|\mathbf{X}) \propto f^\varepsilon(\mathbf{X}, \mathbf{W}) \Omega_{\mathcal{P}}(\mathbf{W}) \prod_{(i,j) \in [n]^2} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \quad (13)$$

$$\xrightarrow{\varepsilon \rightarrow 0} \Omega_{\mathcal{P}}(\mathbf{W}) \prod_{(i,j) \in [n]^2} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \quad (14)$$

which ends the proof. As a complement, we now explicit the simple forms taken by the posterior limit graph in each case.

**B-Prior** Recall that in this case the prior reads:

$$\mathbb{P}_B^\varepsilon(\mathbf{W}; \boldsymbol{\pi}, 1) \propto \mathcal{C}_k^\varepsilon(\mathbf{W}) \prod_{(i,j) \in [n]^2} \pi_{ij}^{W_{ij}} \mathbb{1}_{W_{ij} \leq 1}.$$

Therefore the posterior limit graph has the distribution:

$$\begin{aligned} \mathbb{P}_B(\mathbf{W}; \boldsymbol{\pi} \odot \mathbf{K}_X) &= \frac{\prod_{(i,j) \in [n]^2} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \mathbb{1}_{W_{ij} \leq 1}}{\sum_{\mathbf{W} \in \mathcal{S}_W} \prod_{(i,j) \in [n]^2} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \mathbb{1}_{W_{ij} \leq 1}} \\ &= \prod_{(i,j) \in [n]^2} \left( \frac{\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)}{1 + \pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)} \right)^{W_{ij}} \left( \frac{1}{1 + \pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)} \right)^{1 - W_{ij}} \mathbb{1}_{W_{ij} \leq 1}. \end{aligned}$$

This distribution amounts to:  $\forall (i,j) \in [n]^2, \quad \mathbf{W}_{ij} \stackrel{\perp}{\sim} \mathcal{B} \left( \frac{\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)}{1 + \pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)} \right)$ .

**D-Prior** The prior writes:

$$\mathbb{P}_D^\varepsilon(\mathbf{W}; \boldsymbol{\pi}, 1) \propto \mathcal{C}_k^\varepsilon(\mathbf{W}) \prod_{(i,j) \in [n]^2} \pi_{ij}^{W_{ij}} \mathbb{1}_{W_{i+}=1}.$$

The distribution of the posterior limit then becomes:

$$\begin{aligned} \mathbb{P}_D(\mathbf{W}; \boldsymbol{\pi} \odot \mathbf{K}_X) &= \frac{\prod_{(i,j) \in [n]^2} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \mathbb{1}_{W_{i+}=1}}{\sum_{\mathbf{W} \in \mathcal{S}_W} \prod_{(i,j) \in [n]^2} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \mathbb{1}_{W_{i+}=1}} \\ &= \frac{\prod_{(i,j) \in [n]^2} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \mathbb{1}_{W_{i+}=1}}{\prod_{i \in [n]} \sum_{\ell \in [n]} \pi_{i\ell} k(\mathbf{X}_i - \mathbf{X}_\ell)} \\ &= \prod_{(i,j) \in [n]^2} \left( \frac{\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)}{\sum_{\ell \in [n]} \pi_{i\ell} k(\mathbf{X}_i - \mathbf{X}_\ell)} \right)^{W_{ij}} \mathbb{1}_{W_{i+}=1}. \end{aligned}$$

This distribution amounts to:  $\forall i \in [n], \quad \mathbf{W}_i \stackrel{\perp}{\sim} \mathcal{M} \left( 1, \left( \frac{\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)}{\sum_{\ell \in [n]} \pi_{i\ell} k(\mathbf{X}_i - \mathbf{X}_\ell)} \right)_{j \in [n]} \right)$ .

**E-Prior** In this case the prior reads:

$$\mathbb{P}_E^\varepsilon(\mathbf{W}; \boldsymbol{\pi}, 1) \propto \mathcal{C}_k^\varepsilon(\mathbf{W}) \prod_{(i,j) \in [n]^2} \frac{\pi_{ij}^{W_{ij}}}{W_{ij}!} \mathbb{1}_{W_{++}=n}.$$

Finally, deriving the distribution of the posterior graph limit:

$$\begin{aligned} \mathbb{P}_E(\mathbf{W}; \boldsymbol{\pi} \odot \mathbf{K}_X) &= \frac{\prod_{(i,j) \in [n]^2} (W_{ij}!)^{-1} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \mathbb{1}_{W_{++}=n}}{\sum_{\mathbf{W} \in \mathcal{S}_W} \prod_{(i,j) \in [n]^2} (W_{ij}!)^{-1} (\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j))^{W_{ij}} \mathbb{1}_{W_{++}=n}} \\ &= n! \prod_{(i,j) \in [n]^2} (W_{ij})^{-1} \left( \frac{\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)}{\sum_{(\ell,t) \in [n]^2} \pi_{\ell t} k(\mathbf{X}_\ell - \mathbf{X}_t)} \right)^{W_{ij}} \mathbb{1}_{W_{++}=n}. \end{aligned}$$

This distribution amounts to:  $\mathbf{W} \sim \mathcal{M} \left( n, \left( \frac{\pi_{ij} k(\mathbf{X}_i - \mathbf{X}_j)}{\sum_{(\ell,t) \in [n]^2} \pi_{\ell t} k(\mathbf{X}_\ell - \mathbf{X}_t)} \right)_{(i,j) \in [n]^2} \right)$ .

### A.3 Proof of theorem 2

We consider the following hierarchical model, for  $\nu_X, \nu_Z \geq n$ :

$$\begin{aligned} \boldsymbol{\Theta}_X &\sim \mathcal{W}(\nu_X, \mathbf{I}_n) \\ \text{vec}(\mathbf{X}) | \boldsymbol{\Theta}_X &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Theta}_X^{-1} \otimes \mathbf{I}_p) \\ \boldsymbol{\Theta}_Z &\sim \mathcal{W}(\nu_Z, \mathbf{I}_n) \\ \text{vec}(\mathbf{Z}) | \boldsymbol{\Theta}_Z &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Theta}_Z^{-1} \otimes \mathbf{I}_q). \end{aligned}$$

With this at hand, the posteriors for  $\boldsymbol{\Theta}_X$  and  $\boldsymbol{\Theta}_Z$  can be derived in closed form:

$$\begin{aligned} \boldsymbol{\Theta}_X | \mathbf{X} &\sim \mathcal{W}(\nu_X + p, (\mathbf{I}_n + \mathbf{X} \mathbf{X}^\top)^{-1}) \\ \boldsymbol{\Theta}_Z | \mathbf{Z} &\sim \mathcal{W}(\nu_Z + q, (\mathbf{I}_n + \mathbf{Z} \mathbf{Z}^\top)^{-1}). \end{aligned}$$

Keeping terms of  $-\mathbb{E}_{\boldsymbol{\Theta}_X} [\log \mathbb{P}(\boldsymbol{\Theta}_Z = \boldsymbol{\Theta}_X | \mathbf{Z}) | \mathbf{X}]$  that depends on  $\mathbf{Z}$ , one has the optimization problem:

$$\min_{\mathbf{Z} \in \mathbb{R}^{n \times q}} \frac{\nu_X + p}{2} \text{tr}(\mathbf{Z}^\top (\mathbf{I}_n + \mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{Z}) - \frac{\nu_Z + q}{2} \log |\mathbf{I}_n + \mathbf{Z} \mathbf{Z}^\top|$$

Our strategy is to first find the optimal sample covariance matrix  $\mathbf{Z} \mathbf{Z}^\top$  and then focus on the solution in  $\mathbf{Z}$ . To that extent, consider the eigendecomposition of the sample covariance matrices:  $\mathbf{X} \mathbf{X}^\top = \mathbf{V} \mathbf{D} \mathbf{V}^\top$  and  $\mathbf{Z} \mathbf{Z}^\top = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$  where  $\mathbf{D} = \text{diag}(\mathbf{d})$  and  $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\lambda})$  such that  $d_1 \geq \dots \geq d_n$  and  $\lambda_1 \geq \dots \geq \lambda_n$ . Denoting  $\gamma = (\nu_X + q)/(\nu_Z + p)$ , we consider the following problem:

$$\min_{\mathbf{U} \in \mathcal{O}(n), \boldsymbol{\Lambda}} \text{tr}(\mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top \mathbf{V} (\mathbf{I}_n + \mathbf{D})^{-1} \mathbf{V}^\top) - \gamma \log |\mathbf{I}_n + \boldsymbol{\Lambda}| \quad (15)$$

$$\text{s.t. } \boldsymbol{\Lambda} \succcurlyeq \mathbf{0} \quad (16)$$

$$\text{rank}(\boldsymbol{\Lambda}) \leq q \quad (17)$$

The above problem is non-convex because of the rank constraint (17). Nonetheless it can be simplified as we now show.

We focus on finding the optimal eigenvectors first. To that extent, let us denote,  $\mathbf{R} = \mathbf{U}^\top \mathbf{V}$ . Only the left term in (18) depends on  $\mathbf{R}$ . The optimization problem for eigenvectors writes:

$$\min_{\mathbf{R} \in \mathcal{O}(n)} \text{tr}(\mathbf{R}^\top \boldsymbol{\Lambda} \mathbf{R} (\mathbf{I}_n + \mathbf{D})^{-1}) \quad (18)$$

The objective (18) can be expressed as:  $\sum_{(i,j) \in [n]^2} \lambda_i (1 + d_j)^{-1} R_{ij}^2$ . Now one can notice that since  $\mathbf{R}$  is orthogonal,  $\mathbf{R} \odot \mathbf{R}$  is doubly stochastic (*i.e.* sum of coefficients on each row and

column is equal to one). Therefore thanks to the Birkhoff–von Neumann theorem, there exists  $\theta_1, \dots, \theta_L \geq 0$ ,  $\sum_{\ell \in [L]} \theta_\ell = 1$  and permutation matrices  $\mathbf{P}_1, \dots, \mathbf{P}_L$  such that:

$$\mathbf{R} \odot \mathbf{R} = \sum_{\ell \in [L]} \theta_\ell \mathbf{P}_\ell$$

where for all  $\ell \in [L]$ , there exists a permutation  $\sigma_\ell$  of  $[n]$  such that  $P_{\ell,ij} = \mathbb{1}_{\sigma_\ell(i)=j}$  for  $(i, j) \in [n]^2$ .

With this at hand, objective (18) writes:  $\sum_{\ell \in [L]} \theta_\ell \sum_{i \in [n]} \lambda_i (1 + d_{\sigma_\ell(i)})^{-1}$ . There exists a permutation  $\sigma^*$  such that the quantity  $\sum_{i \in [n]} \lambda_i (1 + d_{\sigma_\ell(i)})^{-1}$  is minimal. Note that the identity permutation *i.e.* for  $i \in [n]$ ,  $\sigma(i) = i$  is optimal in this case as the  $(\lambda_i)_{i \in [n]}$  and the  $(d_i)_{i \in [n]}$  are in decreasing order. Then choosing for  $\ell \in [L]$ ,  $\theta_\ell = \sigma_\ell = \sigma^*$  minimizes the latter quantity. Therefore the solution of (18)  $\mathbf{R}^*$  is such that for  $(i, j) \in [n]^2$ ,  $R_{ij}^* = \pm \mathbb{1}_{\sigma^*(i)=j}$ . Thus an optimum in  $\mathbf{U}$  of (18) is such that  $\mathbf{U}^* = \mathbf{V} \mathbf{R}^*$ .

Hence  $\mathbf{U} = \mathbf{V}$ , in particular, is optimal. We will choose this  $\mathbf{U}$  in what follows as the sign of the axes do not influence the characterization of the final result in  $\mathbf{Z}$  as a PCA embedding. Such a choice gives  $\mathbf{Z} \mathbf{Z}^\top = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$ .

Now it remains to find the optimal eigenvalues  $(\lambda_i)_{i \in [n]}$ . The rank constraint (17) can be easily dealt with: since the eigenvalues are sorted in decreasing order, the constraint implies that for  $i \geq q$ ,  $\lambda_i = 0$ . Thus the eigenvalue problem can be formulated in  $\mathbb{R}^q$ :

$$\min_{\boldsymbol{\lambda} \in \mathbb{R}^q} \boldsymbol{\lambda}^\top (\mathbf{1} + \mathbf{d})^{-1} - \gamma \mathbf{1}^\top \log(\mathbf{1} + \boldsymbol{\lambda}) \quad (19)$$

$$\text{s.t. } \forall i \in [q], \quad \lambda_i \geq 0, \quad \lambda_1 \geq \dots \geq \lambda_q \quad (20)$$

where (20) accounts for (16). The above is convex. (19) is minimized for  $\boldsymbol{\lambda} = \gamma(\mathbf{1} + \mathbf{d}) - \mathbf{1}$ . Taking the feasibility constraint (20) into account one has a solution  $\boldsymbol{\lambda}^*$  such that:

$$\forall i \in [n], \quad \lambda_i^* = \begin{cases} \max(0, \gamma(1 + d_i) - 1) & \text{if } i \leq q \\ 0 & \text{otherwise.} \end{cases}$$

Note that this solution is not unique if there are repeated eigenvalues. Notice also that one has the freedom to choose the Wishart prior parameters such that  $\gamma = 1$ . Doing so, the solution satisfies  $\mathbf{Z}^* \mathbf{Z}^{*\top} = \mathbf{V}_{[:,q]} \mathbf{D}_{[q,q]} \mathbf{V}_{[q,:]}^\top$ . Therefore there exists  $\mathbf{R}$  an orthogonal matrix of size  $q$  such that  $\mathbf{Z}^* = \mathbf{V}_{[:,q]} \mathbf{D}_{[q,q]}^{\frac{1}{2}} \mathbf{R}$ . The latter is the output of a PCA model of  $\mathbf{X}$  with  $q$  components, which is defined up to a rotation.

#### A.4 Proof of Corollary 1

With the presented hierarchical model (fig. 3), the coupling problem is the following:

$$\min_{\mathbf{Z} \in \mathcal{S}_M^q} \text{tr}(\mathbf{U}_{[:,R]} \mathbf{Z}^\top (\mathbf{I}_R + \varepsilon \mathbf{U}_{[:,R]}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}_{[:,R]})^{-1} \mathbf{U}_{[:,R]} \mathbf{Z}) - \log |\mathbf{I}_R + \varepsilon \mathbf{U}_{[:,R]}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{U}_{[:,R]}| \quad (21)$$

where  $\mathbf{U}_{[:,R]}$  are the eigenvectors associated to the Laplacian null-space of  $\overline{\mathbf{W}}_X$ .

Let us denote  $\bar{\mathbf{Z}} = \mathbf{U}_{[:,R]}^\top \mathbf{Z} \in \mathbb{R}^{R \times q}$  and  $\bar{\mathbf{X}} = \mathbf{U}_{[:,R]}^\top \mathbf{X} \in \mathbb{R}^{R \times p}$ . Note that  $\mathbf{Z} \rightarrow \mathbf{U}_{[:,R]}^\top \mathbf{Z}$  is a bijective linear map from  $\mathcal{S}_M^q$  to  $\mathbb{R}^{R \times q}$  with inverse  $\bar{\mathbf{Z}} \rightarrow \mathbf{U}_{[:,R]} \bar{\mathbf{Z}}$  (and equivalently for  $\mathbb{R}^{R \times p}$ ). Hence (21) is equivalent to:

$$\min_{\bar{\mathbf{Z}} \in \mathbb{R}^{R \times q}} \text{tr}(\bar{\mathbf{Z}}^\top (\mathbf{I}_R + \varepsilon \bar{\mathbf{X}} \bar{\mathbf{X}}^\top)^{-1} \bar{\mathbf{Z}}) - \log |\mathbf{I}_R + \varepsilon \bar{\mathbf{Z}} \bar{\mathbf{Z}}^\top| \quad (22)$$

According to theorem 2, the solution of problem (22) is such that there exists  $\mathbf{R}$  orthogonal,  $\bar{\mathbf{Z}}^* = \mathbf{V}_{[:,q]} \mathbf{S}_{[q,q]} \mathbf{R}$  where  $\bar{\mathbf{X}} \bar{\mathbf{X}}^\top = \mathbf{V} \mathbf{S}^2 \mathbf{V}^\top$  is the eigendecomposition in an orthogonal basis of the among-row covariance matrix of  $\bar{\mathbf{X}}$ . Note that the solution does not depend on  $\varepsilon$ .

Therefore (21) is solved for  $\mathbf{Z}^* = \mathbf{U}_{[:,R]} \mathbf{V}_{[:,q]} \mathbf{S}_{[q,q]} \mathbf{R}$ . One can notice that since the singular value decomposition (*i.e.* SVD) of  $\mathbf{U}_{[:,R]}^\top \mathbf{X}$  takes the form  $\mathbf{V} \mathbf{S} \mathbf{B}$  where  $\mathbf{B}$  is a semi-orthogonal matrix of size  $p$ , then  $\mathbf{U}_{[:,R]} \mathbf{U}_{[:,R]}^\top \mathbf{X} = \mathbf{U}_{[:,R]} \mathbf{V} \mathbf{S} \mathbf{B}$ . Noticing that  $\mathbf{V}' = \mathbf{U}_{[:,R]} \mathbf{V}$  is orthogonal, one has that  $\mathbf{V}' \mathbf{S} \mathbf{B}$  is a compact SVD of  $\mathbf{U}_{[:,R]} \mathbf{U}_{[:,R]}^\top \mathbf{X}$ . Therefore, since  $\mathbf{Z}^* = \mathbf{V}' \mathbf{S}$ ,  $\mathbf{Z}^*$  is a PCA embedding of  $\mathbf{U}_{[:,R]} \mathbf{U}_{[:,R]}^\top \mathbf{X}$ .

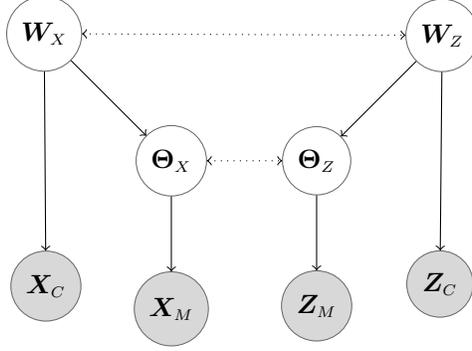


Figure 3: Graphical representation of the hierarchical model considered in section 4.2. Plain directed arrows represent conditional dependencies while dotted arrows represent the coupling links. Corollary 1 provides a solution for the coupling between  $\Theta_X$  and  $\Theta_Z$ .

## B Experiments Supplementary Material

### B.1 Experimental Setup and Details About *ccPCA*

**Implementation of existing methods.** For t-SNE, we rely on the openTSNE implementation [34] for both computing the kernel  $\mathbf{K}_X$  with appropriate bandwidths and running the tSNE algorithm. We keep the training default parameters and 1000 iterations of gradient descent. For all experiments, the default perplexity of 30 was used to set the kernel bandwidths. For UMAP, we use the default Python implementation of [30] with default parameters. For PCA and Laplacian eigenmaps, the scikit-learn implementation is used [32] with default parameters as well.

***ccPCA*.** The pseudo code of the algorithm is given in algorithm 1. CCs’ memberships (*i.e.* eigenvectors  $\mathbf{U}_{[R]}$ ) are computed using igraph [13]. Regarded the time complexity of *ccPCA*, one can sample the posterior graph with constant time if  $\mathcal{P}_X = E$ , linear time if  $\mathcal{P}_X = D$  and quadratic time if  $\mathcal{P}_X = B$ . Moreover, computing  $\mathbf{U}_{[R]}$  can be done with linear complexity *w.r.t.* the number of nodes. Hence the time complexity is  $O(N \times n)$  for  $E$  and  $D$  priors and  $O(N \times n^2)$  for the  $B$  prior, where  $N$  is the number of Monte Carlo samples. In practice we found that  $N \approx 100$  Monte Carlo samples produce a consistent *ccPCA* embedding for  $n \approx 10000$ . Note that the time complexity of PCA is  $O(\min(p^3, n^3))$  where  $p$  is the dimensionality (*i.e.* number of columns) of  $\mathbf{X}$ . Hence in most common applications involving images or biological sequencing data (where  $p$  is very large), the additional time complexity brought by *ccPCA* compared to PCA is negligible.

---

#### Algorithm 1 *ccPCA*

---

**Input:**  $\mathbf{K}_X, \mathcal{P}_X, N$   
**for**  $\ell = 1$  **to**  $N$  **do**  
    Sample  $\mathbf{W}^\ell \sim \mathbb{P}_{\mathcal{P}_X}(\cdot; \mathbf{K}_X)$   
    Compute CCs’ memberships  $\mathbf{U}_{[R]}^\ell$  of  $\mathbf{W}^\ell$   
**end for**  
**Output:** PCA of  $\left(N^{-1} \sum_{\ell \in [N]} \mathbf{U}_{[R]}^\ell \mathbf{U}_{[R]}^{\ell T}\right) \mathbf{X}$

---

All experiments are performed on a machine with four Intel Core i5 processors and 16 GB memory.

### B.2 *ccPCA* with Varying Perplexity Values

Recall that the *ccPCA* algorithm retrieves the same latent graph as neighbor embedding methods. As shown in section 3.1 these graphs’ distributions depend on the type of prior considered, and take simple forms as follows, when  $\boldsymbol{\pi}_X = \mathbf{1}$  :

- if  $\mathcal{P} = B$ ,  $\forall (i, j) \in [n]^2$ ,  $W_{ij} \stackrel{\perp}{\sim} \mathcal{B}(K_{X,ij}/(1 + K_{X,ij}))$
- if  $\mathcal{P} = D$ ,  $\forall i \in [n]$ ,  $\mathbf{W}_i \stackrel{\perp}{\sim} \mathcal{M}(1, \mathbf{K}_{X,i}/K_{X,i+})$
- if  $\mathcal{P} = E$ ,  $\mathbf{W} \sim \mathcal{M}(n, \mathbf{K}_X/K_{X,++})$

and  $\mathbf{K}_X$  is the kernel matrix evaluated on the data such that:

$$\forall (i, j) \in [n]^2, \quad \mathbf{K}_{X,ij} = k((\mathbf{X}_i - \mathbf{X}_j)/\tau_i)$$

where  $\tau \in \mathbb{R}^n$  is set using an heuristic depending on the method considered [38, 30, 36]. In fig. 4, we focus on the effect of the kernel bandwidths on *ccPCA*, choosing the example of t-SNE.

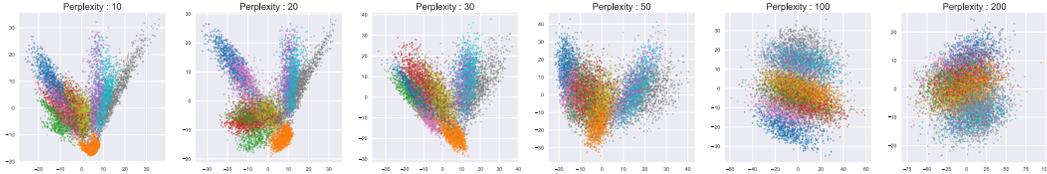


Figure 4: *ccPCA* launched for different values of the perplexity parameter. The latter determines the kernel bandwidths and can be interpreted as the number of effective neighbors of each point [38]. As the perplexity grows, the probability of connecting different clusters of digit by sampling through the graph posterior  $\mathbb{P}_{\mathcal{P}_X}(\cdot; \mathbf{K}_X)$  increases. Therefore clusters are less and less identifiable as the perplexity increases.

From fig. 4, one can notably notice that using a high perplexity leads to a more connected graph and therefore a PCA-like embedding with less degeneracy and no clustering effect. Recall that *ccPCA* computes the same clusters as t-SNE through the CCs of the latent MRF and manage to position t-SNE clusters by focusing on their relative positions (that are filtered by t-SNE). In the case of a connected graph (high perplexity), *ccPCA* will show little advantage over classical PCA since there will not be any cluster to position. Note that this discussion can be extended to other neighbor embedding methods equivalently. Therefore, our probabilistic framework allows us to indentify which part of information is filtered by the posterior graphs with given kernel bandwidth.

### B.3 Quantitative Evaluation of *ccPCA*

For quantitative assessment of *ccPCA*, we focused on t-SNE [38] and UMAP [30] which are the most popular neighbor embedding methods. Note that for these algorithms the initialization is crucial for the global structure of the embeddings as shown in [21]. In addition to MNIST [14], we considered the datasets cifar-10, cifar-100 [22], fashion-MNIST [42] as well as the CD8+ T lymphocytes single cell RNA-seq dataset from [24].

We used the quantitative criterion of [25] to assess the quality of the embeddings. As mentioned in this paper, the use of this criterion appears as the general consensus in dimension reduction, a field in which building meaningful criteria is tedious. The criterion measures the rescaled average agreement between the K-ary neighbourhoods in the input and output spaces. It is constructed as follows.

We first define the following quantities for  $(i, j) \in [n]^2$ ,  $\rho_{ij} = |\{k : \|\mathbf{X}_i - \mathbf{X}_k\|_2^2 < \|\mathbf{X}_i - \mathbf{X}_j\|_2^2\}|$ ,  $r_{ij} = |\{k : \|\mathbf{Z}_i - \mathbf{Z}_k\|_2^2 < \|\mathbf{Z}_i - \mathbf{Z}_j\|_2^2\}|$ ,  $\nu_i^K = \{j : 1 \leq \rho_{ij} \leq K\}$  and  $\gamma_i^K = \{j : 1 \leq r_{ij} \leq K\}$ . The average K-ary neighbourhood preservation is rescaled to indicate the improvement over a random embedding such that:

$$R_n(K) = \frac{(n-1)Q_n(K) - K}{n-1-K} \quad (23)$$

where  $Q_n(K) = \frac{1}{Kn} \sum_{i=1}^n |\nu_i^K \cap \gamma_i^K|$ ,  $n$  is the number of data points and  $K$  is the hyperparameter that adjusts the scale at which we are looking.

To focus on large-scale structure,  $K$  was chosen as either  $n/4$  or  $n/2$ . As summarized by [21], current practice consists in using PCA or Laplacian eigenmaps as initialization for these

algorithms, thus we compare to these strategies. Results are displayed in table 2 and table 3, each entry being an average over 5 random seeds, with standard deviation displayed below each entry. Note that when not specified, tSNE and UMAP are initialized with an isotropic Gaussian variable.

These results show that using *ccPCA* is a reliable alternative to PCA and Laplacian eigenmaps for reproducing large-scale neighborhoods.

Table 2:  $100 \times R_n(K)$  (23) for embeddings produced using t-SNE with various initializations.

K of K-ary	tSNE		PCA + tSNE		LE + tSNE		ccPCA + tSNE	
	n/4	n/2	n/4	n/2	n/4	n/2	n/4	n/2
MNIST	18.7 $\pm 2.2$	7.4 $\pm 5.1$	28.4 $\pm 0.3$	21.9 $\pm 0.2$	26.7 $\pm 0.7$	18.5 $\pm 0.4$	<b>31.3</b> $\pm 0.4$	<b>28.5</b> $\pm 1.2$
cifar-10	20.3 $\pm 3.2$	16.4 $\pm 4.8$	<b>36.9</b> $\pm 0.6$	41.9 $\pm 1.1$	25.8 $\pm 0.6$	24.1 $\pm 1.5$	36.4 $\pm 0.4$	<b>43.4</b> $\pm 1.6$
cifar-100	21.6 $\pm 3.6$	18.2 $\pm 5.5$	38.1 $\pm 0.4$	<b>47.5</b> $\pm 0.4$	23.3 $\pm 1.5$	26.5 $\pm 1.8$	<b>39.6</b> $\pm 0.7$	43.6 $\pm 1.1$
fashion-MNIST	27.2 $\pm 4.3$	12.3 $\pm 7.8$	36.9 $\pm 0.1$	28.5 $\pm 0.2$	32.0 $\pm 0.8$	25.1 $\pm 2.2$	<b>41.6</b> $\pm 0.9$	<b>35.7</b> $\pm 1.5$
Single Cell data	25.7 $\pm 4.8$	22.4 $\pm 10.6$	37.7 $\pm 2.7$	29.0 $\pm 4.7$	28.1 $\pm 1.5$	31.5 $\pm 1.4$	<b>40.1</b> $\pm 1.7$	<b>34.6</b> $\pm 2.6$

Table 3:  $100 \times R_n(K)$  (23) for embeddings produced using UMAP with various initializations.

K of K-ary	UMAP		PCA + UMAP		LE + UMAP		ccPCA + UMAP	
	n/4	n/2	n/4	n/2	n/4	n/2	n/4	n/2
MNIST	29.5 $\pm 1.4$	22.7 $\pm 2.2$	<b>36.6</b> $\pm 0.2$	31.1 $\pm 0.5$	34.6 $\pm 0.2$	24.9 $\pm 0.7$	33.4 $\pm 0.3$	<b>32.3</b> $\pm 0.5$
cifar-10	39.2 $\pm 2.6$	47.6 $\pm 1.1$	44.3 $\pm 0.2$	<b>53.4</b> $\pm 0.1$	44.2 $\pm 0.1$	52.6 $\pm 0.2$	<b>44.6</b> $\pm 0.2$	53.2 $\pm 0.2$
cifar-100	41.6 $\pm 1.8$	42.2 $\pm 0.9$	45.4 $\pm 0.2$	45.2 $\pm 0.1$	44.2 $\pm 0.3$	43.4 $\pm 0.1$	<b>49.9</b> $\pm 0.4$	<b>52.9</b> $\pm 0.6$
fashion-MNIST	48.7 $\pm 2.6$	33.6 $\pm 9.5$	56.2 $\pm 0.5$	54.3 $\pm 0.6$	58.1 $\pm 0.5$	53.4 $\pm 0.6$	<b>58.9</b> $\pm 0.5$	<b>55.8</b> $\pm 0.3$
Single Cell data	39.5 $\pm 1.4$	34.3 $\pm 6.1$	52.3 $\pm 0.8$	47.2 $\pm 6.9$	<b>55.9</b> $\pm 0.3$	45.7 $\pm 0.9$	53.6 $\pm 0.3$	<b>53.9</b> $\pm 1.3$