

## Supplementary Material

### The sufficiently scattered condition

In some prior works (Huang et al., 2014; 2016; 2018), the sufficiently scattered condition is presented as follows:

**Definition 2** (sufficiently scattered [dual cone form]). Let  $\text{cone}(\mathbf{M})^*$  denote the polyhedral cone  $\{\mathbf{x} : \mathbf{M}^\top \mathbf{x} \geq 0\}$  and  $\mathcal{K}^*$  denote the elliptical cone  $\{\mathbf{x} : \|\mathbf{x}\| \leq \mathbf{1}^\top \mathbf{x}\}$ . Matrix  $\mathbf{M}$  is sufficiently scattered if

- i)  $\text{cone}(\mathbf{M})^* \subseteq \mathcal{K}^*$ ,
- ii)  $\text{cone}(\mathbf{M})^* \cap \text{bd}\mathcal{K}^* = \{\alpha \mathbf{e}_\ell : \alpha \geq 0, \ell = 1, \dots, k\}$ .

As for our definition, if we drop the  $\mathbf{1}^\top \mathbf{x} = 1$  constraint, it can be equivalently written in the following cone form:

**Definition 3** (sufficiently scattered [primal cone form]). Let  $\text{cone}(\mathbf{M})$  denote the polyhedral cone  $\{\mathbf{M}\boldsymbol{\beta} : \boldsymbol{\beta} \geq 0\}$  and  $\mathcal{K}$  denote the elliptical cone  $\{\mathbf{x} \in \mathbb{R}^k : \|\mathbf{x}\| \leq \frac{1}{\sqrt{k-1}} \mathbf{1}^\top \mathbf{x}\}$ . Matrix  $\mathbf{M}$  is sufficiently scattered if

- i)  $\text{cone}(\mathbf{M}) \supseteq \mathcal{K}$ ,
- ii)  $\text{bd}\text{cone}(\mathbf{M}) \cap \mathcal{K} = \{\alpha(\mathbf{1} - \mathbf{e}_\ell) : \alpha \geq 0, \ell = 1, \dots, k\}$ .

To understand the equivalence of Definitions 2 and 3, it is essential to invoke the concept of the dual cone, denoted with a superscript \*

$$\mathcal{C}^* = \{\mathbf{x} : \mathbf{y}^\top \mathbf{x} \geq 0, \forall \mathbf{y} \in \mathcal{C}\}.$$

One can verify that  $\text{cone}(\mathbf{M})$  and  $\text{cone}(\mathbf{M})^*$  are indeed dual to each other, and so are  $\mathcal{K}$  and  $\mathcal{K}^*$ .

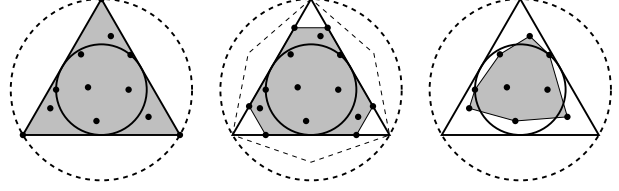
We now prove that Definitions 2 and 3 are indeed primal-dual representations of the same condition. For two convex cones  $\mathcal{C}_1$  and  $\mathcal{C}_2$ , if  $\mathcal{C}_1 \subseteq \mathcal{C}_2$ , then  $\mathcal{C}_1^* \supseteq \mathcal{C}_2^*$ . It clearly shows the equivalence of the first requirements in Definitions 2 and 3. As for the second requirement, we claim the following:

- Requirement ii) in Definition 2 asks that all extreme rays of  $\text{cone}(\mathbf{M})^*$  lie strictly inside  $\mathcal{K}^*$ , except for the coordinate directions  $\alpha \mathbf{e}_\ell$ 's, which lie on the boundary of  $\mathcal{K}^*$ ;
- Requirement ii) in Definition 3 asks that all facets of  $\text{cone}(\mathbf{M})$  lie strictly outside  $\mathcal{K}$ , except for the facets spanned by  $k-1$  coordinate vectors, which touch the boundary of  $\mathcal{K}$  at  $\alpha(\mathbf{1} - \mathbf{e}_\ell)$ .

There is a one-to-one correspondence between extreme rays of  $\text{cone}(\mathbf{M})^*$  and facets of  $\text{cone}(\mathbf{M})$ , which are both defined by  $k-1$  columns of  $\mathbf{M}$ . Let  $\mathbf{v}$  denote an extreme ray of  $\text{cone}(\mathbf{M})^*$ , then  $\mathbf{M}^\top \mathbf{v} \geq 0$  and  $k-1$  of them holds as equalities; those  $k-1$  columns defines a facet of  $\text{cone}(\mathbf{M})$ , and any point  $\mathbf{x}$  in that facet satisfies that  $\mathbf{x}^\top \mathbf{v} = 0$ .

Now if  $\mathbf{v}$  is not a coordinate direction, Definition 2 asks that  $\|\mathbf{v}\| < \mathbf{1}^\top \mathbf{v}$ . It can be rearranged as

$$\frac{\mathbf{1}^\top \mathbf{v}}{\|\mathbf{1}\| \|\mathbf{v}\|} > \frac{1}{\sqrt{k}},$$



(a) Pure node (b) Sufficiently scattered (c) Not identifiable

Figure 5. Same illustration as in Figure 3 with dual cones added (in dash).

which means the angle between  $\mathbf{v}$  and  $\mathbf{1}$  is less than  $\arccos(1/\sqrt{k})$ . Therefore the angle between  $\mathbf{1}$  and any point on the corresponding facet is greater than  $\pi - \arccos(1/\sqrt{k})$ , or equivalently

$$\frac{\mathbf{1}^\top \mathbf{x}}{\|\mathbf{1}\| \|\mathbf{x}\|} < \sqrt{\frac{k-1}{k}} \iff \|\mathbf{x}\| > \frac{1}{\sqrt{k-1}} \mathbf{1}^\top \mathbf{x},$$

meaning all points on that facet lie strictly outside  $\mathcal{K}$ .

The other direction is true as well: If the smallest angle between  $\mathbf{1}$  and any point on the facet is greater than  $\pi - \arccos(1/\sqrt{k})$ , then the angle between  $\mathbf{1}$  and  $\mathbf{v}$  is less than  $\arccos(1/\sqrt{k})$ , meaning  $\mathbf{v}$  lie strictly inside  $\mathcal{K}^*$ . This shows the equivalence of the second requirement in Definitions 2 and 3. A geometric illustration with the dual cones shown is given in Figure 5.

### Proof of Theorem 1

Denote an optimal solution of (8) as  $(\mathbf{E}_\star, \mathbf{M}_\star)$ , then clearly

$$|\det \mathbf{E}_\star| \leq |\det \tilde{\mathbf{E}}^h| \iff |\det \mathbf{E}_\star^{-1} \tilde{\mathbf{E}}^h| \geq 1. \quad (13)$$

Furthermore, since both  $(\mathbf{E}_\star, \mathbf{M}_\star)$  and  $(\tilde{\mathbf{E}}^h, \mathbf{M}_2^h)$  are feasible, we have

$$\begin{aligned} \mathbf{E}_\star \mathbf{M}_\star &= \tilde{\mathbf{E}}^h \mathbf{M}_2^h, \mathbf{M}_\star \geq 0 \\ \implies \mathbf{E}_\star^{-1} \tilde{\mathbf{E}}^h \mathbf{M}_2^h &\geq 0 \end{aligned} \quad (14)$$

$$\begin{aligned} \mathbf{e}_k^\top \mathbf{E}_\star &= \mathbf{1}^\top, \mathbf{e}_k^\top \tilde{\mathbf{E}}^h = \mathbf{1}^\top, \\ \implies \mathbf{1}^\top \mathbf{E}_\star^{-1} \tilde{\mathbf{E}}^h &= \mathbf{1}^\top. \end{aligned} \quad (15)$$

Denote  $\mathbf{V} = \mathbf{E}_\star^{-1} \tilde{\mathbf{E}}^h$  with  $\ell$ -th row denoted as  $\mathbf{v}_\ell^\top$ . Then (14) means  $\mathbf{v}_\ell$ 's lie in  $\text{cone}(\mathbf{M}_2^h)$ ; since  $\mathbf{M}_2^h$  is sufficiently scattered,  $\mathbf{v}_\ell \in \mathcal{K}^*$ , meaning

$$\|\mathbf{v}_\ell\| \leq \mathbf{1}^\top \mathbf{v}_\ell.$$

Using Hadamard's inequality,

$$|\det \mathbf{B}| \leq \prod_{\ell=1}^k \|\mathbf{v}_\ell\|,$$

and we further have that

$$\prod_{\ell=1}^k \|\mathbf{v}_\ell\| \leq \prod_{\ell=1}^k \mathbf{1}^\top \mathbf{v}_\ell \quad (16a)$$

$$\leq \left( \frac{\sum_{\ell=1}^k \mathbf{1}^\top \mathbf{v}_\ell}{k} \right)^k \quad (16b)$$

$$= \left( \frac{\mathbf{1}^\top \mathbf{V} \mathbf{1}}{k} \right)^k = 1, \quad (16c)$$

where (16a) stems from  $\mathbf{v}_\ell \in \mathcal{K}^*$ , (16b) is due to arithmetic-geometric mean inequality, and (16c) is due to (15). This means

$$|\det \mathbf{\Xi}_\star^{-1} \tilde{\mathbf{\Xi}}^{\natural}| \leq 1. \quad (17)$$

Combining (17) and (13), we conclude that

$$|\det \mathbf{\Xi}_\star^{-1} \tilde{\mathbf{\Xi}}^{\natural}| = 1.$$

Furthermore, all inequalities in (16) hold as equalities, which means  $\mathbf{v}_\ell$ 's lie on the boundary of  $\mathcal{K}^*$ . The second requirement of sufficiently scattered (dual form Definition 2) imply that  $\mathbf{v}_\ell$ 's can only take coordinate vectors, meaning  $V$  is a permutation matrix  $\mathbf{H}$ , and we have

$$\mathbf{M}_2^{\natural} = \mathbf{H} \mathbf{M}_\star, \quad \tilde{\mathbf{\Xi}}^{\natural} = \mathbf{\Xi}_\star \mathbf{H}^\top.$$

**Q.E.D.**

### Proof of Theorem 2

Define a vector  $\tilde{\mathbf{f}} \in \mathbb{R}^k$  for a specific  $X$  with the  $m$ -th element equal to

$$\tilde{f}_m = (-1)^{\ell+m} \det X_{\ell m},$$

then the co-factor expansion tells us

$$\det X = \mathbf{x}_\ell^\top \tilde{\mathbf{f}},$$

where  $\mathbf{x}_\ell^\top$  is the  $\ell$ -th row of  $X$ . Therefore, at a particular point  $X$ , and we look at problem (10) at the  $\ell$ -th row of  $X$ , the subproblem is

$$\begin{aligned} & \underset{\mathbf{z}}{\text{maximize}} \quad (\tilde{\mathbf{f}}^\top \mathbf{z})^2 \\ & \text{subject to} \quad \mathbf{z}^\top \tilde{\mathbf{Y}} \geq 0, \mathbf{z}^\top \tilde{\mathbf{Y}} \mathbf{1} = 1. \end{aligned}$$

The objective is a convex quadratic, so if we take a linear approximation at  $\mathbf{x}_\ell$  it defines a global lowerbound to the objective

$$(\mathbf{x}_\ell^\top \tilde{\mathbf{f}}) \tilde{\mathbf{f}}^\top \mathbf{z} = (\det X) \tilde{\mathbf{f}}^\top \mathbf{z},$$

with equality holds when  $\mathbf{z} = \mathbf{x}_\ell$ . This is where the linear programming sub-problem (12) comes from.

Cyclically solving (12) with respect to each row of  $X$  falls into the framework of BSUM, proposed by Razaviyayn et al. (2013). Specifically, we have that the constraint set decouples over the rows of  $X$ , and the objective is a tight lowerbound  $(\det X)^2$  when restricted to the  $\ell$ -th row of  $X$ . Because  $(\det X)$  is a smooth function, it automatically satisfies that the directional derivative of  $\mathbf{f}^\top \mathbf{z}$  is equal to that of  $(\det X)^2$  at everywhere. Furthermore, it is easy to see that the constraint set is a compact set. According to Razaviyayn et al. (2013), the proposed iterative algorithm is guaranteed to converge to a stationary point as long as each of the LP sub-problems (12) has a unique solution.

Regarding the uniqueness of  $\arg \max(12)$ , we note that for any LP, the constraint set defines a polyhedron, and a solution always exists on a vertex. A solution is not unique if the objective direction  $\mathbf{f}$  happens to be normal to one of the edges of the constraint polyhedron. This means there exists a set of  $k-1$  columns in  $\tilde{\mathbf{Y}}$  such that  $\mathbf{f}$  lies in their range. As long as columns of  $\tilde{\mathbf{Y}}$  appears somewhat incoherent, this will happen with a very small probability. **Q.E.D.**

### Proof of Theorem 3

Suppose  $\tilde{\mathbf{Y}} = \tilde{\mathbf{\Xi}} \mathbf{M}_2$  where  $\mathbf{M}_2$  satisfies the separability assumption. This means all the coordinate vectors  $\mathbf{e}_1, \dots, \mathbf{e}_k$  exists in its columns. This means columns of  $\tilde{\mathbf{\Xi}}$  exists in columns of  $\tilde{\mathbf{Y}}$ , which corresponds to the coordinate vectors in the columns of  $\mathbf{M}_2$ . Then  $\mathbf{z}^\top \tilde{\mathbf{\Xi}} \geq 0$  implies  $\mathbf{z}^\top \tilde{\mathbf{Y}} \geq 0$ , since  $\mathbf{M}_2 \geq 0$ , which means  $\mathbf{z}^\top \tilde{\mathbf{Y}} \geq 0$  consists of a lot of redundant constraints, and (12) is equivalent to

$$\begin{aligned} & \underset{\mathbf{z}}{\text{maximize}} \quad \mathbf{f}^\top \mathbf{z} \\ & \text{subject to} \quad \mathbf{z}^\top \tilde{\mathbf{\Xi}} \geq 0, \mathbf{z}^\top \mathbf{b} = 1. \end{aligned}$$

Let us denote the dual variable with respect to the inequality constraint as  $\boldsymbol{\mu}$ , and the equality constraint as  $\lambda$ , then the KKT condition implies

$$\lambda \mathbf{b} - \mathbf{f} = \tilde{\mathbf{\Xi}} \boldsymbol{\mu}.$$

Because we assume  $\tilde{\mathbf{\Xi}}$  is non-singular, it is equivalent to

$$\lambda \tilde{\mathbf{\Xi}}^{-1} \mathbf{b} - \tilde{\mathbf{\Xi}}^{-1} \mathbf{f} = \boldsymbol{\mu} \geq 0.$$

Therefore,  $\lambda$  should be chosen to make  $\lambda \tilde{\mathbf{\Xi}}^{-1} \mathbf{b} - \tilde{\mathbf{\Xi}}^{-1} \mathbf{f} \geq 0$ , and for a scalar  $\lambda$ , the resulting  $\boldsymbol{\mu}$  will have only one zero almost surely. According to complementary slackness, then  $\mathbf{z}^\top \tilde{\mathbf{\Xi}}$  will have only one nonzero entry. This implies that the solution to (12) should be a row of  $\tilde{\mathbf{\Xi}}^{-1}$ , which is exactly what we want.

Now after running one iteration of CD-MVSI, we would obtain  $k$  rows of  $\tilde{\mathbf{\Xi}}$ . If all  $k$  rows of them are present, then we have successfully recovered the ground truth. The only non-ideal case is if one row of  $\tilde{\mathbf{\Xi}}$  appear multiple times. However,

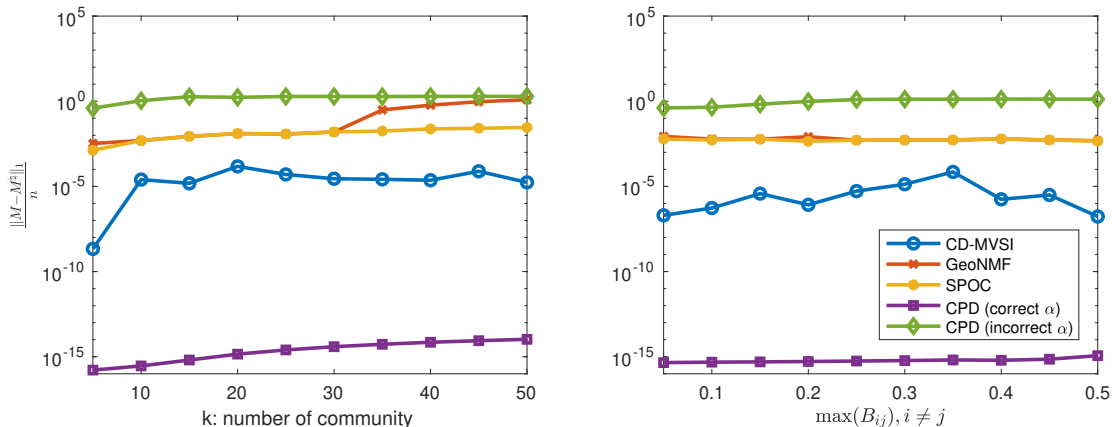


Figure 6. Algorithm performance on synthetic data. Note that the performance of CPD heavily relies on the knowledge of the underlying Dirichlet parameter  $\alpha$ , which may not be practical.

if that is the case, the resulting  $\det \mathbf{X} = 0$ . Suppose we initialize with a non-singular matrix, say  $\mathbf{X} = \text{Diag}(\mathbf{b})^{-1}$ , which is feasible, CD-MVSI is guaranteed to monotonically increase the objective value. Therefore  $\det \mathbf{X}$  is impossible to be equal to zero, and we indeed manage to recover all  $k$  rows of  $\hat{\mathbf{E}}^{-1}$ . **Q.E.D.**

### Additional synthetic experiments

Here we validate the correctness of various algorithms, *given exact statistics ideal for each method*. We start by generating  $k$ -dimensional membership coefficients from a Dirichlet distribution with parameter  $(1/k)\mathbf{1}$ . Since the Dirichlet parameter for all the components are less than one, it tends to generate points that lie on the boundary of  $\Delta$ . However, we note that none of the components will be exactly equal to one, which means the sufficiently scattered condition will not be satisfied exactly. Nevertheless, since the points are well spread out, the recovery result is very close to optimal. On the other hand, the separability / pure-node will be grossly violated as  $k$  goes larger, and we will see that pure-node-based methods are very vulnerable to such assumption violation.

We fix the number of nodes to be 1000. For GeoNMF and SPOC, the entire underlying matrix  $\mathbf{M}^T \mathbf{B} \mathbf{M}$  is given. For the tensor method, the nodes are divided into three groups of size 300, 300, and 400, and the corresponding  $300 \times 300 \times 400$  moment tensor is given. For the tensor method we consider two scenarios: one with the correct Dirichlet parameter  $(1/k)\mathbf{1}$  and one with an incorrect Dirichlet parameter  $\mathbf{1}$ . The goal is to test how sensitive is the method to the accurate knowledge of the prior distribution.

Figure 6 shows the performance of different algorithms under various scenarios. We show the normalized  $L_1$  norm of the estimation error  $\|\hat{\mathbf{M}} - \mathbf{M}^h\|_1/n$ . On the left panel,

we fix  $\mathbf{B} = \mathbf{I}$ , and increase  $k$  from 5 to 50. As we can see, CD-MVSI gives acceptable recovery result, even though the sufficiently scattered condition is not exactly satisfied. Both GeoNMF and SPOC give much worse performance compared to CD-MVSI. Tensor method with the correct Dirichlet parameter gives nearly perfect recovery, thanks to its nice identifiability guarantees. However, a correct Dirichlet is absolutely necessary, as an incorrect one gives the worst performance. On the right panel, we fix  $k = 10$  and let  $\mathbf{B}$  take nonzero off-diagonal values, and similar patterns present. We remark that it is relatively easy to approximately satisfy the sufficiently scattered condition, whereas knowing exactly the Dirichlet parameter of the prior is somewhat impractical.