

Supplementary Material: A probabilistic framework for multi-view feature learning with many-to-many associations via neural networks

A. Proof of Theorem 5.1

Since $g_* : [-M', M']^{2K_*} \rightarrow \mathbb{R}$ is a positive definite kernel on a compact set, it follows from Mercer's theorem that there exist positive eigenvalues $\{\lambda_k\}_{k=1}^\infty$ and continuous eigenfunctions $\{\phi_k\}_{k=1}^\infty$ such that

$$g_*(\mathbf{y}_*, \mathbf{y}'_*) = \sum_{k=1}^\infty \lambda_k \phi_k(\mathbf{y}_*) \phi_k(\mathbf{y}'_*), \quad \mathbf{y}_*, \mathbf{y}'_* \in [-M', M']^{K_*},$$

where the convergence is absolute and uniform (Minh et al., 2006). The uniform convergence implies that for any $\varepsilon_1 > 0$ there exists $K_0 \in \mathbb{N}$ such that

$$\sup_{(\mathbf{y}_*, \mathbf{y}'_*) \in [-M', M']^{2K_*}} \left| g_*(\mathbf{y}_*, \mathbf{y}'_*) - \sum_{k=1}^K \lambda_k \phi_k(\mathbf{y}_*) \phi_k(\mathbf{y}'_*) \right| < \varepsilon_1, \quad K \geq K_0.$$

This means $g_*(\mathbf{y}_*, \mathbf{y}'_*) \approx \langle \Phi_K(\mathbf{y}_*), \Phi_K(\mathbf{y}'_*) \rangle$ for a feature map $\Phi_K(\mathbf{y}_*) = (\sqrt{\lambda_k} \phi_k(\mathbf{y}_*))_{k=1}^K$.

We fix K and consider approximation of $h_k^{(d)}(\mathbf{x}) := \sqrt{\lambda_k} \phi_k(f_*^{(d)}(\mathbf{x}))$ below. Since $h_k^{(d)}$ are continuous functions on a compact set, there exists $C = C(K) > 0$ such that

$$\sup_{\mathbf{x} \in [-M, M]^{p_d}} |h_k^{(d)}(\mathbf{x})| < C, \quad k = 1, \dots, K, d = 1, \dots, D.$$

Let us write the neural networks as $f_\psi^{(d)} = (f_1^{(d)}, \dots, f_K^{(d)})$, where $f_k^{(d)} : \mathbb{R}^{p_d} \rightarrow \mathbb{R}$, $d = 1, \dots, D$, $k = 1, \dots, K$, are two-layer neural networks with T hidden units. Since $h_k^{(d)}$ are continuous functions, it follows from the universal approximation theorem (Cybenko, 1989; Telgarsky, 2017) that for any $\varepsilon_2 > 0$, there exists $T_0(K) \in \mathbb{N}$ such that

$$\sup_{\mathbf{x} \in [-M, M]^{p_d}} |h_k^{(d)}(\mathbf{x}) - f_k^{(d)}(\mathbf{x})| < \varepsilon_2, \quad k = 1, \dots, K, d = 1, \dots, D$$

for $T \geq T_0(K)$. Therefore, for all $d, e \in \{1, 2, \dots, D\}$, we have

$$\begin{aligned} & \sup_{(\mathbf{x}, \mathbf{x}') \in [-M, M]^{p_d + p_e}} \left| g_* \left(f_*^{(d)}(\mathbf{x}), f_*^{(e)}(\mathbf{x}') \right) - \sum_{k=1}^K f_k^{(d)}(\mathbf{x}) f_k^{(e)}(\mathbf{x}') \right| \\ & \leq \sup_{(\mathbf{x}, \mathbf{x}') \in [-M, M]^{p_d + p_e}} \left| g_* \left(f_*^{(d)}(\mathbf{x}), f_*^{(e)}(\mathbf{x}') \right) - \sum_{k=1}^K h_k^{(d)}(\mathbf{x}) h_k^{(e)}(\mathbf{x}') \right| \\ & \quad + \sup_{(\mathbf{x}, \mathbf{x}') \in [-M, M]^{p_d + p_e}} \left| \sum_{k=1}^K h_k^{(d)}(\mathbf{x}) \left(h_k^{(e)}(\mathbf{x}') - f_k^{(e)}(\mathbf{x}') \right) \right| \\ & \quad + \sup_{(\mathbf{x}, \mathbf{x}') \in [-M, M]^{p_d + p_e}} \left| \sum_{k=1}^K \left(h_k^{(d)}(\mathbf{x}) - f_k^{(d)}(\mathbf{x}) \right) f_k^{(e)}(\mathbf{x}') \right| \\ & \leq \sup_{\mathbf{y}_*, \mathbf{y}'_* \in [-M', M']^{K_*}} \left| g_*(\mathbf{y}_*, \mathbf{y}'_*) - \sum_{k=1}^K \lambda_k \phi_k(\mathbf{y}_*) \phi_k(\mathbf{y}'_*) \right| \\ & \quad + \sum_{k=1}^K \sup_{\mathbf{x} \in [-M, M]^{p_d}} |h_k^{(d)}(\mathbf{x})| \sup_{\mathbf{x}' \in [-M, M]^{p_e}} |h_k^{(e)}(\mathbf{x}') - f_k^{(e)}(\mathbf{x}')| \\ & \quad + \sum_{k=1}^K \sup_{\mathbf{x} \in [-M, M]^{p_d}} |h_k^{(d)}(\mathbf{x}) - f_k^{(d)}(\mathbf{x})| \sup_{\mathbf{x}' \in [-M, M]^{p_e}} |f_k^{(e)}(\mathbf{x}')| \\ & < \varepsilon_1 + KC\varepsilon_2 + K\varepsilon_2(C + \varepsilon_2). \end{aligned}$$

By letting $\varepsilon_1 = \varepsilon/2, \varepsilon_2 = \min(C, \varepsilon/(6KC))$, the last formula becomes smaller than ε , thus proving

$$\sup_{(\mathbf{x}, \mathbf{x}') \in [-M, M]^{p_d + p_e}} \left| g_* \left(f_*^{(d)}(\mathbf{x}), f_*^{(e)}(\mathbf{x}') \right) - \sum_{k=1}^K f_k^{(d)}(\mathbf{x}) f_k^{(e)}(\mathbf{x}') \right| < \varepsilon, \quad d, e = 1, \dots, D.$$

□

B. Consistency of MLE in PMvGE

In this section, we provide technical details of the argument of Section 5.2.

For proving the consistency of MLE, we introduce the following generative model. Let $d_i, i = 1, \dots, n$, be random variables independently distributed with the probability $\mathbb{P}(d_i = d) = \eta^{(d)} \in (0, 1)$ where $\sum_{d=1}^D \eta^{(d)} = 1$. Data vectors are also treated as random variables. The conditional distribution of \mathbf{x}_i given d_i is

$$\mathbf{x}_i \mid d_i \stackrel{\text{indep.}}{\sim} q^{(d_i)}, \quad i = 1, \dots, n,$$

where $q^{(d_i)}$ is a distribution on a compact support in $\mathbb{R}^{p_{d_i}}$. Let us denote eq. (3) as $\mu_{ij}(\mathbf{x}_i, \mathbf{x}_j, d_i, d_j \mid \boldsymbol{\alpha}, \boldsymbol{\psi})$ for indicating the dependency on $(\mathbf{x}_i, \mathbf{x}_j, d_i, d_j)$. The conditional distributions of link weights are already specified in (1) as

$$w_{ij} \mid \mathbf{x}_i, \mathbf{x}_j, d_i, d_j \stackrel{\text{indep.}}{\sim} \text{Po}(\mu_{ij}^*), \quad i, j = 1, \dots, n,$$

where $\mu_{ij}^* = \mu_{ij}(\mathbf{x}_i, \mathbf{x}_j, d_i, d_j \mid \boldsymbol{\alpha}_*, \boldsymbol{\psi}_*)$ with a true parameter $(\boldsymbol{\alpha}_*, \boldsymbol{\psi}_*)$. Due to the constraints $\boldsymbol{\alpha} = \boldsymbol{\alpha}^\top$ and $w_{ij} = 0$ ($(d_i, d_j) \notin \mathcal{D}$), the vector of free parameters in $\boldsymbol{\alpha}$ is $\boldsymbol{\alpha}_{\mathcal{D}} := \{\alpha^{(d,e)}\}_{(d,e) \in \mathcal{D}, d \leq e} \in \mathbb{R}_{\geq 0}^{|\mathcal{D}|}$, and we write $\boldsymbol{\theta} := (\boldsymbol{\alpha}_{\mathcal{D}}, \boldsymbol{\psi})$. Let $\tilde{n} := |\mathcal{I}_n| = O(n^2)$ denote the number of terms in the sum of $\ell_n(\boldsymbol{\theta})$ in eq. (6). Then the expected value of $\tilde{n}^{-1} \ell_n(\boldsymbol{\theta})$ under the generative model with the true parameter $\boldsymbol{\theta}_*$ is expressed as

$$\ell(\boldsymbol{\theta}) := E_{\mathbf{x}_1, \mathbf{x}_2, d_1, d_2} \left[\mu_{12}(\mathbf{x}_1, \mathbf{x}_2, d_1, d_2 \mid \boldsymbol{\theta}_*) \log \mu_{12}(\mathbf{x}_1, \mathbf{x}_2, d_1, d_2 \mid \boldsymbol{\theta}) - \mu_{12}(\mathbf{x}_1, \mathbf{x}_2, d_1, d_2 \mid \boldsymbol{\theta}) \right].$$

If it were the case of i.i.d. observations of sample size n , we would have that $\tilde{n}^{-1} \ell_n(\boldsymbol{\theta})$ converges to $\ell(\boldsymbol{\theta})$ as $n \rightarrow \infty$ from the law of large numbers. In Theorem B.1, we actually prove the uniform convergence in probability, but we have to pay careful attention to the fact that n^2 observations of $(\mathbf{x}_i, \mathbf{x}_j, d_i, d_j)$ are not independent when indices overlap.

Theorem B.1 Let us assume that the parameter space of $\boldsymbol{\theta}$ is $\Theta := [\delta, 1/\delta]^{|\mathcal{D}|} \times \Psi$, where $\delta \in (0, 1)$ is a sufficiently small constant and $\Psi \subset \mathbb{R}^q$ is a compact set. Assume also that the transformations $f_{\boldsymbol{\psi}}^{(d)}(\mathbf{x}), d = 1, \dots, D$, are Lipschitz continuous with respect to $(\boldsymbol{\psi}, \mathbf{x})$. Then we have, as $n \rightarrow \infty$,

$$\sup_{\boldsymbol{\theta} \in \Theta} |\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta})| \xrightarrow{P} 0. \quad (13)$$

Proof B.1 We refer to Corollary 2.2 in Newey (1991). This corollary shows $\sup_{\boldsymbol{\theta} \in \Theta} |\hat{Q}_n(\boldsymbol{\theta}) - \bar{Q}_n(\boldsymbol{\theta})| = o_p(1)$ under general setting of $\hat{Q}_n(\boldsymbol{\theta})$ and $\bar{Q}_n(\boldsymbol{\theta})$. Here we consider the case of $\hat{Q}_n(\boldsymbol{\theta}) = \tilde{n}^{-1} \ell_n(\boldsymbol{\theta})$ and $\bar{Q}_n(\boldsymbol{\theta}) = \ell(\boldsymbol{\theta})$. For showing (13), the four conditions of the corollary are written as follows. (i) Θ is compact, (ii) $\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}) \xrightarrow{P} 0$ for each $\boldsymbol{\theta} \in \Theta$, (iii) $\ell(\boldsymbol{\theta})$ is continuous, (iv) there exists $B_n = O_p(1)$ such that $|\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) - \tilde{n}^{-1} \ell_n(\boldsymbol{\theta}')| \leq B_n \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2$ for all $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \Theta$. The two conditions (i) and (iii) hold obviously, and thus we verify (ii) and (iv) below.

Before verifying (ii) and (iv), we first consider an array $\mathbf{Z} := (Z_{ij})$ of random variables $Z_{ij} \in \mathcal{Z}, (i, j) \in \mathcal{I}_n$, and a bounded and continuous function $h : \mathcal{Z} \rightarrow \mathbb{R}$. We assume that $\mathcal{Z} \subset \mathbb{R}$ is a compact set, and Z_{ij} is independent of Z_{kl} if

$k, l \in \mathcal{R}_n(i, j) := \{(k, l) \in \mathcal{I}_n \mid k, l \in \{1, \dots, n\} \setminus \{i, j\}\}$ for all $(i, j) \in \mathcal{I}_n$. Then we have

$$\begin{aligned} V_{\mathbf{Z}} \left[\frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} h(Z_{ij}) \right] &= E_{\mathbf{Z}} \left[\left(\frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} h(Z_{ij}) \right)^2 \right] - E_{\mathbf{Z}} \left[\frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} h(Z_{ij}) \right]^2 \\ &= \frac{1}{\tilde{n}^2} \left\{ \sum_{(i,j) \in \mathcal{I}_n} \sum_{(k,l) \in \mathcal{I}_n} E_{\mathbf{Z}} [h(Z_{ij})h(Z_{kl})] - \left(\sum_{(i,j) \in \mathcal{I}_n} E_{\mathbf{Z}} [h(Z_{ij})] \right)^2 \right\} \\ &= \frac{1}{\tilde{n}^2} \sum_{(i,j) \in \mathcal{I}_n} \sum_{(k,l) \in \mathcal{I}_n \setminus \mathcal{R}_n(i,j)} (E_{\mathbf{Z}} [h(Z_{ij})h(Z_{kl})] - E_{\mathbf{Z}} [h(Z_{ij})]E_{\mathbf{Z}} [h(Z_{kl})]). \end{aligned}$$

By considering $|\mathcal{I}_n \setminus \mathcal{R}_n(i, j)| = O(n)$, the last formula is $O(\tilde{n}^{-2} \cdot \tilde{n} \cdot n) = O(n^{-1})$. Therefore,

$$V_{\mathbf{Z}} \left[\frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} h(Z_{ij}) \right] = O(n^{-1}). \quad (14)$$

Next we evaluate the variance of $\tilde{n}^{-1} \ell_n(\boldsymbol{\theta})$ to show (ii). Denoting $\mathbf{W} := (w_{ij})$, $\mathbf{X} := (\mathbf{x}_i)$, $\mathbf{d} := (d_i)$,

$$\begin{aligned} V_{\mathbf{W}, \mathbf{X}, \mathbf{d}}[\tilde{n}^{-1} \ell_n(\boldsymbol{\theta})] &= E_{\mathbf{X}, \mathbf{d}}[V_{\mathbf{W}}[\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) \mid \mathbf{X}, \mathbf{d}]] + V_{\mathbf{X}, \mathbf{d}}[E_{\mathbf{W}}[\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) \mid \mathbf{X}, \mathbf{d}]] \\ &= E_{\mathbf{X}, \mathbf{d}} \left[\frac{1}{\tilde{n}^2} \sum_{(i,j) \in \mathcal{I}_n} \mu_{ij}(\boldsymbol{\theta}_*) (\log \mu_{ij}(\boldsymbol{\theta}))^2 \right] + V_{\mathbf{X}, \mathbf{d}} \left[\frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} (\mu_{ij}(\boldsymbol{\theta}_*) \log \mu_{ij}(\boldsymbol{\theta}) - \mu_{ij}(\boldsymbol{\theta})) \right], \end{aligned}$$

for every $\boldsymbol{\theta} \in \Theta$. The first term in the last formula is $O(\tilde{n}^{-1} \cdot \tilde{n}) = O(\tilde{n}^{-1}) = o(1)$, and the second term is $O(n^{-1}) = o(1)$ by applying eq. (14) with $Z_{ij} := (\mathbf{x}_i, \mathbf{x}_j, d_i, d_j)$, $h(Z_{ij}) = \mu_{ij}(\boldsymbol{\theta}_*) \log \mu_{ij}(\boldsymbol{\theta}) - \mu_{ij}(\boldsymbol{\theta})$. Therefore, $V_{\mathbf{W}, \mathbf{X}, \mathbf{d}}[\tilde{n}^{-1} \ell_n(\boldsymbol{\theta})] = o(1)$ and Chebyshev's inequality implies the pointwise convergence $\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) \xrightarrow{p} \ell(\boldsymbol{\theta})$ for every $\boldsymbol{\theta} \in \Theta$ where $\ell(\boldsymbol{\theta}) = E_{\mathbf{W}, \mathbf{X}, \mathbf{d}}[\tilde{n}^{-1} \ell_n(\boldsymbol{\theta})] = E_{\mathbf{x}_1, \mathbf{x}_2, d_1, d_2}[\mu_{12}(\boldsymbol{\theta}_*) \log \mu_{12}(\boldsymbol{\theta}) - \mu_{12}(\boldsymbol{\theta})]$. Thus, condition (ii) holds.

Finally, we work on condition (iv). Since $\mu_{ij}(\boldsymbol{\theta})$ is a composite function of C^1 -functions on Θ , $\mu_{ij}(\boldsymbol{\theta})$ is Lipschitz continuous. The Lipschitz continuity of $\mu_{ij}(\boldsymbol{\theta})$ and $\mu_{ij}(\boldsymbol{\theta}) > 0$ ($\boldsymbol{\theta} \in \Theta$) indicates the Lipschitz continuity of $\log \mu_{ij}(\boldsymbol{\theta})$. Therefore, there exist $M_1, M_2 > 0$ such that

$$\begin{aligned} |\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) - \tilde{n}^{-1} \ell_n(\boldsymbol{\theta}')| &\leq \left| \frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} w_{ij} (\log \mu_{ij}(\boldsymbol{\theta}) - \log \mu_{ij}(\boldsymbol{\theta}')) - \frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} (\mu_{ij}(\boldsymbol{\theta}) - \mu_{ij}(\boldsymbol{\theta}')) \right| \\ &\leq \frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} w_{ij} |\log \mu_{ij}(\boldsymbol{\theta}) - \log \mu_{ij}(\boldsymbol{\theta}')| + \frac{1}{\tilde{n}} \sum_{(i,j) \in \mathcal{I}_n} |\mu_{ij}(\boldsymbol{\theta}) - \mu_{ij}(\boldsymbol{\theta}')| \\ &\leq M_1 \left(\tilde{n}^{-1} \sum_{(i,j) \in \mathcal{I}_n} w_{ij} \right) \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2 + M_2 \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2. \end{aligned}$$

Denoting by $B_n := M_1 \cdot \tilde{n}^{-1} \sum_{(i,j) \in \mathcal{I}_n} w_{ij} + M_2$, we have

$$|\tilde{n}^{-1} \ell_n(\boldsymbol{\theta}) - \tilde{n}^{-1} \ell_n(\boldsymbol{\theta}')| \leq B_n \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2.$$

Since $\tilde{n}^{-1} \sum_{(i,j) \in \mathcal{I}_n} w_{ij} = O_p(1)$, the law of large numbers indicates $B_n = O_p(1)$. Thus, condition (iv) holds. \square

Noticing that $\boldsymbol{\theta}_*$ is a maximizer of $\ell(\boldsymbol{\theta})$ and $\hat{\boldsymbol{\theta}}_n$ is a maximizer of $\ell_n(\boldsymbol{\theta})$, we would have the desired result $\hat{\boldsymbol{\theta}}_n \xrightarrow{p} \boldsymbol{\theta}_*$ by combining Theorem B.1 and continuity of $\ell(\boldsymbol{\theta})$. However it does not hold unfortunately. Instead, we define the set of parameter values equivalent to $\boldsymbol{\theta}_*$ as $\Theta_* := \{\boldsymbol{\theta} \in \Theta \mid \ell(\boldsymbol{\theta}) = \ell(\boldsymbol{\theta}_*)\}$. Every $\boldsymbol{\theta} \in \Theta_*$ gives the correct probability of link weights, because $\ell(\boldsymbol{\theta}) = \ell(\boldsymbol{\theta}_*)$ holds if and only if $\mu_{12}(\boldsymbol{\theta}) = \mu_{12}(\boldsymbol{\theta}_*)$ almost surely w.r.t. $(\mathbf{x}_1, \mathbf{x}_2, d_1, d_2)$. With this setting, the theorem below states that $\hat{\Theta}_n$ converges to Θ_* in probability. This indicates that, $\hat{\boldsymbol{\theta}}_n$ will represent the true probability model for sufficiently large n .

Theorem B.2 Let $d_H(\cdot, \cdot)$ denote the Hausdorff distance defined as max-min L^2 -distance between two sets. We assume the same conditions as in Theorem B.1. Then we have, as $n \rightarrow \infty$,

$$d_H(\hat{\Theta}_n, \Theta_*) \xrightarrow{P} 0. \quad (15)$$

Proof B.2 We refer to the case (1) of Theorem 3.1 in Chernozhukov et al. (2007) with $\hat{c} = 1$ under the condition C.1 This theorem shows, for general setting of $\hat{\Theta}_I, \Theta_I$, that

$$d_H(\hat{\Theta}_I, \Theta_I) \xrightarrow{P} 0. \quad (16)$$

Here $\hat{\Theta}_I := \{\theta \in \Theta \mid \hat{Q}_n(\theta) \leq 1/a_n\}$ and $\Theta_I := \arg \inf_{\theta \in \Theta} Q(\theta)$, where $\hat{Q}_n(\theta), Q(\theta)$ are general functions satisfying $\sup_{\theta \in \Theta} Q_n(\theta) = o_p(1/a_n)$, and $a_n \rightarrow \infty$.

For proving (15), we consider the case of $\hat{Q}_n(\theta) = -\tilde{n}^{-1}\ell_n(\theta) + \sup_{\theta \in \Theta} \tilde{n}^{-1}\ell_n(\theta)$, $Q(\theta) = -\ell(\theta) + \sup_{\theta \in \Theta} \ell(\theta)$. The condition C.1 for (16) is re-written as follows. (i) Θ is a (non-empty) compact set, (ii) $\tilde{n}^{-1}\ell_n(\theta)$ and $\ell(\theta)$ are continuous, (iii) $\sup_{\theta \in \Theta} |\tilde{n}^{-1}\ell_n(\theta) - \ell(\theta)| \xrightarrow{P} 0$, and (iv) $\sup_{\theta \in \Theta_*} (-\tilde{n}^{-1}\ell_n(\theta) + \sup_{\theta \in \Theta} \tilde{n}^{-1}\ell_n(\theta)) \xrightarrow{P} 0$. The conditions (i), (ii) are obvious. (iii) is shown in Theorem B.1. (iv) is verified by

$$\sup_{\theta \in \Theta_*} \left(-\tilde{n}^{-1}\ell_n(\theta) + \sup_{\theta \in \Theta} \tilde{n}^{-1}\ell_n(\theta) \right) = - \inf_{\theta \in \Theta_I} \tilde{n}^{-1}\ell_n(\theta) + \sup_{\theta \in \Theta} \tilde{n}^{-1}\ell_n(\theta) \xrightarrow{P} -\ell(\theta_*) + \ell(\theta_*) = 0,$$

where θ_* is an element of Θ_I . Thus, (16) holds.

Next, we consider two sets $\hat{\Theta}_n = \arg \sup_{\theta \in \Theta} \tilde{n}^{-1}\ell_n(\theta) = \{\theta \in \Theta \mid Q_n(\theta) = 0\}$ and $\hat{\Theta}_I$. Since these sets satisfy $Q_n(\theta) = 0$ ($\theta \in \hat{\Theta}_n$), $Q_n(\theta') \leq 1/a_n$ ($\theta' \in \hat{\Theta}_I$) and $1/a_n \rightarrow 0$ as $n \rightarrow \infty$, we have $d_H(\hat{\Theta}_n, \hat{\Theta}_I) \xrightarrow{P} 0$. It follows from this convergence and (16) that, by noticing $\Theta_* = \Theta_I$,

$$d_H(\hat{\Theta}_n, \Theta_*) = d_H(\hat{\Theta}_n, \Theta_I) \leq d_H(\hat{\Theta}_n, \hat{\Theta}_I) + d_H(\hat{\Theta}_I, \Theta_I) \xrightarrow{P} 0,$$

thus (15) holds. \square

C. CDMCA is approximated by PMvGE with linear transformations

We argue an approximate relation between CDMCA and PMvGE, which is briefly explained in Section 3.6. In the below, we will derive the solution $\hat{\psi}_{\text{CDMCA}}$ of a slightly modified version of CDMCA, and an approximate solution $\hat{\psi}_{\text{Apr.PMvGE}}$ of PMvGE with linear transformations. We then show that these two solutions are equivalent up to a scaling in each axis of the shared space.

C.1. Solution of a modified CDMCA

The original CDMCA imposes the quadratic constraint (8) for maximizing the objective function (7). Here we replace w_{ij} in (8) with δ_{ij} so that the constraint becomes

$$\sum_{i=1}^n \psi^{(d_i)\top} \mathbf{x}_i \mathbf{x}_i^\top \psi^{(d_i)} = \mathbf{I}.$$

This modification changes the scaling in the solution, but the computation below is essentially the same as that in Shimodaira (2016). Let us define the augmented data vector, called ‘‘simple coding’’ (Shimodaira, 2016), $\tilde{\mathbf{x}}_i := (\mathbf{0}_{p_1}, \dots, \mathbf{0}_{p_{d_i-1}}, \mathbf{x}_i, \mathbf{0}_{p_{d_i+1}}, \dots, \mathbf{0}_{p_D}) \in \mathbb{R}^p$ where $p := p_1 + p_2 + \dots + p_D$. Now, data matrix is $\mathbf{X} := (\tilde{\mathbf{x}}_1^\top, \tilde{\mathbf{x}}_2^\top, \dots, \tilde{\mathbf{x}}_n^\top)^\top \in \mathbb{R}^{n \times p}$, and the parameter matrix is $\psi := (\psi^{(1)\top}, \psi^{(2)\top}, \dots, \psi^{(D)\top})^\top \in \mathbb{R}^{p \times K}$. With this augmented representation, the D -view embedding is now interpreted as a 1-view embedding. CDMCA maximizes the objective function

$$\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_{ij} \langle \psi^{(d_i)\top} \mathbf{x}_i, \psi^{(d_j)\top} \mathbf{x}_j \rangle = \text{tr} \left(\psi^\top \mathbf{H} \psi \right) \quad (\mathbf{H} := \mathbf{X}^\top \mathbf{W} \mathbf{X}),$$

with respect to ψ under constraint $\psi^\top \mathbf{G} \psi = \mathbf{I}$ where $\mathbf{G} := \mathbf{X}^\top \mathbf{X}$. Let \mathbf{U}_K be the matrix composed of the top- K eigenvectors of $\mathbf{G}^{-1/2} \mathbf{H} \mathbf{G}^{-1/2}$. Then the solution of the modified CDMCA is

$$\hat{\psi}_{\text{CDMCA}} := \mathbf{G}^{-1/2} \mathbf{U}_K.$$

C.2. Approximate solution of PMvGE with linear transformations

MLE of PMvGE maximizes $\ell_n(\boldsymbol{\alpha}, \boldsymbol{\psi})$ defined in (6). Here we modify it by adding an extra term as $\tilde{\ell}_n(\boldsymbol{\alpha}, \boldsymbol{\psi}) := \ell_n(\boldsymbol{\alpha}, \boldsymbol{\psi}) - \frac{1}{2} \sum_{i:(d_i, d_i) \in \mathcal{D}} \mu_{ii}(\boldsymbol{\alpha}, \boldsymbol{\psi})$. The difference approaches zero for large n , because $|\ell_n(\boldsymbol{\alpha}, \boldsymbol{\psi}) - \tilde{\ell}_n(\boldsymbol{\alpha}, \boldsymbol{\psi})|/|\ell_n(\boldsymbol{\alpha}, \boldsymbol{\psi})| = O(n^{-1})$. Since the parameter $\boldsymbol{\alpha}$ is not considered in CDMCA, we assume $\mathcal{D} := \{\text{all pairs of views}\}$ and $\bar{\boldsymbol{\alpha}} = (\bar{\alpha}^{(de)})$, $\bar{\alpha}^{(de)} \equiv \alpha_0 > 0 (\forall d, e)$. We further assume that the transformation of PMvGE is linear: $f_{\boldsymbol{\psi}}^{(d)}(\mathbf{x}) = \boldsymbol{\psi}^{(d)\top} \mathbf{x} (\forall d)$, and data vectors in each view are centered. With this setting, we will show that the maximizer of a quadratic approximation of $\tilde{\ell}_n$ is equivalent to CDMCA.

To rewrite the likelihood function as $\tilde{\ell}_n(\bar{\boldsymbol{\alpha}}, \boldsymbol{\psi}) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n S_{ij}(g_{ij}(\boldsymbol{\psi}))$, we define $g_{ij}(\boldsymbol{\psi}) := \langle \boldsymbol{\psi}^{(d_i)\top} \mathbf{x}_i, \boldsymbol{\psi}^{(d_j)\top} \mathbf{x}_j \rangle$ and $S_{ij}(g) := w_{ij} \log(\alpha_0 \exp(g)) - \alpha_0 \exp(g)$, $g \in \mathbb{R}$. Since $S_{ij}(g)$ is approximated quadratically around $g = 0$ by

$$S_{ij}^Q(g) = \alpha_0 \left\{ -\frac{1}{2} g^2 + \left(\frac{w_{ij}}{\alpha_0} - 1 \right) g \right\} + S_{ij}(0),$$

$\tilde{\ell}_n(\bar{\boldsymbol{\alpha}}, \boldsymbol{\psi})$ is approximated quadratically by

$$\begin{aligned} \tilde{\ell}_n^Q(\boldsymbol{\psi}) &:= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n S_{ij}^Q(g_{ij}(\boldsymbol{\psi})) \\ &= \alpha_0 \left\{ -\frac{1}{2} \text{tr} \left((\boldsymbol{\psi}^\top \mathbf{G} \boldsymbol{\psi})^2 \right) + \frac{1}{\alpha_0} \text{tr} \left(\boldsymbol{\psi}^\top \mathbf{H} \boldsymbol{\psi} \right) - \underbrace{\text{tr} \left(\boldsymbol{\psi}^\top \mathbf{X}^\top \mathbf{1}_n \mathbf{1}_n^\top \mathbf{X} \boldsymbol{\psi} \right)}_{=0 (\because \{\mathbf{x}_i\} \text{ is centered.})} \right\} + \underbrace{\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n S_{ij}(0)}_{\text{Const.}}, \end{aligned} \quad (17)$$

where $\mathbf{G} = \mathbf{X}^\top \mathbf{X}$, $\mathbf{H} = \mathbf{X}^\top \mathbf{W} \mathbf{X}$. The function $\tilde{\ell}_n^Q(\boldsymbol{\psi})$ has rotational degrees of freedom: $\tilde{\ell}_n^Q(\boldsymbol{\psi}) = \tilde{\ell}_n^Q(\boldsymbol{\psi} \mathbf{O})$ for any orthogonal matrix $\mathbf{O} \in \mathbb{R}^{K \times K}$. Thus, we impose an additional constraint $\boldsymbol{\psi}^\top \mathbf{G} \boldsymbol{\psi} = \boldsymbol{\Gamma}_K := \text{diag}(\gamma_1, \gamma_2, \dots, \gamma_K)$ for any $(\gamma_1, \dots, \gamma_K) \in \mathbb{R}_{\geq 0}^K$. $\boldsymbol{\psi}$ satisfying this constraint is written as $\boldsymbol{\psi} = \mathbf{G}^{-1/2} \mathbf{V}_K \boldsymbol{\Gamma}_K^{1/2}$ where $\mathbf{V}_K \in \mathbb{R}^{P \times K}$ is a column-orthogonal matrix such that $\mathbf{V}_K^\top \mathbf{V}_K = \mathbf{I}$. By substituting $\boldsymbol{\psi} = \mathbf{G}^{-1/2} \mathbf{V}_K \boldsymbol{\Gamma}_K^{1/2}$ into eq. (17), we have

$$\tilde{\ell}_n^Q(\boldsymbol{\psi}) = \alpha_0 \left\{ -\frac{1}{2} \text{tr} \left(\boldsymbol{\Gamma}_K^2 \right) + \frac{1}{\alpha_0} \text{tr} \left(\boldsymbol{\Gamma}_K \mathbf{S}_K \right) \right\} + \text{Const.} = \frac{\alpha_0}{2} \left\{ \left\| \frac{1}{\alpha_0} \mathbf{S}_K \right\|_{\text{F}}^2 - \left\| \boldsymbol{\Gamma}_K - \frac{1}{\alpha_0} \mathbf{S}_K \right\|_{\text{F}}^2 \right\} + \text{Const.}, \quad (18)$$

where $\mathbf{S}_K = \mathbf{V}_K^\top \mathbf{G}^{-1/2} \mathbf{H} \mathbf{G}^{-1/2} \mathbf{V}_K$ and $\|\cdot\|_{\text{F}}$ denotes the Frobenius norm. This objective function is maximized when $\boldsymbol{\Gamma}_K = \frac{1}{\alpha_0} \mathbf{S}_K$ and $\mathbf{V}_K = \mathbf{U}_K$, because $\min_{\boldsymbol{\Gamma}_K} \|\boldsymbol{\Gamma}_K - \frac{1}{\alpha_0} \mathbf{S}_K\|_{\text{F}}^2 = 0$ is achieved by $\boldsymbol{\Gamma}_K = \frac{1}{\alpha_0} \mathbf{S}_K$, and $\max_{\mathbf{V}_K} \|\mathbf{S}_K\|_{\text{F}}^2$ is achieved by $\mathbf{V}_K = \mathbf{U}_K$ where \mathbf{U}_K is the matrix composed of the top- K eigenvectors of $\mathbf{G}^{-1/2} \mathbf{H} \mathbf{G}^{-1/2}$. Therefore, $\tilde{\ell}_n^Q(\boldsymbol{\psi})$ is maximized by

$$\hat{\boldsymbol{\psi}}_{\text{Apr.PMvGE}} = \mathbf{G}^{-1/2} \mathbf{U}_K \boldsymbol{\Gamma}_K^{1/2}.$$

By substituting $\mathbf{V}_K = \mathbf{U}_K$ into $\boldsymbol{\Gamma}_K = \frac{1}{\alpha_0} \mathbf{S}_K$, we verify that $\boldsymbol{\Gamma}_K$ is a diagonal matrix with $\gamma_k := \lambda_k / \alpha_0$ where λ_k is the k -th largest eigenvalue of $\mathbf{G}^{-1/2} \mathbf{H} \mathbf{G}^{-1/2}$ ($k = 1, 2, \dots, K$).

C.3. Equivalence of the two solutions up to a scaling

By comparing the two solutions, we have

$$\hat{\boldsymbol{\psi}}_{\text{Apr.PMvGE}} = \hat{\boldsymbol{\psi}}_{\text{CDMCA}} \boldsymbol{\Gamma}_K^{1/2}. \quad (19)$$

This simply means that each axis in the shared space is scaled by the factor $\sqrt{\gamma_k}$, $k = 1, \dots, K$. Let $\hat{\mathbf{y}}, \hat{\mathbf{y}}'$ be feature vectors in the shared space computed by the approximate PMvGE with linear transformations, and \mathbf{y}, \mathbf{y}' be feature vectors in the shared space computed by the modified CDMCA. Then the inner product is weighted in PMvGE as

$$\langle \hat{\mathbf{y}}, \hat{\mathbf{y}}' \rangle = \sum_{k=1}^K \hat{y}_k \hat{y}'_k = \sum_{k=1}^K \gamma_k y_k y'_k.$$

References

- Chernozhukov, V., Hong, H., and Tamer, E. (2007). Estimation and Confidence Regions for Parameter Sets in Econometric Models. *Econometrica*, 75(5):1243–1284.
- Cybenko, G. (1989). Approximation by Superpositions of a Sigmoidal Function. *Mathematics of Control, Signals, and Systems (MCSS)*, 2(4):303–314.
- Minh, H. Q., Niyogi, P., and Yao, Y. (2006). Mercer’s Theorem, Feature Maps, and Smoothing. In *International Conference on Computational Learning Theory*, pages 154–168. Springer.
- Newey, W. K. (1991). Uniform Convergence in Probability and Stochastic Equicontinuity. *Econometrica*, pages 1161–1167.
- Shimodaira, H. (2016). Cross-validation of matching correlation analysis by resampling matching weights. *Neural Networks*, 75:126–140.
- Telgarsky, M. (2017). Neural networks and rational functions. In Precup, D. and Teh, Y. W., editors, *Proceedings of the International Conference on Machine Learning (ICML)*.