
Deep Clustering with Incomplete Noisy Pairwise Annotations: A Geometric Regularization Approach

Tri Nguyen¹ Shahana Ibrahim¹ Xiao Fu¹

Abstract

The recent integration of deep learning and pairwise similarity annotation-based constrained clustering—i.e., *deep constrained clustering* (DCC)—has proven effective for incorporating weak supervision into massive data clustering: Less than 1% of pair similarity annotations can often substantially enhance the clustering accuracy. However, beyond empirical successes, there is a lack of understanding of DCC. In addition, many DCC paradigms are sensitive to annotation noise, but performance-guaranteed noisy DCC methods have been largely elusive. This work first takes a deep look into a recently emerged logistic loss function of DCC, and characterizes its theoretical properties. Our result shows that the logistic DCC loss ensures the identifiability of data membership under reasonable conditions, which may shed light on its effectiveness in practice. Building upon this understanding, a new loss function based on geometric factor analysis is proposed to fend against noisy annotations. It is shown that even under *unknown* annotation confusions, the data membership can still be *provably* identified under our proposed learning criterion. The proposed approach is tested over multiple datasets to validate our claims.

1. Introduction

Clustering is one of the most prominent unsupervised learning tasks (Jain & Dubes, 1988). Classic clustering paradigms, e.g., K-means and spectral clustering, are designed to work without any label information. In practice, it was observed that limited supervision may significantly boost the clustering performance. The so-called *constrained*

clustering (CC) method (Wagstaff & Cardie, 2000) is one of such weak (or semi-)supervised approaches. The CC paradigm has many variants (see (Basu et al., 2008; Schultz & Joachims, 2003; Yeung & Yeo, 1996; Davidson et al., 2010; Zhang et al., 2019; 2021)). The arguably most widely used paradigm annotates the similarity (using binary codes) of data pairs and restrains the clustering outcomes to respect these pairwise constraints; see, e.g., (Basu et al., 2004a; Segal et al., 2003; Wagstaff et al., 2001; Givoni & Frey, 2009; Hsu & Kira, 2015; Manduchi et al., 2021). Note that annotating similarity is considerably easier than annotating the exact class labels of data, yet it was observed that such “simple” annotations could boost the performance of data categorization. These findings may assist designing economical annotation mechanisms in the era of data explosion.

Early CC approaches are often combined with existing clustering modules like K-means (Basu et al., 2004a; Wagstaff et al., 2001; Bilenko et al., 2004; Li et al., 2009; Kamvar et al., 2003; Wang et al., 2014; Cucuringu et al., 2016). The pairwise annotations are used to construct regularization terms that are similar to graph-Laplacian based regularization. Variants using Bayesian perspectives were also proposed (Basu et al., 2004b; Law et al., 2005). These methods were observed to largely outperform their unsupervised counterparts, e.g., K-means and spectral clustering, even only as few as 1% of the pairs are annotated.

In recent years, clustering modules and deep neural network-based feature extractors are proposed to be learned jointly—leading to the so-called end-to-end deep clustering approaches; see, e.g., (Yang et al., 2017; Caron et al., 2018). These methods use clustering structures in the latent space to regularize the neural networks, and the neural networks offers enhanced transformation power to find latent spaces where the data can be well clustered. The CC methods also benefited from similar ideas. The so-called *deep constrained clustering* (DCC) approaches were proposed in (Zhang et al., 2019; 2021; Manduchi et al., 2021), and substantial performance enhancement relative to classic CC methods was observed.

Although the DCC methods have enjoyed empirical successes, some notable challenges remain. First, there is a general lack of understanding to the effectiveness of DCC. It is

¹School of Electrical Engineering and Computer Science, Oregon State University, OR, USA. Correspondence to: Xiao Fu <xiao.fu@oregonstate.edu>.

unclear under what conditions the loss functions constructed for DCC could succeed or fail in finding the ground-truth cluster membership of the data entities. However, understanding the *identifiability* of the membership is critical for designing principled and robust DCC systems. The interplay of key aspects, e.g., neural network complexity and the generalization ability of the learned feature extractor, in the context of DCC is also of great interest—yet no pertinent study exists. Second, most of the existing (D)CC methods (implicitly) assumed that annotations are accurate; see, e.g., (Basu et al., 2004a; Wagstaff et al., 2001; Li et al., 2009; Zhang et al., 2019; Ren et al., 2019; Hsu et al., 2018). Hence, many of these methods may not be robust to annotation noise. This is particularly detrimental to DCC methods as large over-parameterized neural models easily overfit (Du et al., 2019). Some works took noisy annotations into consideration (e.g., (Luo et al., 2018; Manduchi et al., 2021)), but no guarantees of recovering the cluster membership.

Contributions. In this work, we take a deeper look at the DCC problem from a membership-identifiability analysis viewpoint. Our contribution is twofold:

First, we re-examine a recently emerged effective loss function of DCC, namely, the logistic loss-based DCC criterion, which was proposed by (Hsu et al., 2018; Zhang et al., 2021). We show that, if the pairwise annotations are generated following a model that is reminiscent of the mixed-membership stochastic blockmodel (MMSB) (Airoldi et al., 2008; Huang & Fu, 2019), then, the logistic loss can provably recover (i) the data entities’ cluster membership and (ii) the nonlinear function that maps the data features to the membership indicator for both seen and unseen data—and thus generalization of the learned neural network is guaranteed.

Second, using our understanding, we propose a noisy annotation-robust version of logistic loss for DCC. We explicitly model the annotators’ confusion as a probability transition matrix, which is inspired by classic noisy label analysis such as the Dawid-Skene model (Dawid & Skene, 1979; Ghosh et al., 2011; Zhang et al., 2014). We propose a geometric factor analysis (Fu et al., 2018; 2019) based learning criterion to *provably* ensure the identifiability of the ground-truth cluster membership, in the presence of annotation confusion.

We test our method over a series of DCC tasks and observe that the proposed approach significantly improves the performance over existing paradigms, especially when annotation noise exists. Our finding shows the significance of identifiability in DCC, echoing observations made in similar semi-supervised/unsupervised problems, e.g., (Arora et al., 2013; Kumar et al., 2013; Anandkumar et al., 2014; Zhang et al., 2014). We also evaluate the algorithms using real data collected through the Amazon Me-

chanical Turk (AMT) platform. The code is published at github.com/ductri/VolMaxDCC.

Notation. We use x , \mathbf{x} and \mathbf{X} to denote scalar, vector and matrix, respectively; both $[\mathbf{x}]_k, x_k$ refer to the k th element of vector \mathbf{x} ; X_{ij} (and $[\mathbf{X}]_{i,j}$) is the element in the i th row and j th column of \mathbf{X} ; \mathbf{I}_K denotes the identity matrix of size K ; $\mathbf{0}$ and $\mathbf{1}$ are all-zero and all-one matrices with proper sizes; $\langle \mathbf{x}, \mathbf{y} \rangle$ and $\langle \mathbf{X}, \mathbf{Y} \rangle$ denote dot products between two vectors and two matrices, respectively; $\|\mathbf{x}\|$ denotes the ℓ_2 -norm; $\|\mathbf{X}\|_F$ and $\|\mathbf{X}\|_2$ denote the Frobenius norm and spectral norm of \mathbf{X} , respectively; $\sigma_{\max}(\mathbf{X})$, $\sigma_{\min}(\mathbf{X})$, and $\sigma_i(\mathbf{X})$ represent the largest, the smallest, and the i th singular value of matrix \mathbf{X} , respectively; $[N]$ is the set of natural numbers from 1 to N , i.e., $[N] = \{1, \dots, N\}$; $[N] \times [N]$ denotes the set of all possible pairs of (i, j) where $i \in [N], j \in [N]$; $\text{cone}(\mathbf{X})$ to denote conic hull of the column vectors of \mathbf{X} , i.e., $\text{cone}(\mathbf{X}) = \{\mathbf{y} \mid \mathbf{y} = \mathbf{X}\boldsymbol{\theta}, \boldsymbol{\theta} \geq \mathbf{0}\}$.

2. Background

Problem Setting. We consider the CC setting as follows: There are N data samples $\mathbf{x}_1, \dots, \mathbf{x}_N$. Each sample belongs to one or multiple clusters of in total K clusters. The association of \mathbf{x}_n with the clusters is represented by a vector $\mathbf{m}_n \in \mathbb{R}^K$. The element $[\mathbf{m}_n]_k$ represents the probability that data n belongs to cluster k . Note that in hard clustering, $[\mathbf{m}_n]_k \in \{0, 1\}$ for all $k \in [K]$, which means the clusters have no overlaps. In more general cases, we have

$$\mathbf{1}^\top \mathbf{m}_n = 1, \mathbf{m}_n \geq \mathbf{0}. \quad (1)$$

A collection of M pairwise annotations are available, which are denoted by $(i_1, j_1, y_1), \dots, (i_M, j_M, y_M)$. Here,

$$y_m = \begin{cases} 1, & \mathbf{x}_{i_m}, \mathbf{x}_{j_m} \text{ are "similar",} \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

where the similarity of the membership of \mathbf{x}_{i_m} and \mathbf{x}_{j_m} is often deemed by an annotator. Note that there are in total $N(N-1)/2$ such data pairs, and we often have

$$M \ll N(N-1)/2;$$

that is, only a small portion of the data pairs are annotated. The objective of pairwise annotation-based CC is to find the cluster membership vector \mathbf{m}_n of each \mathbf{x}_n using the data and the M annotations.

Early CC Methods. The task of clustering with the pairwise constraints can be dated back to the early 2000s, where (Wagstaff & Cardie, 2000; Wagstaff et al., 2001) used the pairwise annotations to impose extra constraints of the classic K-means iterations. The work (Basu et al., 2004a) considered using pairwise annotation-induced “soft” constraints

(or, regularization terms) to modify K-means; see similar ideas in (Bilenko et al., 2004). Instead of modifying K-means, another line of approaches proposed to work with pairwise constraints under the spectral clustering framework. The idea is to incorporate the pairwise annotation-based constraints into the construction of the graph affinity matrix; see, e.g., (Kulis et al., 2005; Lu & Carreira-Perpinan, 2008; Li et al., 2009; Cucuringu et al., 2016).

DCC Developments. In recent years, deep neural network-based feature extractors were proposed to combine with CC—leading to the *deep constrained clustering* (DCC) paradigms, e.g., (Manduchi et al., 2021; Luo et al., 2018; Zhang et al., 2019; 2021). In DCC, a deep neural network (DNN) $f_\theta(\cdot)$ is used to link the data vector \mathbf{x}_n with its membership, i.e.,

$$\mathbf{m}_n = f_\theta(\mathbf{x}_n).$$

The use of DNN helps nonlinearly transform the data to spaces that are “friendly” to clustering (Yang et al., 2017). Learning $f_\theta(\cdot)$ also allows the neural network to generalize to unseen data. By (1),

$$\mathbf{m}_{i_m}^\top \mathbf{m}_{j_m} \in [0, 1]$$

can be considered as the probability that \mathbf{x}_{i_m} and \mathbf{x}_{j_m} belong to the same cluster, i.e., $y_m \sim \text{Bernoulli}(\mathbf{m}_{i_m}^\top \mathbf{m}_{j_m})$. From this perspective, recent works have proposed an extension of logistic regression to incorporate pairwise annotations (Hsu et al., 2018; Zhang et al., 2019; 2021):

$$\text{Loss}_{\text{cc}}(\theta) = \frac{1}{M} \sum_{m=1}^M \left(y_m \log \frac{1}{\mathbf{f}_\theta(\mathbf{x}_{i_m})^\top \mathbf{f}_\theta(\mathbf{x}_{j_m})} + (1 - y_m) \log \frac{1}{1 - \mathbf{f}_\theta(\mathbf{x}_{i_m})^\top \mathbf{f}_\theta(\mathbf{x}_{j_m})} \right). \quad (3)$$

In the literature, the Loss_{cc} term is sometimes used with other loss functions; e.g., (Zhang et al., 2019; 2021) used an overall loss function consisting of two terms:

$$\text{minimize } \text{Loss}_{\text{recon}} + \lambda \text{Loss}_{\text{cc}}, \quad (4)$$

where $\lambda \geq 0$ and the reconstruction loss $\text{Loss}_{\text{recon}}$ was realized using an autoencoder. Such combination is advocated for practical reasons, e.g., utilizing all available data. Nonetheless, as we will show, Loss_{cc} itself suffices to offer strong guarantees under reasonable conditions.

Challenges. Although there have been abundant empirical evidence demonstrating the effectiveness of CC and DCC, theoretical understanding has been largely behind. This is particularly obvious for the DCC case, where aspects such as the identifiability of \mathbf{m}_n and the generalization ability of the learned $f_\theta(\cdot)$ are of great interest—yet no theoretical support exists, to our best knowledge.

Another challenge lies in noise robustness. Although it has been widely observed that noisy annotations could greatly impact the performance of CC and DCC (Liu et al., 2017; Covoes et al., 2013; Pelleg & Baras, 2007; Manduchi et al., 2021; Luo et al., 2018; Chang et al., 2017b; Zhang et al., 2021; Zhu et al., 2015), effective solutions—especially performance-guaranteed ones—for handling this problem have been largely lacking. Many works did test their algorithms with noisy labels (see, e.g., (Cucuringu et al., 2016; Wang et al., 2014)), but no special care was taken to alleviate the impact of such noisy labels. Some heuristics—such as pre-processing the data (Yi et al., 2012), modeling annotation uncertainties (Manduchi et al., 2021), and introducing concepts reflecting human behaviors (Luo et al., 2018; Chang et al., 2017b) and annotators’ accuracy (Luo et al., 2018; Chang et al., 2017b)—were also proposed in the literature. However, performance guarantees have been elusive.

3. DCC Loss Revisited: Identifiability and Generalization

In this section, we take a deeper look into the DCC loss in (3) and understand its theoretical properties. Such understanding will allow us to design a performance-guaranteed new DCC loss in the presence of noisy pairwise annotations.

3.1. A Generative Model of Annotations

To better present the results, in this section, we use the superscript “ \natural ” to denote all the ground-truth terms; e.g., $\mathbf{f}^\natural(\cdot)$ denotes the ground-truth nonlinear mapping from data to membership and $\mathbf{m}_n^\natural = \mathbf{f}^\natural(\mathbf{x}_n)$ denotes the ground-truth membership vector of sample n . We propose to employ the following generative model of y_m : Given $\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathcal{P}_{\mathcal{X}}$, $\mathbf{x}_n \in \mathcal{X}$,

$$i, j \text{ are sampled over } [N] \times [N]; \quad (5a)$$

$$\mathbf{m}_i^\natural = \mathbf{f}^\natural(\mathbf{x}_i) \text{ and } \mathbf{m}_j^\natural = \mathbf{f}^\natural(\mathbf{x}_j); \quad (5b)$$

$$y_{i,j} \sim \text{Bernoulli}(\langle \mathbf{m}_i^\natural, \mathbf{m}_j^\natural \rangle). \quad (5c)$$

Note that the logistic loss in (3) is the maximum likelihood estimator (MLE) of the parameters in the generative model. The model is reminiscent of the classic generative models of logistic regression and network analysis, particularly, MMSB (Airoldi et al., 2008). In MMSB, the nonlinear mapping from the data features to the membership vectors were not considered. Incorporating the nonlinear mapping \mathbf{f}^\natural follows the ideas from supervised learning, where the relationship between the data features and the data class is often modeled as the following conditional probability represented by \mathbf{f}^\natural (Shalev-Shwartz & Ben-David, 2014):

$$[\mathbf{m}_n^\natural]_k = \Pr(y = k | \mathbf{x}_n) = [\mathbf{f}^\natural(\mathbf{x}_n)]_k.$$

Once \mathbf{f}^{\natural} is learned, it can be used as a multi-class classifier. This perspective was also mentioned in (Hsu et al., 2018; Zhang et al., 2019; 2021)—for algorithm design purpose. However, membership identifiability was not addressed.

In this section, we will show that the logistic loss (3) ensures identifying the membership vectors \mathbf{m}_n under the generative model in (5). It also ensures that $\mathbf{f}_{\theta} \approx \mathbf{f}^{\natural}$, under reasonable conditions. These findings may shed some light onto the effectiveness and good generalization performance of DCC using (3) in the literature.

3.2. Performance Analysis

Finite-Sample Identifiability and Generalization. We first show that Loss_{cc} is a sound criterion for identifying \mathbf{m}_n^{\natural} and $\mathbf{f}^{\natural}(\cdot)$ by itself. Specifically, let us denote

$$\theta^* = \arg \min_{\theta} \text{Loss}_{\text{cc}}(\theta) \quad (6)$$

and $\mathbf{f}^* = \mathbf{f}_{\theta^*}$. Here, \mathbf{f}_{θ} is represented by a deep neural network. To be more specific, we consider \mathbf{f}_{θ} belonging to a function class \mathcal{F} defined by

$$\mathcal{F} \triangleq \{\text{softmax}(\text{net}(\mathbf{x}; \theta)) \mid \forall \mathbf{x} \in \mathcal{X}\},$$

where $\text{net}(\cdot; \theta)$ is a neural network that maps \mathbf{x} to \mathbb{R}^K , and $[\text{softmax}(\mathbf{x})]_k \triangleq \exp(x_k) / \sum_{\ell=1}^K \exp(x_{\ell})$ is imposed onto the output layer of the neural network, which is used to reflect the constraints in (1).

We will show that $\mathbf{f}^* \approx \mathbf{f}^{\natural}$ and $\mathbf{m}_n^* = \mathbf{f}^*(\mathbf{x}_n) \approx \mathbf{m}_n^{\natural}$ under reasonable conditions. To proceed, let us invoke the following assumptions:

Assumption 3.1 (Anchor Sample Condition (ASC)). Let $\mathbf{m}_n^{\natural} = \mathbf{f}^{\natural}(\mathbf{x}_n)$, $\mathbf{M}^{\natural} = [\mathbf{m}_1^{\natural}, \dots, \mathbf{m}_N^{\natural}]$. \mathbf{M}^{\natural} satisfies ASC if there exists a set \mathcal{K} of K indices such that $\mathbf{M}^{\natural}[:, \mathcal{K}] = \mathbf{I}$. Accordingly, define $\mathbf{V}^{\natural} \triangleq \mathbf{M}^{\natural}[:, \mathcal{K}^c]$, $\mathcal{K}^c \triangleq [N] \setminus \mathcal{K}$.

Assumption 3.2. (Function Class) There exist $0 < \nu < 1$ and $\tilde{\mathbf{f}} \in \mathcal{F}$ such that

$$\|\tilde{\mathbf{f}}(\mathbf{x}) - \mathbf{f}^{\natural}(\mathbf{x})\| \leq \nu, \forall \mathbf{x} \in \mathcal{X}. \quad (7)$$

In addition, $\alpha < \mathbf{f}(\mathbf{x})^{\top} \mathbf{f}(\mathbf{y}) < 1 - \alpha$, $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}, \forall \mathbf{f} \in \mathcal{F}$, for some $0 < \alpha < 1$. The complexity measure of the neural network $\text{net}(\mathbf{x}; \theta)$ is R_{NET} .

The ASC means that there exist samples that solely belong to a single cluster, which are called the anchor samples. This assumption is reminiscent of the anchor point assumption in the community detection literature (Panov et al., 2017; Mao et al., 2017). Condition (7) takes the approximation error of the employed neural network class \mathcal{F} into consideration. The assumption on α is a regularity condition that prevents pathological unbounded cases of the logistic loss from happening. The constant R_{NET} is proportional to the upper

bound of the so-called spectral complexity in (Bartlett et al., 2017). A formal definition of R_{NET} is given in Lemma A.10. The parameters ν and R_{NET} present a tradeoff: Roughly speaking, if one has a deeper and wider neural network, then ν is smaller—but R_{NET} is bigger.

Define $\mathbf{P}^{\natural} \in \mathbb{R}^{N \times N}$ such that $P_{ij}^{\natural} = \langle \mathbf{f}^{\natural}(\mathbf{x}_i), \mathbf{f}^{\natural}(\mathbf{x}_j) \rangle$, $\mathbf{P}^* \in \mathbb{R}^{N \times N}$ such that $P_{ij}^* = \langle \mathbf{f}^*(\mathbf{x}_i), \mathbf{f}^*(\mathbf{x}_j) \rangle$, and $\mathbf{S}_{\mathcal{X}} = [\mathbf{x}_{i_1}, \mathbf{x}_{j_1}, \dots, \mathbf{x}_{i_M}, \mathbf{x}_{j_M}]$. We first show that the matrix \mathbf{P}^{\natural} is approximately recoverable via minimizing (3):

Lemma 3.3. Let $S = \{(i_1, j_1, y_1), \dots, (i_M, j_M, y_M)\}$, where $(i_1, j_1), \dots, (i_M, j_M)$ are drawn independently and uniformly at random from $[N] \times [N]$. Suppose that $y_m \mid (i_m, j_m) \sim \text{Bernoulli}(P_{i_m, j_m}^{\natural})$ following the generative model in (5) and that \mathcal{F} satisfies Assumption 3.2. Then, with probability at least $1 - \delta$, we have

$$\frac{1}{N^2} \|\mathbf{P}^* - \mathbf{P}^{\natural}\|_{\text{F}}^2 \leq \epsilon(M, \delta)^2,$$

where $\epsilon(M, \delta)^2$ is defined as follows:

$$\begin{aligned} \epsilon(M, \delta)^2 \triangleq & \frac{64 \log(1/\alpha)}{M} + \frac{96\sqrt{2} \log M}{\alpha M \log 2} \|\mathbf{S}_{\mathcal{X}}\|_{\text{F}} \sqrt{R_{\text{NET}}} \\ & + 64 \log(1/\alpha) \sqrt{\frac{2 \log(4/\delta)}{M}} + \frac{16\nu}{\alpha}. \end{aligned} \quad (8)$$

The proof of Lemma 3.3 is relegated to Appendix A. It is not surprising that \mathbf{P}^{\natural} can be recovered from minimizing the logistic loss, as the problem of recovering \mathbf{P}^{\natural} can be regarded as a generalized 1-bit matrix completion (MC) problem. Unlike the conventional 1-bit MC frameworks that leveraged the low-rank structure of the complete data (see, e.g., (Davenport et al., 2014)), here we exploit the neural generative model in (5), which is also a low-dimensional model, as long as R_{NET} is sufficiently small (compared to M and N).

Before proceeding to showing recovery of \mathbf{M}^{\natural} , let us observe the following fact based on Lemma 3.3: It can be seen in (8) that $\epsilon(M, M^{-0.5})^2$ is decreasing with a rate of $\mathcal{O}(\sqrt{(\log M)/M})$, hence there exists $M_0 \in \mathbb{N}$ independent to N such that $\forall M > M_0$,

$$\epsilon'(M)^2 \triangleq M^{0.25} \epsilon(M, M^{-0.5})^2 + M^{-0.25} \leq \frac{1}{8K^2}. \quad (9)$$

Using the above fact, we show that

Theorem 3.4. Under the same assumptions as in Lemma 3.3, further assume that Assumption 3.1 is satisfied. Let $\mathbf{M}^* = [\mathbf{f}^*(\mathbf{x}_1), \dots, \mathbf{f}^*(\mathbf{x}_N)]$, and consider $M > M_0$, then there exists a permutation matrix $\mathbf{\Pi}^*$ such that

$$\begin{aligned} \frac{1}{NK} \|\mathbf{\Pi}^* \mathbf{M}^* - \mathbf{M}^{\natural}\|_{\text{F}}^2 & \leq \frac{4N}{K} \epsilon(M, \delta)^2 \\ & + \frac{2K}{N} (1 + 8\sigma_{\max}^2(\mathbf{V}^{\natural})) \epsilon'(M), \end{aligned}$$

holds with probability at least $1 - \delta - K^2/M^{0.25}$, where $\epsilon'(M)$ is defined in (9).

The proof of Theorem 3.4 is in Appendix B. We should mention that the sample complexity in Theorem 3.4 is based on a worst-case analysis, and thus the bound tends to be pessimistic—it starts to make sense when

$$N \leq \mathcal{O}(K\sqrt{M}),$$

as reflected in the first term on the right hand side. In practice, Loss_{cc} can work fairly well using a much smaller M . Nonetheless, Theorem 3.4 for the first time shows the soundness of using Loss_{cc} in (3), in terms of being able to guarantee the identifiability of M^\natural .

With Theorem 3.4, it is readily to show the following:

Theorem 3.5 (Generalization). *Under the same conditions as in Theorem 3.4, with probability at least $1 - \delta - K^2/M^{1/4}$,*

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_x} \left[\frac{1}{K} \|\Pi^* \mathbf{f}^*(\mathbf{x}) - \mathbf{f}^\natural(\mathbf{x})\|^2 \right] &\leq \frac{4N}{K} \epsilon(M, \delta)^2 \\ &+ \frac{2K}{N} (1 + 8\sigma_{\max}^2(\mathbf{V}^\natural)) \epsilon'(M) \\ &+ \frac{4}{NK} + \frac{12\sqrt{2R_{\text{NET}}} \|\mathbf{X}\|_{\text{F}} \log N}{NK \log 2} + 8\sqrt{\frac{2 \log(4/\delta)}{NK^2}}. \end{aligned}$$

See the proof in Appendix C. Theorem 3.5 confirms that the learned \mathbf{f}^* via minimizing Loss_{cc} can generalize to classify unseen data, which was observed in the literature (Hsu et al., 2018; Zhang et al., 2019; 2021) but never formally shown.

Enhanced Identifiability with Large Sample. In Theorem 3.4, the identifiability of M^\natural was established under finite samples. The proof relies on the ASC. One may argue that the anchor samples may not exist in some cases, as some data naturally exhibit mixed membership, e.g., social network users (Airoldi et al., 2008) and multi-topic documents (Blei et al., 2003). Here, we relax the finite sample assumption to show that at the limit of large M , the logistic loss enjoys enhanced membership identifiability that does not hinge on the ASC.

To see this, let us introduce the following condition:

Assumption 3.6 (Sufficiently Scattered Condition (SSC)). Let $\mathbf{m}_n^\natural = \mathbf{f}^\natural(\mathbf{x}_n)$, $M^\natural = [\mathbf{m}_1^\natural, \dots, \mathbf{m}_N^\natural]$. M^\natural satisfies SSC if there exists a subset $\mathcal{S} \in [N]$ such that

$$\mathcal{C} \subseteq \text{cone}(M^\natural[:, \mathcal{S}]),$$

where $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^K \mid \sqrt{K-1} \|\mathbf{x}\| \leq \mathbf{1}^\top \mathbf{x}\}$, and $\text{cone}(M^\natural[:, \mathcal{S}]) \subseteq \text{cone}(\mathbf{Q})$ does not hold for any orthogonal \mathbf{Q} except for the permutation matrices.

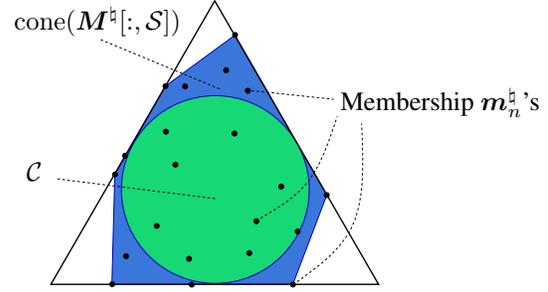


Figure 1: An example of M^\natural satisfying SSC when $K = 3$.

The SSC is illustrated in Fig. 1 with a $K = 3$ case. The triangle represents the nonnegative orthant and the green circle the second-order cone \mathcal{C} . SSC means that the conic hull of a subset of the \mathbf{m}_n^\natural 's (the blue region) covers the green region as a subset. This means that the columns of M^\natural are sufficiently different. The SSC is widely used in the nonnegative matrix factorization literature; see (Fu et al., 2018; 2015b). Note that M^\natural satisfying the SSC is much more relaxed compared to the ASC, which means that $\text{cone}(M^\natural) = \mathbb{R}_+^K$, i.e., the blue region in Fig. 1 covers the entire triangle. In the context of membership learning, an SSC-satisfying M^\natural means that the membership of the data should contain enough diversity, but an ASC-satisfying M^\natural means that some single-membership data samples need to exist—which is more stringent.

Using the SSC, we show that the following holds:

Theorem 3.7. (Enhanced Identifiability) *Under the same assumptions as in Lemma 3.3, suppose that M^\natural satisfies SSC (Assumption 3.6) and that $\mathbf{f}^\natural \in \mathcal{F}$. Then, at the limit of $\max(\log(1/\alpha), \log(M)\sqrt{R_{\text{NET}}}/\alpha)/\sqrt{M} \rightarrow 0$, the following statements hold:*

- (i) *There exists a permutation matrix Π^* such that $\Pi^* M^* = M^\natural$.*
- (ii) *The learned neural network \mathbf{f}^* satisfies*

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_x} \left[\frac{1}{K} \|\Pi^* \mathbf{f}^*(\mathbf{x}) - \mathbf{f}^\natural(\mathbf{x})\|^2 \right] &\leq \frac{8}{NK} + \\ &\frac{12\sqrt{2R_{\text{NET}}} \|\mathbf{X}\|_{\text{F}} \log N}{NK \log 2} + 8\sqrt{\frac{2 \log(4/\delta)}{NK^2}} \end{aligned}$$

with probability at least $1 - \delta$.

The proof of Theorem 3.7 is in Appendix D. Theorem 3.7 has plausible implications in practice: When the sample size M is large and \mathcal{F} is expressive, even no anchor samples exist, identifying M^\natural and finding a generalizable \mathbf{f}^* are still possible.

4. DCC with Noisy Labels: Geometric Regularization

Incorporating Annotation Confusion. We should mention that the generative model in (5) can already model annotation noise to a certain extent: The Bernoulli sampling process can encode some 0-1 flipping probability of observing y_m . Nonetheless, this level of noise consideration is not enough, as annotations could be grossly inaccurate. To take more severe annotation errors into consideration, we modify the generative model as follows. We assume that the annotator confuses class i with class j with probability $\Pr(j|i)$. Let $A_{i,j} = \Pr(i|j)$, we have a “confusion matrix” $\mathbf{A} \in \mathbb{R}^{K \times K}$ where $[\mathbf{A}]_{i,j} = A_{i,j}$. Hence, the annotator’s “confused membership vector” is modeled as

$$\mathbf{m}_n^{\text{confused}} = \mathbf{A}\mathbf{m}_n^{\natural} = \mathbf{A}\mathbf{f}^{\natural}(\mathbf{x}_n). \quad (10)$$

Then, the annotator’s output is sampled from the following:

$$y_m \sim \text{Bernoulli}(\langle \mathbf{A}\mathbf{f}^{\natural}(\mathbf{x}_{i_m}), \mathbf{A}\mathbf{f}^{\natural}(\mathbf{x}_{j_m}) \rangle). \quad (11)$$

Note that using a confusion matrix to model noisy labels’ generating process is widely seen in noisy label learning—but mostly under the supervised learning setting; see, e.g., (Dawid & Skene, 1979; Liu et al., 2012; Zhang et al., 2014; Chu et al., 2021). We argue that this confusion model is also suitable for pairwise annotation. The rationale is that the error happened in comparison is mainly caused by the annotator’s confusion on the membership of \mathbf{x}_{i_m} or that of \mathbf{x}_{j_m} —which is exactly reflected in (11).

Volume Maximization DCC. To proceed, we propose the following modified logistic loss:

$$\text{Loss}'_{\text{cc}}(\boldsymbol{\theta}, \mathbf{B}) = \frac{1}{M} \sum_{m=1}^M \left(y_m \log \frac{1}{\mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_{i_m})^{\top} \mathbf{B} \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_{j_m})} + (1 - y_m) \log \frac{1}{1 - \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_{i_m})^{\top} \mathbf{B} \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_{j_m})} \right),$$

where $\mathbf{B} \in \mathbb{R}^{K \times K}$ satisfies $0 \leq \mathbf{B} \leq 1$, as it is induced by $\mathbf{B} = \mathbf{A}^{\top} \mathbf{A}$. Note that we will use $\boldsymbol{\theta}$ and \mathbf{B} as our optimization variables (instead of \mathbf{A}) as it simplifies the loss function. In addition, as our ultimate goal is to learn \mathbf{M}^{\natural} and \mathbf{f}^{\natural} , the intermediate variable \mathbf{A} does not need to be explicitly estimated.

We show that minimizing $\text{Loss}'_{\text{cc}}(\boldsymbol{\theta}, \mathbf{B})$ provably recovers the data membership and finds a generalizable \mathbf{f}^* , with an additional volume requirement satisfied by the solution. Specifically, we have the following theorem:

Theorem 4.1 (Identifiability of Noisy Case). *Assume that the assumptions in Lemma 3.3 hold, except that the generative model is replaced by (11). Suppose that \mathbf{M}^{\natural} satisfies*

SSC (Assumption 3.6) and that $\mathbf{f}^{\natural} \in \mathcal{F}$. Also assume that $\text{rank}(\mathbf{A}^{\top} \mathbf{A} \mathbf{M}^{\natural}) = K$. Denote

$$(\boldsymbol{\theta}^*, \mathbf{B}^*) = \arg \min \text{Loss}'_{\text{cc}}(\boldsymbol{\theta}, \mathbf{B}) \quad (12)$$

and $\mathbf{f}^ = \mathbf{f}_{\boldsymbol{\theta}^*}$ and $\mathbf{m}_n^* = \mathbf{f}^*(\mathbf{x}_n)$. In addition, assume that Loss'_{cc} is minimized with a solution $\mathbf{M}^*, \mathbf{B}^*$ such that $\log \det(\mathbf{M}^* (\mathbf{M}^*)^{\top})$ is maximized among all possible optimal solutions. Then, at the limit of $\max(\log(1/\alpha), \log(M) \sqrt{R_{\text{NET}}/\alpha}) / \sqrt{M} \rightarrow 0$, the following statements hold:*

(i) *There exists a permutation matrix $\boldsymbol{\Pi}^*$ such that $\boldsymbol{\Pi}^* \mathbf{M}^* = \mathbf{M}^{\natural}$.*

(ii) *The learned neural network \mathbf{f}^* satisfies*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{P}_{\mathbf{x}}} \left[\frac{1}{K} \|\boldsymbol{\Pi}^* \mathbf{f}^*(\mathbf{x}) - \mathbf{f}^{\natural}(\mathbf{x})\|^2 \right] \leq \frac{8}{NK} + \frac{12\sqrt{2R_{\text{NET}}}\|\mathbf{X}\|_{\text{F}} \log N}{NK \log 2} + 8\sqrt{\frac{2 \log(4/\delta)}{NK^2}}$$

with probability at least $1 - \delta$.

The proof of Theorem 4.1 is in Appendix E. The take-home point is that, when one has large samples and an expressive neural network, the membership identifiability and generalization performance of using Loss'_{cc} can be as good as that of using Loss_{cc} —as if there is no annotation confusion.

Of course, there are more requirements to satisfy under Theorem 4.1. Particularly, the maximal $\log \det(\mathbf{M}^* (\mathbf{M}^*)^{\top})$ requirement is nontrivial. In practice, the optimization criterion in Theorem 4.1 can be approximated via a regularized version of Loss'_{cc} as follows:

$$\underset{\boldsymbol{\theta}, 0 \leq \mathbf{B} \leq 1}{\text{minimize}} \text{Loss}'_{\text{cc}} + \text{Loss}_{\text{vol}}, \quad (13)$$

where $\text{Loss}_{\text{vol}} = -\lambda \log \det(\mathbf{M} \mathbf{M}^{\top})$ and $\mathbf{m}_n = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{x}_n)$. Note that $\log \det(\mathbf{M} \mathbf{M}^{\top})$ is proportional to the volume of the Gram matrix $\mathbf{M} \mathbf{M}^{\top}$ (Boyd et al., 2004). Hence, the term can be regarded as a geometry-driven regularization. We name our method using (13) as the *Volume Maximization-Regularized Deep Constrained Clustering* (VolMaxDCC). An overall architecture is shown in Fig. 2.

5. Related Work

Recovering the underlying unseen matrix from incomplete and binary measurement is related to 1-bit low-rank matrix completion (Davenport et al., 2014), which was often studied in the context of recommender systems. The generative model in (5) is reminiscent of the MMSB (Airoldi et al., 2008; Huang & Fu, 2019) that has been a workhorse in overlapped community detection. The MMSB model with

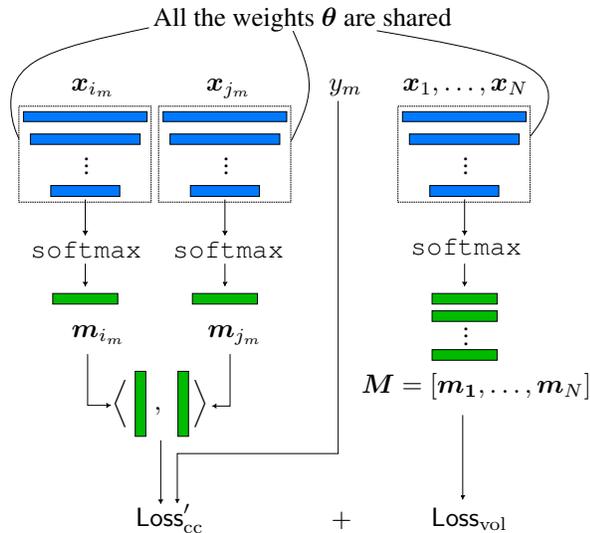


Figure 2: The architecture of the VolMaxDCC approach.

missing links was also used for clustering-related tasks, e.g., crowdclustering (Gomes et al., 2011). The ASC and SSC are commonly seen conditions in identifiability analysis of nonnegative matrix factorization (Fu et al., 2018; 2015a; 2019; 2016; Huang et al., 2014; Donoho & Stodden, 2003; Gillis, 2014; Gillis & Vavasis, 2014; Gillis & Luce, 2014; Gillis, 2020; Nguyen et al., 2022). Using volume maximization/minimization to enhance NMF identifiability appeared as early as 1994 (Craig, 1994) in the context of a blind source separation (BSS) problem in spectral image analysis. The volume-based geometric regularization was connected to ASC- and SSC-like conditions (e.g., the so-called “pure pixel condition” and “local dominance condition”) in (Chan et al., 2009) and (Fu et al., 2015b; 2016; Lin et al., 2015), respectively, to attain uniqueness of matrix factorization models, again, in the context of BSS. All these models do not involve nonlinear function learning or deep neural networks. In addition, the classic geometric factorization models were developed with continuous low-rank matrix data—instead of binary data generated from models involving complex *unknown* nonlinear function. Confusion matrices are often used in supervised noisy label learning and crowdsourcing (Dawid & Skene, 1979; Zhang et al., 2014; Rodrigues & Pereira, 2018), to model probabilistic label transition in the annotating process. The ASC and SSC were also used to establish identifiability in supervised (crowdsourced) noisy label learning (Xia et al., 2019; Li et al., 2021a; Ibrahim et al., 2019; Ibrahim & Fu, 2021; Ibrahim et al., 2023). Incorporating the idea of confusion matrix-based modeling in similarity annotation was not seen before. The proposed generative model connects label confusion matrices, volume maximization (more generally, ASC/SSC-based identifiable factor analysis), and DCC for the first time, to our best

knowledge.

6. Experiments

Datasets. We use STL-10 (Coates et al., 2011), ImageNet-10 (Chang et al., 2017a), and CIFAR-10 (Krizhevsky et al., 2009). For three datasets, we use $N = 10000$ samples as the seen data, and set aside 2000, 2000 and 45000 unseen data, respectively, for testing the generalization performance. In all experiments, $M = 10,000$ pairwise constraints are uniformly and randomly drawn from $[N] \times [N]$. There are no more than 0.02% of the total number of pairs annotated for all three datasets.

Baselines. We compare our method with several baselines, including the classic CC methods, i.e., PCKMeans (Basu et al., 2004a), COP-KMeans (Wagstaff et al., 2001), and the DCC methods, namely, DC-GMM (Manduchi et al., 2021) and C-IDEAC (Zhang et al., 2019; 2021). We also include the plain-vanilla K-means as a reference. We use a validation set for the baselines whenever proper for parameter tuning and algorithm stopping. The sizes of the validation sets are $N_{\text{valid}} = 1000$ for STL-10 and ImageNet-10 and $N_{\text{valid}} = 5000$ for CIFAR-10. For C-IDEAC, the regularization parameters is chosen from $\{0, 1e-1, 1e-2, 1e-3, 1e-4, 1e-5\}$. For DC-GMM, we use their heuristic to set the constraint violation penalty with the true (oracle) label flipping rate. For the proposed VolMaxDCC, we also choose λ among $\lambda \in \{0, 1e-1, 1e-2, 1e-3, 1e-4, 1e-5\}$. We also include the result of using the simple logistic loss Loss_{cc} in (3), which is referred to as VanillaDCC.

Neural Network Settings. For all the DCC methods, we employ the unsupervised pre-training method by (Li et al., 2021b) to convert the images to feature vectors $\{x_1, \dots, x_N\} \subseteq \mathbb{R}^{512}$ (Li et al., 2021b). The feature vectors are then fed to a two-hidden-layer fully connected neural network $f_{\theta}(\cdot)$, where each hidden layer has 512 ReLU activation functions. The output layer of $f_{\theta}(\cdot)$ has $K = 10$ dimensions with the softmax constraints. The classic methods also work with the pre-trained feature vectors.

Algorithm Implementation. To tackle the proposed criterion in (13), we first parameterize \mathbf{B} such that each element $B_{ij} = 1/(1 + \exp(-B'_{ij}))$, where $B'_{ij} \in \mathbb{R}$ is a trainable parameter. By doing so, (13) becomes an unconstrained optimization problem. We then employ the commonly used stochastic gradient-based solvers to tackle the re-parameterized problem. In our implementation, we use stochastic gradient descent with a batch size of 128. We set the learning rate for \mathbf{B}' and θ to be 0.1 and 0.5, respectively. The initialization of θ is chosen randomly following uniform distributions with parameters depending on output dimension of each layer. To initialize \mathbf{B}' , we make the

Table 1: Clustering performance of (seen data, unseen data) on STL10; $N_{\text{unseen}} = 2000$.

Noise level	Methods		Kmeans	COP-Kmeans	PCKmeans	DC-GMM	C-IDEAC	VanillaDCC	VolMaxDCC
	ACC	NMI							
0.0%	ACC	0.71, —	0.66, —	0.70, —	0.88, 0.87	0.89, 0.88	0.93, 0.91	0.91, 0.89	
	NMI	0.75, —	0.67, —	0.72, —	0.82, 0.80	0.81, 0.80	0.84, 0.81	0.82, 0.80	
	ARI	0.54, —	0.52, —	0.55, —	0.80, 0.78	0.79, 0.77	0.85, 0.83	0.84, 0.82	
8.3%	ACC	—	0.70, —	0.70, —	0.75, 0.76	0.77, 0.79	0.78, 0.80	0.80, 0.81	
	NMI	—	0.64, —	0.71, —	0.67, 0.69	0.67, 0.69	0.59, 0.62	0.64, 0.65	
	ARI	—	0.51, —	0.56, —	0.57, 0.59	0.59, 0.61	0.69, 0.71	0.73, 0.74	
10.3%	ACC	—	0.62, —	0.69, —	0.70, 0.72	0.70, 0.71	0.72, 0.73	0.79, 0.81	
	NMI	—	0.59, —	0.73, —	0.62, 0.64	0.60, 0.62	0.50, 0.51	0.68, 0.70	
	ARI	—	0.44, —	0.55, —	0.51, 0.52	0.50, 0.52	0.62, 0.64	0.77, 0.78	
15.0%	ACC	—	0.62, —	0.64, —	0.60, 0.61	0.57, 0.57	0.56, 0.58	0.79, 0.81	
	NMI	—	0.54, —	0.72, —	0.54, 0.55	0.50, 0.50	0.33, 0.35	0.68, 0.69	
	ARI	—	0.41, —	0.52, —	0.38, 0.39	0.38, 0.39	0.50, 0.51	0.76, 0.77	

Table 2: Clustering performance of (seen data, unseen data) on CIFAR10; $N_{\text{unseen}} = 45000$.

Noise level	Methods		Kmeans	COP-Kmeans	PCKmeans	DC-GMM	C-IDEAC	VanillaDCC	VolMaxDCC
	ACC	NMI							
0.0%	ACC	0.78, —	0.67, —	0.67, —	0.91, 0.89	0.90, 0.89	0.92, 0.90	0.91, 0.90	
	NMI	0.71, —	0.66, —	0.71, —	0.83, 0.81	0.83, 0.81	0.84, 0.80	0.83, 0.80	
	ARI	0.62, —	0.54, —	0.55, —	0.82, 0.79	0.81, 0.79	0.85, 0.81	0.84, 0.82	
4.9%	ACC	—	0.75, —	0.70, —	0.86, 0.86	0.86, 0.86	0.86, 0.86	0.86, 0.86	
	NMI	—	0.69, —	0.69, —	0.77, 0.77	0.77, 0.77	0.73, 0.73	0.73, 0.74	
	ARI	—	0.60, —	0.57, —	0.73, 0.74	0.73, 0.74	0.77, 0.77	0.77, 0.77	
8.7%	ACC	—	0.64, —	0.72, —	0.76, 0.76	0.76, 0.76	0.76, 0.77	0.83, 0.83	
	NMI	—	0.63, —	0.69, —	0.71, 0.71	0.70, 0.70	0.58, 0.59	0.70, 0.70	
	ARI	—	0.49, —	0.59, —	0.58, 0.59	0.59, 0.60	0.70, 0.70	0.75, 0.75	
10.9%	ACC	—	0.68, —	0.70, —	0.74, 0.74	0.73, 0.73	0.68, 0.69	0.82, 0.82	
	NMI	—	0.61, —	0.68, —	0.67, 0.68	0.65, 0.66	0.48, 0.49	0.68, 0.68	
	ARI	—	0.50, —	0.57, —	0.55, 0.55	0.56, 0.57	0.62, 0.63	0.74, 0.74	

diagonal elements to be 1 and the other elements -1 . The baselines are handled by their respective author-provided code for optimization.

6.1. Noisy Machine Annotations.

Noisy Annotation Settings. We first conduct experiments under noise-controlled settings. To be specific, we divide the experiments into two cases, where the annotations are accurate and noisy, respectively. In the former case, the annotations are set to 1 if the pair of data are from the same class, and set to 0 otherwise. In the latter case, to generate noisy pairwise constraints, several supervised machine annotators (i.e., classifiers) are trained using different number of training samples. By doing this, the trained classifiers have different prediction errors. The pairs are then annotated based on the class predictions made by these imperfect classifiers. If the pair of samples share the same predicted membership by the machine annotator, the annotation is set to be 1 (and 0 otherwise). The annotations are noisy as the machine annotators are far from perfect. The annotation noise level can be controlled by tuning the prediction error of the classifiers, via using various amounts of training samples. The statistics of the annotation errors are provided with our experiment results (see the “noise level” column in the tables).

Results. We report average performance on the seen set $\{x_1, \dots, x_N\}$ and the unseen data over 5 random trials in Tables 1, 2, and 3—which correspond to STL10, CIFAR10, and ImageNet10, respectively. For K-means, COP-Kmeans, and PCKMeans, we only report results on the seen set since

Table 3: Clustering performance of (seen data, unseen data) on ImageNet10; $N_{\text{unseen}} = 2000$.

Noise level	Methods		Kmeans	COP-Kmeans	PCKmeans	DC-GMM	C-IDEAC	VanillaDCC	VolMaxDCC
	ACC	NMI							
0.0%	ACC	0.85, —	0.79, —	0.70, —	0.97, 0.96	0.97, 0.96	0.97, 0.96	0.97, 0.96	
	NMI	0.80, —	0.77, —	0.75, —	0.93, 0.91	0.93, 0.92	0.94, 0.91	0.94, 0.91	
	ARI	0.68, —	0.66, —	0.55, —	0.94, 0.92	0.94, 0.92	0.93, 0.91	0.93, 0.91	
3.4%	ACC	—	0.79, —	0.66, —	0.93, 0.92	0.93, 0.93	0.92, 0.91	0.94, 0.94	
	NMI	—	0.75, —	0.73, —	0.86, 0.85	0.87, 0.86	0.83, 0.82	0.88, 0.87	
	ARI	—	0.65, —	0.52, —	0.84, 0.83	0.86, 0.85	0.84, 0.84	0.89, 0.88	
6.9%	ACC	—	0.74, —	0.72, —	0.84, 0.84	0.88, 0.88	0.84, 0.84	0.92, 0.91	
	NMI	—	0.70, —	0.76, —	0.79, 0.79	0.82, 0.82	0.70, 0.70	0.84, 0.83	
	ARI	—	0.58, —	0.59, —	0.58, —	0.71, 0.71	0.77, 0.77	0.88, 0.87	
11.2%	ACC	—	0.72, —	0.62, —	0.71, 0.72	0.80, 0.81	0.65, 0.66	0.91, 0.90	
	NMI	—	0.64, —	0.73, —	0.68, 0.70	0.74, 0.76	0.49, 0.52	0.83, 0.82	
	ARI	—	0.54, —	0.51, —	0.56, 0.58	0.66, 0.68	0.62, 0.63	0.87, 0.86	

Table 4: Clustering performance of (seen data, unseen data) on ImageNet10 with pairwise annotations acquired from the AMT platform; $N_{\text{unseen}} = 2000$, noise level: 23.09%.

Metrics	Methods		Kmeans	COP-Kmeans	PCKmeans	DC-GMM	C-IDEAC	VanillaDCC	VolMaxDCC
	ACC	NMI							
ACC	0.84, —	0.68, —	0.84, —	0.87, 0.85	0.91, 0.90	0.95, 0.94			
NMI	0.79, —	0.49, —	0.79, —	0.88, 0.86	0.86, 0.85	0.89, 0.87			
ARI	0.67, —	0.42, —	0.67, —	0.82, 0.79	0.83, 0.82	0.89, 0.88			

these methods do not have the notion of generalization. The performance is measured by three commonly seen metrics, namely, clustering accuracy (ACC) (Cai et al., 2010), normalized mutual information (NMI) (Cai et al., 2010), and adjusted rank index (ARI) (Yeung & Ruzso, 2001). For all metrics, a higher score indicates a better performance.

In case where accurate pairwise constraints are used (i.e., the rows in all the tables corresponding to “Noise Level = 0%”), most methods work reasonably well. As expected, all the DCC methods exhibit tangible edges over the CC methods that do not use deep neural networks. This is consistent with the observations made from previous DCC works (Manduchi et al., 2021; Zhang et al., 2019; 2021). The good performance of VanillaDCC on both training and testing set are also as expected, per our identifiability analysis.

The rows in the tables associated with nonzero noise levels show that the performance of DC-GMM, C-IDEAC, and VanillaDCC drops quickly. For example, in Table 2, the ACC of DC-GMM drops from (0.91,0.89) to (0.74,0.74) when the noise level changes from 0% to 10.9%. Similar performance degradation is observed for C-IDEAC and VanillaDCC, which are both DCC methods that do not explicitly consider annotation noise. Nonetheless, the proposed VolMaxDCC’s performance decline is much more graceful on all three datasets. In particular, Table 3 shows that the ACC of VolMaxDCC is still at (0.91,0.90) when the noise level reaches 11.2%, while the baselines have a best ACC of (0.80,0.81). The results show the usefulness of our confusion model, as well as the effectiveness of our identifiability-driven loss function design. More experiment results can be seen in Appendix F.

6.2. Noisy AMT Annotations.

Data Acquisition. In addition to using machine classifier-annotated data, we also conduct experiments using pairwise annotations that are obtained from the AMT platform. We uploaded 8994 data pairs to AMT, where the samples are from the ImageNet10 dataset. The annotators were asked to provide their judgement on the similarity of the pairs. Recall that there are $N = 10000$ samples in the ImageNet10 dataset, which means that 0.018% of all the pairs were annotated. We manually checked the error rate, which was found to be 23.09%. The annotated pairs are also released with the code.

Results. Table 4 shows the results on this AMT dataset. As before, we use the available pairs to learn the membership of training data and observe the testing accuracy over $N_{\text{unseen}} = 2,000$ samples. One can see that the proposed VolMaxDCC exhibits the highest clustering accuracy over the seen and unseen data. The clustering accuracy of the second best baseline is 4% lower than that of VolMaxDCC over both seen and unseen data. The margins of the proposed method over the baselines in terms of NMI and ARI are also obvious. The performance on the noisy AMT data speaks for the usefulness and effectiveness of the proposed method in real-world scenarios.

7. Conclusion

We revisited the pairwise annotation-based DCC problem from a membership identifiability viewpoint. We showed that a recently emerged logistic DCC loss is a sound criterion in terms of model identification—if the annotations are generated following a model that is reminiscent of the MMSB and deep learning based classifier learning. Based on our understanding to the vanilla logistic loss, we moved forward to consider the noisy annotation case and proposed a confusion-matrix based generative model. We proposed a modified logistic loss with a geometric regularization for provable membership identification—whose identifiability guarantee is the first of the kind under noisy annotations-based DCC, to our best knowledge. We tested our new design over various datasets under multiple noisy levels. We observed tangible improvements over all cases, showing our confusion-based modeling and identifiability-driven design are promising.

Limitations. The proposed approach has a couple of notable limitations. First, the model used for annotator noise relies on a confusion matrix model, assuming uniform confusion across all data samples, which may not always hold true in practical scenarios. Developing a framework that takes into account more realistic confusion models could lead to further improvements in performance. Second, the

establishment of membership identifiability under the SSC assumption lacks finite sample analysis (cf. Theorem 3.7 and Theorem 4.1). The assumption that M reaches infinity can never be met in practice. It is of great interest to show how the performance of the volume-based criterion scales with different sample sizes.

Acknowledgement. This work is supported in part by the National Science Foundation (NSF) under project NSF IIS-2007836.

References

- Airoldi, E. M., Blei, D., Fienberg, S., and Xing, E. Mixed membership stochastic blockmodels. *Advances in Neural Information Processing Systems*, 21, 2008.
- Anandkumar, A., Ge, R., Hsu, D., and Kakade, S. M. A tensor approach to learning mixed membership community models. *The Journal of Machine Learning Research*, 15 (1):2239–2312, 2014.
- Arora, S., Ge, R., Halpern, Y., Mimno, D., Moitra, A., Sontag, D., Wu, Y., and Zhu, M. A practical algorithm for topic modeling with provable guarantees. In *International Conference on Machine Learning*, pp. 280–288. PMLR, 2013.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. Spectrally-normalized margin bounds for neural networks. *Advances in Neural Information Processing Systems*, 30, 2017.
- Basu, S., Banerjee, A., and Mooney, R. J. Active semi-supervision for pairwise constrained clustering. In *Proceedings of the 2004 SIAM International Conference on Data Mining*, pp. 333–344, 2004a.
- Basu, S., Bilenko, M., and Mooney, R. J. A probabilistic framework for semi-supervised clustering. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 59–68, 2004b.
- Basu, S., Davidson, I., and Wagstaff, K. *Constrained clustering: Advances in algorithms, theory, and applications*. CRC Press, 2008.
- Bilenko, M., Basu, S., and Mooney, R. J. Integrating constraints and metric learning in semi-supervised clustering. In *Proceedings of the Twenty-First International Conference on Machine Learning*, pp. 11, 2004.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3 (Jan):993–1022, 2003.
- Boyd, S., Boyd, S. P., and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.

- Cai, D., He, X., and Han, J. Locally consistent concept factorization for document clustering. *IEEE Transactions on Knowledge and Data Engineering*, 23(6):902–913, 2010.
- Cai, T. T. and Zhou, W.-X. Matrix completion via max-norm constrained optimization. *Electronic Journal of Statistics*, 10(1):1493 – 1525, 2016.
- Caron, M., Bojanowski, P., Joulin, A., and Douze, M. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 132–149, 2018.
- Chan, T.-H., Chi, C.-Y., Huang, Y.-M., and Ma, W.-K. A convex analysis-based minimum-volume enclosing simplex algorithm for hyperspectral unmixing. *IEEE Transactions on Signal Processing*, 57(11):4418–4432, 2009.
- Chang, J., Wang, L., Meng, G., Xiang, S., and Pan, C. Deep adaptive image clustering. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5879–5887, 2017a.
- Chang, Y., Chen, J., Cho, M. H., Castaldi, P. J., Silverman, E. K., and Dy, J. G. Multiple clustering views from multiple uncertain experts. In *International Conference on Machine Learning*, pp. 674–683, 2017b.
- Chu, Z., Ma, J., and Wang, H. Learning from crowds by modeling common confusions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 5832–5840, 2021.
- Coates, A., Ng, A., and Lee, H. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15, pp. 215–223. PMLR, 11–13 Apr 2011.
- Covoes, T. F., Hruschka, E. R., and Ghosh, J. A study of K-means-based algorithms for constrained clustering. *Intelligent Data Analysis*, 17(3):485–505, 2013.
- Craig, M. D. Minimum-volume transforms for remotely sensed data. *IEEE Transactions on Geoscience and Remote Sensing*, 32(3):542–552, 1994.
- Cucuringu, M., Koutis, I., Chawla, S., Miller, G., and Peng, R. Simple and scalable constrained clustering: A generalized spectral method. In *Artificial Intelligence and Statistics*, pp. 445–454, 2016.
- Davenport, M. A., Plan, Y., van den Berg, E., and Wootters, M. 1-bit matrix completion. *Information and Inference: A Journal of the IMA*, 3(3):189–223, 2014.
- Davidson, I., Ravi, S., and Shamis, L. A SAT-based framework for efficient constrained clustering. In *Proceedings of the 2010 SIAM International Conference on Data Mining*, pp. 94–105. SIAM, 2010.
- Dawid, A. P. and Skene, A. M. Maximum likelihood estimation of observer error-rates using the EM algorithm. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 28(1):20–28, 1979.
- Donoho, D. and Stodden, V. When does non-negative matrix factorization give a correct decomposition into parts? In *Advances in Neural Information Processing Systems*, volume 16, pp. 1141–1148, 2003.
- Du, S. S., Zhai, X., Poczos, B., and Singh, A. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019.
- Fu, X., Ma, W.-K., Chan, T.-H., and Bioucas-Dias, J. M. Self-dictionary sparse regression for hyperspectral unmixing: Greedy pursuit and pure pixel search are related. *IEEE Journal of Selected Topics in Signal Processing*, 9(6):1128–1141, 2015a.
- Fu, X., Ma, W.-K., Huang, K., and Sidiropoulos, N. D. Blind separation of quasi-stationary sources: Exploiting convex geometry in covariance domain. *IEEE Transactions on Signal Processing*, 63:2306–2320, 2015b.
- Fu, X., Huang, K., Yang, B., Ma, W.-K., and Sidiropoulos, N. D. Robust volume minimization-based matrix factorization for remote sensing and document clustering. *IEEE Transactions on Signal Processing*, 64(23):6254–6268, 2016.
- Fu, X., Huang, K., and Sidiropoulos, N. D. On identifiability of nonnegative matrix factorization. *IEEE Signal Processing Letters*, 25(3):328–332, 2018.
- Fu, X., Huang, K., Sidiropoulos, N. D., and Ma, W.-K. Non-negative matrix factorization for signal and data analytics: Identifiability, Algorithms, and Applications. *IEEE Signal Process. Mag.*, 36(2):59–80, March 2019.
- Ghosh, A., Kale, S., and McAfee, P. Who moderates the moderators? crowdsourcing abuse detection in user-generated content. In *Proceedings of the 12th ACM Conference on Electronic Commerce*, pp. 167–176, 2011.
- Gillis, N. The why and how of nonnegative matrix factorization. In *Regularization, Optimization, Kernels, and Support Vector Machines*, pp. 275–310. Chapman and Hall/CRC, 2014.
- Gillis, N. *Nonnegative matrix factorization*. SIAM, 2020.

- Gillis, N. and Luce, R. Robust near-separable nonnegative matrix factorization using linear optimization. *The Journal of Machine Learning Research*, 15(1):1249–1280, 2014.
- Gillis, N. and Vavasis, S. A. Fast and robust recursive algorithms for separable nonnegative matrix factorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(4):698–714, 2014.
- Givoni, I. and Frey, B. Semi-supervised affinity propagation with instance-level constraints. In *Artificial Intelligence and Statistics*, pp. 161–168. PMLR, 2009.
- Gomes, R., Welinder, P., Krause, A., and Perona, P. Crowd-clustering. *Advances in Neural Information Processing Systems*, 24, 2011.
- Hsu, Y.-C. and Kira, Z. Neural network-based clustering using pairwise constraints. *arXiv preprint arXiv:1511.06321*, 2015.
- Hsu, Y.-C., Lv, Z., Schlosser, J., Odom, P., and Kira, Z. A probabilistic constrained clustering for transfer learning and image category discovery. *CVPR Deep-Vision workshop*, 2018.
- Huang, K. and Fu, X. Detecting overlapping and correlated communities without pure nodes: Identifiability and algorithm. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 2859–2868. PMLR, 09–15 Jun 2019.
- Huang, K., Sidiropoulos, N. D., and Swami, A. Non-negative matrix factorization revisited: Uniqueness and algorithm for symmetric decomposition. *IEEE Transactions on Signal Processing*, 62:211–224, 2014.
- Huang, K., Fu, X., and Sidiropoulos, N. D. Anchor-free correlated topic modeling: Identifiability and algorithm. *Advances in Neural Information Processing Systems*, 29, 2016.
- Ibrahim, S. and Fu, X. Crowdsourcing via annotator co-occurrence imputation and provable symmetric nonnegative matrix factorization. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pp. 4544–4554. PMLR, 18–24 Jul 2021.
- Ibrahim, S., Fu, X., Kargas, N., and Huang, K. Crowdsourcing via pairwise co-occurrences: Identifiability and algorithms. *Advances in Neural Information Processing Systems*, 32, 2019.
- Ibrahim, S., Nguyen, T., and Fu, X. Deep learning from crowdsourced labels: Coupled cross-entropy minimization, identifiability, and regularization. In *The Eleventh International Conference on Learning Representations*, 2023.
- Jain, A. K. and Dubes, R. C. *Algorithms for Clustering Data*. Prentice-Hall, Inc., USA, 1988. ISBN 013022278X.
- Kamvar, K., Sepandar, S., Klein, K., Dan, D., Manning, M., and Christopher, C. Spectral learning. In *International Joint Conference of Artificial Intelligence*. Stanford InfoLab, 2003.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Kulis, B., Basu, S., Dhillon, I., and Mooney, R. Semi-supervised graph clustering: A kernel approach. In *Proceedings of the 22nd International Conference on Machine Learning*, pp. 457–464, 2005.
- Kumar, A., Sindhwani, V., and Kambadur, P. Fast conical hull algorithms for near-separable non-negative matrix factorization. In *International Conference on Machine Learning*, pp. 231–239, 2013.
- Law, M. H. C., Topchy, A., and Jain, A. K. Model-based clustering with probabilistic constraints. In *Proceedings of the 2005 SIAM International Conference on Data Mining (SDM)*, pp. 641–645, 2005.
- Li, X., Liu, T., Han, B., Niu, G., and Sugiyama, M. Provably end-to-end label-noise learning without anchor points. In *International Conference on Machine Learning*, pp. 6403–6413. PMLR, 2021a.
- Li, Y., Hu, P., Liu, Z., Peng, D., Zhou, J. T., and Peng, X. Contrastive clustering. In *2021 AAAI Conference on Artificial Intelligence (AAAI)*, 2021b.
- Li, Z., Liu, J., and Tang, X. Constrained clustering via spectral regularization. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 421–428, 2009.
- Lin, C.-H., Ma, W.-K., Li, W.-C., Chi, C.-Y., and Ambikapathi, A. Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. *IEEE Transactions on Geoscience and Remote Sensing*, 53(10):5530–5546, 2015.
- Liu, H., Tao, Z., and Fu, Y. Partition level constrained clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(10):2469–2483, 2017.
- Liu, Q., Peng, J., and Ihler, A. T. Variational inference for crowdsourcing. *Advances in Neural Information Processing Systems*, 25, 2012.
- Lu, Z. and Carreira-Perpinan, M. A. Constrained spectral clustering through affinity propagation. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8, 2008.

- Luo, Y., Tian, T., Shi, J., Zhu, J., and Zhang, B. Semi-crowdsourced clustering with deep generative models. *Advances in Neural Information Processing Systems*, 31, 2018.
- Manduchi, L., Chin-Cheong, K., Michel, H., Wellmann, S., and Vogt, J. Deep conditional gaussian mixture model for constrained clustering. *Advances in Neural Information Processing Systems*, 34:11303–11314, 2021.
- Mao, X., Sarkar, P., and Chakrabarti, D. On mixed memberships and symmetric nonnegative matrix factorizations. In *International Conference on Machine Learning*, pp. 2324–2333, 2017.
- Nguyen, T., Fu, X., and Wu, R. Memory-efficient convex optimization for self-dictionary separable nonnegative matrix factorization: A frank–wolfe approach. *IEEE Transactions on Signal Processing*, 70:3221–3236, 2022.
- Panov, M., Slavnov, K., and Ushakov, R. Consistent estimation of mixed memberships with successive projections. *International Workshop on Complex Networks and their Applications*, pp. 53–64, 2017.
- Pelleg, D. and Baras, D. K-means with large and noisy constraint sets. In *European Conference on Machine Learning*, pp. 674–682. Springer, 2007.
- Ren, Y., Hu, K., Dai, X., Pan, L., Hoi, S. C., and Xu, Z. Semi-supervised deep embedded clustering. *Neurocomputing*, 325:121–130, 2019.
- Rodrigues, F. and Pereira, F. Deep learning from crowds. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Schultz, M. and Joachims, T. Learning a distance metric from relative comparisons. *Advances in Neural Information Processing Systems*, 16, 2003.
- Segal, E., Wang, H., and Koller, D. Discovering molecular pathways from protein interaction and gene expression data. *Bioinformatics*, 19(suppl_1):i264–i272, 2003.
- Shalev-Shwartz, S. and Ben-David, S. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- Wagstaff, K. and Cardie, C. Clustering with instance-level constraints. *AAAI/IAAI*, 1097:577–584, 2000.
- Wagstaff, K., Cardie, C., Rogers, S., Schrödl, S., et al. Constrained K-means clustering with background knowledge. In *International Conference on Machine Learning*, volume 1, pp. 577–584, 2001.
- Wang, X., Qian, B., and Davidson, I. On constrained spectral clustering and its applications. *Data Mining and Knowledge Discovery*, 28(1):1–30, 2014.
- Weyl, H. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen (mit einer anwendung auf die theorie der hohlraumstrahlung). *Mathematische Annalen*, 71(4):441–479, 1912.
- Xia, X., Liu, T., Wang, N., Han, B., Gong, C., Niu, G., and Sugiyama, M. Are anchor points really indispensable in label-noise learning? *Advances in Neural Information Processing Systems*, 32, 2019.
- Yang, B., Fu, X., Sidiropoulos, N. D., and Hong, M. Towards K-means-friendly spaces: Simultaneous deep learning and clustering. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 3861–3870. PMLR, 06–11 Aug 2017.
- Yeung, K. Y. and Ruzzo, W. L. Details of the adjusted rand index and clustering algorithms, supplement to the paper an empirical study on principal component analysis for clustering gene expression data. *Bioinformatics*, 17(9): 763–774, 2001.
- Yeung, M. M. and Yeo, B.-L. Time-constrained clustering for segmentation of video into story units. In *Proceedings of 13th International Conference on Pattern Recognition*, volume 3, pp. 375–380. IEEE, 1996.
- Yi, J., Jin, R., Jain, S., Yang, T., and Jain, A. Semi-crowdsourced clustering: Generalizing crowd labeling by robust distance metric learning. *Advances in Neural Information Processing Systems*, 25, 2012.
- Zhang, H., Basu, S., and Davidson, I. A framework for deep constrained clustering-algorithms and advances. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 57–72, 2019.
- Zhang, H., Zhan, T., Basu, S., and Davidson, I. A framework for deep constrained clustering. *Data Mining and Knowledge Discovery*, 35(2):593–620, 2021.
- Zhang, Y., Chen, X., Zhou, D., and Jordan, M. I. Spectral methods meet EM: A provably optimal algorithm for crowdsourcing. *Advances in Neural Information Processing Systems*, 27, 2014.
- Zhu, X., Loy, C. C., and Gong, S. Constrained clustering with imperfect oracles. *IEEE Transactions on Neural Networks and Learning Systems*, 27(6):1345–1357, 2015.

A. Proof of Lemma 3.3

Denote (I, J) and Y as two random variables (RVs) such that $\Pr(I, J, Y) = \Pr((I, J)) \Pr(Y|(I, J))$ and

$$\begin{cases} I, J \sim \mathcal{U}([N] \times [N]), \\ Y|(I, J) \sim \text{Bernoulli}(P_{IJ}^{\natural}), \end{cases}$$

where $(i_1, j_1, y_1), \dots, (i_M, j_M, y_M)$ are M i.i.d realizations of (I, J, Y) . Denote $\mathcal{D}_{\mathbf{P}^{\natural}}$ as the joint distribution of (I, J, Y) , i.e., $(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^{\natural}}$. We first note the following relationship, which was also used in 1-bit matrix completion literature (Davenport et al., 2014; Cai & Zhou, 2016):

$$\begin{aligned} D_{\text{kl}}(\mathbf{P}^{\natural} \parallel \mathbf{P}^{\star}) &= \frac{1}{N^2} \sum_{(i,j) \in [N]^2} d_{\text{kl}}(P_{ij}^{\natural} \parallel P_{ij}^{\star}) \quad (\text{by definition in (28)}) \\ &= \mathbb{E}_{(I,J) \sim \mathcal{U}} \left[d_{\text{kl}}(P_{IJ}^{\natural} \parallel P_{IJ}^{\star}) \right] \\ &= \mathbb{E}_{(I,J) \sim \mathcal{U}} \left[P_{IJ}^{\natural} \log \frac{P_{IJ}^{\natural}}{P_{IJ}^{\star}} + (1 - P_{IJ}^{\natural}) \log \frac{1 - P_{IJ}^{\natural}}{1 - P_{IJ}^{\star}} \right] \quad (\text{by definition in (29)}) \\ &= \mathbb{E}_{(I,J) \sim \mathcal{U}} \left[\mathbb{E}_{Y|(I,J) \sim \text{Bern}(P_{IJ}^{\natural})} \left[Y \log \frac{P_{IJ}^{\natural}}{P_{IJ}^{\star}} + (1 - Y) \log \frac{1 - P_{IJ}^{\natural}}{1 - P_{IJ}^{\star}} \right] \right] \\ &= \mathbb{E}_{(I,J,Y) \sim \mathcal{D}_{\mathbf{P}^{\natural}}} \left[Y \log P_{IJ}^{\natural} + (1 - Y) \log P_{IJ}^{\natural} \right] - \mathbb{E}_{(I,J,Y) \sim \mathcal{D}_{\mathbf{P}^{\natural}}} \left[Y \log P_{IJ}^{\star} + (1 - Y) \log P_{IJ}^{\star} \right] \\ &= L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\star}) - L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\natural}), \end{aligned} \quad (14)$$

where we define $L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P})$ for any matrix $0 \leq \mathbf{P} \leq 1$ as

$$\begin{aligned} L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}) &\triangleq \mathbb{E}_{(I,J,Y) \sim \mathcal{D}_{\mathbf{P}^{\natural}}} [\ell(\mathbf{P}, (I, J, Y))], \quad \text{where} \\ \ell(\mathbf{P}, (I, J, Y)) &\triangleq -Y \log P_{IJ} - (1 - Y) \log(1 - P_{IJ}). \end{aligned} \quad (15)$$

Recall $S \triangleq \{(i_1, j_1, y_1), \dots, (i_M, j_M, y_M)\}$. Define $L_S(\mathbf{P}) \triangleq (1/M) \sum_{i=1}^M \ell(\mathbf{P}, (i_m, j_m, y_m))$. By this definition, \mathbf{P}^{\star} is an optimal solution of the following problem

$$\underset{\mathbf{P} \in \mathcal{P}_{\mathcal{F}, \mathbf{X}}}{\text{minimize}} \quad L_S(\mathbf{P}),$$

where $\mathcal{P}_{\mathcal{F}, \mathbf{X}}$ is a set of matrices defined by the function class \mathcal{F} and the data set $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$. Specifically, it is defined as

$$\mathcal{P}_{\mathcal{F}, \mathbf{X}} \triangleq \{ \mathbf{P} \in \mathbb{R}^{N \times N} \mid P_{ij} = \text{dot}(\mathbf{f}'(\mathbf{x}_i), \mathbf{x}_j) \}, \quad (16a)$$

$$\text{dot}([\mathbf{u}; \mathbf{v}]) \triangleq \mathbf{u}^{\top} \mathbf{v}, \quad (16b)$$

$$\mathbf{f}'(\mathbf{x}, \mathbf{y}) \triangleq [\mathbf{f}(\mathbf{x}); \mathbf{f}(\mathbf{y})], \quad \mathbf{f} \in \mathcal{F}. \quad (16c)$$

Define $\tilde{\mathbf{P}} \in \mathbb{R}^{N \times N}$ such that $\tilde{P}_{ij} = \tilde{\mathbf{f}}(\mathbf{x}_i)^{\top} \tilde{\mathbf{f}}(\mathbf{x}_j)$ where $\tilde{\mathbf{f}}$ is defined in Assumption 7, i.e., the ‘‘best approximation’’ for \mathbf{f}^{\natural} by the class \mathcal{F} . Using (14) and the above notations, we have

$$\begin{aligned} D_{\text{kl}}(\mathbf{P}^{\natural} \parallel \mathbf{P}^{\star}) &= L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\star}) - L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\natural}) \\ &= (L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\star}) - L_S(\mathbf{P}^{\star})) + (L_S(\mathbf{P}^{\star}) - L_S(\tilde{\mathbf{P}})) + (L_S(\tilde{\mathbf{P}}) - L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\tilde{\mathbf{P}})) + (L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\tilde{\mathbf{P}}) - L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\natural})) \\ &\leq (L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\star}) - L_S(\mathbf{P}^{\star})) + (L_S(\tilde{\mathbf{P}}) - L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\tilde{\mathbf{P}})) + (L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\tilde{\mathbf{P}}) - L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\natural})) \\ &\leq 2 \sup_{\mathbf{P} \in \mathcal{P}_{\mathcal{F}, \mathbf{X}}} |L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}) - L_S(\mathbf{P})| + |L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\tilde{\mathbf{P}}) - L_{\mathcal{D}_{\mathbf{P}^{\natural}}}(\mathbf{P}^{\natural})|, \end{aligned} \quad (17)$$

where the first inequality holds because $L_S(\mathbf{P}^{\star}) \leq L_S(\tilde{\mathbf{P}})$ which holds by the definition of \mathbf{P}^{\star} , and the second inequality holds as both $\mathbf{P}^{\star}, \tilde{\mathbf{P}} \in \mathcal{P}_{\mathcal{F}, \mathbf{X}}$.

Bound the first term in (17). Notice that $L_{\mathcal{D}_{\mathbf{P}^\natural}}(\mathbf{P})$ and $L_S(\mathbf{P})$ are the expected loss and the empirical loss, respectively. Hence, the difference between these two quantities can be bounded using generalization analysis via concentration inequalities.

To this end, we first determine an upper bound on the loss function $\ell(\mathbf{P}, (i, j, y))$. Under the assumption stated in Assumption 3.2,

$$\alpha < P_{ij} < 1 - \alpha, \quad \forall (i, j) \in [N] \times [N],$$

which leads to

$$\ell(\mathbf{P}, (i, j, y)) = y \log \frac{1}{P_{ij}} + (1-y) \log \frac{1}{1-P_{ij}} \leq \log \frac{1}{P_{ij}} + \log \frac{1}{1-P_{ij}} \leq 2 \log \frac{1}{\alpha}, \quad \forall (i, j) \in [N]^2, y \in \{0, 1\}, \quad (18)$$

As $\ell(\mathbf{P}, (i, j, y))$ is bounded, the sample set S contains M samples of (i_m, j_m, y_m) 's, where each sample (i_m, j_m, y_m) is generated independently with the same distribution $\mathcal{D}_{\mathbf{P}^\natural}$, one can readily apply the following generalization bound (Shalev-Shwartz & Ben-David, 2014, Theorem 26.5):

$$\begin{aligned} \sup_{\mathbf{P} \in \mathcal{P}_{\mathcal{F}, \mathbf{X}}} |L_{\mathcal{D}_{\mathbf{P}^\natural}}(\mathbf{P}) - L_S(\mathbf{P})| &\leq 2\mathcal{R}(\ell \circ \mathcal{P}_{\mathcal{F}, \mathbf{X}} \circ S) + 8 \log(1/\alpha) \sqrt{\frac{2 \log(4/\delta)}{M}} \\ &\leq \frac{8 \log(1/\alpha)}{M} + \frac{12\sqrt{2} \log M}{\alpha M \log 2} \|\mathbf{S}_{\mathbf{X}}\|_{\mathbb{F}} \sqrt{R_{\text{NET}}} + 8 \log(1/\alpha) \sqrt{\frac{2 \log(4/\delta)}{M}} \end{aligned} \quad (19)$$

holds with probability at least $1 - \delta$, where we have applied Lemma A.3 to reach the final bound.

Bound the second term in (17). We have

$$\begin{aligned} |L_{\mathcal{D}_{\mathbf{P}^\natural}}(\tilde{\mathbf{P}}) - L_{\mathcal{D}_{\mathbf{P}^\natural}}(\mathbf{P}^\natural)| &= \left| \mathbb{E}_{(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^\natural}} \left[\ell(\tilde{\mathbf{P}}, (I, J, Y)) - \ell(\mathbf{P}^\natural, (I, J, Y)) \right] \right| \\ &= \left| \mathbb{E}_{(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^\natural}} \left[Y \left(\log P_{IJ}^\natural - \log \tilde{P}_{IJ} \right) + (1-Y) \left(\log(1 - P_{IJ}^\natural) - \log(1 - \tilde{P}_{IJ}) \right) \right] \right| \\ &\leq \mathbb{E}_{(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^\natural}} \left[\left| Y \left(\log P_{IJ}^\natural - \log \tilde{P}_{IJ} \right) + (1-Y) \left(\log(1 - P_{IJ}^\natural) - \log(1 - \tilde{P}_{IJ}) \right) \right| \right] \\ &\leq \mathbb{E}_{(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^\natural}} \left[Y \left| \log \frac{P_{IJ}^\natural}{\tilde{P}_{IJ}} \right| + (1-Y) \left| \log \frac{1 - P_{IJ}^\natural}{1 - \tilde{P}_{IJ}} \right| \right] \\ &\leq \mathbb{E}_{(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^\natural}} \left[\left| \log \frac{P_{IJ}^\natural}{\tilde{P}_{IJ}} \right| + \left| \log \frac{1 - P_{IJ}^\natural}{1 - \tilde{P}_{IJ}} \right| \right] \quad (\text{since } Y \in \{0, 1\}) \\ &\leq^{(a)} \mathbb{E}_{(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^\natural}} \left[\left| \frac{P_{IJ}^\natural - \tilde{P}_{IJ}}{\tilde{P}_{IJ}} \right| + \left| \frac{\tilde{P}_{IJ} - P_{IJ}^\natural}{1 - \tilde{P}_{IJ}} \right| \right] \\ &\leq \frac{1}{\alpha} \mathbb{E}_{(I, J, Y) \sim \mathcal{D}_{\mathbf{P}^\natural}} \left[\left| P_{IJ}^\natural - \tilde{P}_{IJ} \right| \right], \end{aligned} \quad (20)$$

where ^(a) holds by applying inequality $\log(x+1) \leq x$, $\forall x > -1$.

In order to bound (20), consider arbitrary $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,

$$\begin{aligned}
 & \left| \mathbf{f}^{\natural}(\mathbf{x})^\top \mathbf{f}^{\natural}(\mathbf{y}) - \tilde{\mathbf{f}}(\mathbf{x})^\top \tilde{\mathbf{f}}(\mathbf{y}) \right| \\
 &= \left| \left(\mathbf{f}^{\natural}(\mathbf{x}) - \tilde{\mathbf{f}}(\mathbf{x}) + \tilde{\mathbf{f}}(\mathbf{x}) \right)^\top \left(\mathbf{f}^{\natural}(\mathbf{y}) - \tilde{\mathbf{f}}(\mathbf{y}) + \tilde{\mathbf{f}}(\mathbf{y}) \right) - \tilde{\mathbf{f}}(\mathbf{x})^\top \tilde{\mathbf{f}}(\mathbf{y}) \right| \\
 &= \left| \left(\mathbf{f}^{\natural}(\mathbf{x}) - \tilde{\mathbf{f}}(\mathbf{x}) \right)^\top \left(\mathbf{f}^{\natural}(\mathbf{y}) - \tilde{\mathbf{f}}(\mathbf{y}) \right) + \left(\mathbf{f}^{\natural}(\mathbf{x}) - \tilde{\mathbf{f}}(\mathbf{x}) \right)^\top \tilde{\mathbf{f}}(\mathbf{y}) + \tilde{\mathbf{f}}(\mathbf{x})^\top \left(\mathbf{f}^{\natural}(\mathbf{y}) - \tilde{\mathbf{f}}(\mathbf{y}) \right) \right| \\
 &\leq \left\| \mathbf{f}^{\natural}(\mathbf{x}) - \tilde{\mathbf{f}}(\mathbf{x}) \right\| \left\| \mathbf{f}^{\natural}(\mathbf{y}) - \tilde{\mathbf{f}}(\mathbf{y}) \right\| + \left\| \mathbf{f}^{\natural}(\mathbf{x}) - \tilde{\mathbf{f}}(\mathbf{x}) \right\| \left\| \tilde{\mathbf{f}}(\mathbf{y}) \right\| + \left\| \tilde{\mathbf{f}}(\mathbf{x}) \right\| \left\| \mathbf{f}^{\natural}(\mathbf{y}) - \tilde{\mathbf{f}}(\mathbf{y}) \right\| \\
 &\leq \nu^2 + 2\nu \quad (\text{by Eq. (7)}) \\
 &< 4\nu, \quad (\text{since } \nu \leq \sqrt{2} \text{ due to the fact that all } \mathbf{f}(\mathbf{x}) \text{ are in the probability simplex})
 \end{aligned} \tag{21}$$

where we have used $\|\tilde{\mathbf{f}}(\mathbf{x})\|_2 \leq \|\tilde{\mathbf{f}}(\mathbf{x})\|_1 = 1$ in the second inequality, as the neural network uses a softmax output layer.

Substituting (21) into (20), we have

$$L_S(\tilde{\mathbf{P}}) - L_S(\mathbf{P}^{\natural}) \leq \frac{4\nu}{\alpha}. \tag{22}$$

Putting together. Combining (17), (19), (22) gives

$$D_{\text{kl}}(\mathbf{P}^{\natural} \parallel \mathbf{P}^*) \leq \frac{16 \log(1/\alpha)}{M} + \frac{24\sqrt{2} \log M}{\alpha M \log 2} \|\mathbf{S}_{\mathbf{X}}\|_{\text{F}} \sqrt{R_{\text{NET}}} + 16 \log(1/\alpha) \sqrt{\frac{2 \log(4/\delta)}{M}} + \frac{4\nu}{\alpha}$$

By Lemma A.4, we can upper bound our quantity of interest by the following:

$$\|\mathbf{P}^* - \mathbf{P}^{\natural}\|_{\text{F}}^2 \leq 4N^2 D_{\text{kl}}(\mathbf{P}^{\natural} \parallel \mathbf{P}^*). \tag{23}$$

That is,

$$\frac{1}{N^2} \|\mathbf{P}^* - \mathbf{P}^{\natural}\|_{\text{F}}^2 \leq \frac{64 \log(1/\alpha)}{M} + \frac{96\sqrt{2} \log M}{\alpha M \log 2} \|\mathbf{S}_{\mathbf{X}}\|_{\text{F}} \sqrt{R_{\text{NET}}} + 64 \log(1/\alpha) \sqrt{\frac{2 \log(4/\delta)}{M}} + \frac{16\nu}{\alpha}$$

hold with probability at least $1 - \delta$ over $S = \{(i_1, j_1, y_1) \dots, (i_M, j_M, y_M)\}$.

A.1. Supporting lemmas for the proof of Theorem 3.3

A.1.1. RADEMACHER COMPLEXITY BOUND FOR $\mathcal{P}_{\mathcal{F}, \mathbf{X}}$

Definition A.1 (Rademacher complexity). Let $A \subset \mathbb{R}^m$ be a set of vectors. We define Rademacher complexity of A as

$$R(A) \triangleq \frac{1}{m} \mathbb{E} \left[\sup_{\mathbf{a} \in A} \sum_{i=1}^m \sigma_i a_i \right],$$

where $\sigma_1, \dots, \sigma_m \in \{-1, 1\}$ are i.i.d. distributed according to $\Pr(\sigma_i = 1) = \Pr(\sigma_i = -1) = 0.5$.

Definition A.2 (Covering number). Let $A \subset \mathbb{R}^m$ be a set of vectors. We say that A is r -covered by a set A' , with respect to metric d if for all $\mathbf{a} \in A$, there exists $\mathbf{a}' \in A'$ with $d(\mathbf{a}, \mathbf{a}') \leq r$. We define by $\mathcal{N}(r, A, d)$ the cardinality of the smallest A' that r -covers A .

Lemma A.3 (Rademacher complexity bound). Consider function class \mathcal{F} satisfying Assumptions 3.2, and that $\alpha < \mathbf{f}(\mathbf{x})^\top \mathbf{f}(\mathbf{y}) < 1 - \alpha, \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}, \forall \mathbf{f} \in \mathcal{F}$. Then, we have

$$\mathcal{R}(\ell \circ \mathcal{P}_{\mathcal{F}, \mathbf{X}} \circ S) \leq \frac{4 \log(1/\alpha)}{M} + \frac{6\sqrt{2} \log M}{\alpha M \log 2} \|\mathbf{S}_{\mathbf{X}}\|_{\text{F}} \sqrt{R_{\text{NET}}},$$

where $\ell \circ \mathcal{P}_{\mathcal{F}, \mathbf{X}} \circ S \triangleq \{\ell(\mathbf{P}, (i_1, j_1, y_1)), \dots, \ell(\mathbf{P}, (i_M, j_M, y_M))\} \in \mathbb{R}^M \mid \mathbf{P} \in \mathcal{P}_{\mathcal{F}, \mathbf{X}}\}$, $\mathcal{R}(A)$ denotes Rademacher complexity of set $A \subseteq \mathbb{R}^M$, $\mathcal{P}_{\mathcal{F}, \mathbf{X}}$ was defined in (16), $S = \{(i_1, j_1, y_1), \dots, (i_M, j_M, y_M)\}$, and the loss function ℓ was defined in (15).

Proof. We start with deriving the covering number of an ϵ -net of the set $\ell \circ \mathcal{P}_{\mathbf{X}} \circ S$:

$$\begin{aligned} \mathcal{N}(\epsilon, \ell \circ \mathcal{P}_{\mathbf{X}} \circ S, \|\cdot\|) &= \mathcal{N}(\epsilon, \ell \circ \text{dot} \circ \mathbf{f}' \circ S, \|\cdot\|) \\ &\leq \mathcal{N}(\epsilon/L_\ell, \text{dot} \circ \mathbf{f}' \circ S, \|\cdot\|) \end{aligned} \quad (24)$$

$$\leq \mathcal{N}(\epsilon/(2L_\ell), \mathbf{f}' \circ S, \|\cdot\|_{\text{F}}) \quad (25)$$

$$\leq \mathcal{N}(\epsilon/(L_\ell\sqrt{2}), \mathbf{f} \circ S_{\mathbf{X}}^{(1)}, \|\cdot\|_{\text{F}}) \mathcal{N}(\epsilon/(L_\ell\sqrt{2}), \mathbf{f} \circ S_{\mathbf{X}}^{(2)}, \|\cdot\|_{\text{F}}) \quad (26)$$

$$= \mathcal{N}(\epsilon/(L_\ell\sqrt{2}), \text{softmax} \circ \text{net} \circ S_{\mathbf{X}}^{(1)}, \|\cdot\|_{\text{F}}) \mathcal{N}(\epsilon/(L_\ell\sqrt{2}), \text{softmax} \circ \text{net} \circ S_{\mathbf{X}}^{(2)}, \|\cdot\|_{\text{F}})$$

$$\leq \mathcal{N}(\epsilon/(L_\ell\sqrt{2}), \text{net} \circ S_{\mathbf{X}}^{(1)}, \|\cdot\|_{\text{F}}) \mathcal{N}(\epsilon/(L_\ell\sqrt{2}), \text{net} \circ S_{\mathbf{X}}^{(2)}, \|\cdot\|_{\text{F}})$$

$$\leq \exp\left(\frac{\|S_{\mathbf{X}}^{(1)}\|_{\text{F}}^2}{\epsilon^2/(2L_\ell^2)} R_{\text{NET}} + \frac{\|S_{\mathbf{X}}^{(2)}\|_{\text{F}}^2}{\epsilon^2/(2L_\ell^2)} R_{\text{NET}}\right) \quad (\text{By Lemma A.10})$$

$$= \exp\left(\left(\|S_{\mathbf{X}}^{(1)}\|_{\text{F}}^2 + \|S_{\mathbf{X}}^{(2)}\|_{\text{F}}^2\right) \frac{2R_{\text{NET}}L_\ell^2}{\epsilon^2}\right) \quad (27)$$

Eq. (24) holds by applying Lemma A.7 and the fact that function $\phi_y(P_W) \triangleq -y \log(P_W) - (1-y) \log(1-P_W)$ is L_ℓ -Lipschitz continuous. Specifically, $L_\ell = (1/\alpha)$ since for $\alpha \leq p_1, p_2 \leq 1-\alpha$,

$$\begin{aligned} |\phi_y(p_1) - \phi_y(p_2)| &= \left| -y \log \frac{p_1}{p_2} - (1-y) \log \frac{1-p_1}{1-p_2} \right| \\ &\leq \left| \log \frac{p_1}{p_2} \right| + \left| \log \frac{1-p_1}{1-p_2} \right| \\ &\leq \left| \frac{p_1-p_2}{p_2} \right| + \left| \frac{p_2-p_1}{1-p_2} \right| \quad (\text{by } \log(1+x) \leq x) \\ &\leq \frac{1}{\alpha} |p_1 - p_2| = L_\ell |p_1 - p_2|. \end{aligned}$$

Similarly, Eq. (25) holds by Lemma A.7 and the fact that for $\forall \mathbf{u}, \mathbf{v} \in \{\mathbf{x} \mid \mathbf{x} \geq 0, \mathbf{1}^\top \mathbf{x} = 1\}$, function $\text{dot}([\mathbf{u}; \mathbf{v}])$ defined in (16b) is 2-Lipschitz continuous. To see this, consider the following:

$$\|\nabla \text{dot}([\mathbf{u}; \mathbf{v}])\| = \left\| \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \right\| \leq 2.$$

Eq. (26) holds as a consequence of applying Lemma A.9 to the set $\mathbf{f}' \circ S$ which is

$$\begin{aligned} \mathbf{f}' \circ S &= \left\{ [\mathbf{Z}_1, \mathbf{Z}_2] \in \mathbb{R}^{M \times 2K} \mid \mathbf{Z}_1 \in \mathbf{f} \circ S_{\mathbf{X}}^{(1)}, \mathbf{Z}_2 \in \mathbf{f} \circ S_{\mathbf{X}}^{(2)} \right\}, \\ S_{\mathbf{X}}^{(1)} &= [\mathbf{x}_{\omega_1^1}, \mathbf{x}_{\omega_2^1}, \dots, \mathbf{x}_{\omega_M^1}] \in \mathbb{R}^{D \times M}, \\ S_{\mathbf{X}}^{(2)} &= [\mathbf{x}_{\omega_1^2}, \mathbf{x}_{\omega_2^2}, \dots, \mathbf{x}_{\omega_M^2}] \in \mathbb{R}^{D \times M}. \end{aligned}$$

Now we can proceed to bound $\mathcal{R}(\ell \circ \mathcal{P}_{\mathbf{X}} \circ S)$ with the help of Dudley's entropy integral.

Apply Lemma A.6 onto the set $\ell \circ \mathcal{P}_{\mathbf{X}} \circ S$ with $c = 2 \log(1/\alpha) \sqrt{M}$ (by (18)), for any integer $T > 0$,

$$\begin{aligned}
 \mathcal{R}(\ell \circ \mathcal{P}_{\mathbf{X}} \circ S) &\leq \frac{c2^{-T}}{\sqrt{M}} + \frac{6c}{M} \sum_{t=1}^T 2^{-t} \sqrt{\log(\mathcal{N}(c2^{-t}, A))} \quad (\text{By Lemma A.6}) \\
 &\leq \frac{c2^{-T}}{\sqrt{M}} + \frac{6c}{M} \sum_{t=1}^T 2^{-t} \sqrt{\left(\|\mathbf{S}_{\mathbf{X}}^{(1)}\|_{\text{F}}^2 + \|\mathbf{S}_{\mathbf{X}}^{(2)}\|_{\text{F}}^2 \right) \frac{2R_{\text{NET}}L_{\ell}^2}{c^22^{-2t}}} \quad (\text{By (27)}) \\
 &\leq \frac{c2^{-T}}{\sqrt{M}} + \frac{6T}{M} \sqrt{\left(\|\mathbf{S}_{\mathbf{X}}^{(1)}\|_{\text{F}}^2 + \|\mathbf{S}_{\mathbf{X}}^{(2)}\|_{\text{F}}^2 \right) 2R_{\text{NET}}L_{\ell}^2} \\
 &= \frac{c2^{-T}}{\sqrt{M}} + \frac{6\sqrt{2}T}{M} \|\mathbf{S}_{\mathbf{X}}\|_{\text{F}} L_{\ell} \sqrt{R_{\text{NET}}} \\
 &= 2 \log(1/\alpha) 2^{-T} + \frac{6\sqrt{2}T}{M} \|\mathbf{S}_{\mathbf{X}}\|_{\text{F}} L_{\ell} \sqrt{R_{\text{NET}}} \quad (\text{substitute } c).
 \end{aligned}$$

Lastly, choosing $T = \lfloor \log_2 M \rfloor$ and substituting $L_{\ell} = 1/\alpha$ concludes the proof,

$$\begin{aligned}
 \mathcal{R}(\ell \circ \mathcal{P}_{\mathcal{F}, \mathbf{X}} \circ S) &\leq 2 \log(1/\alpha) 2^{-\lfloor \log_2 M \rfloor} + \frac{6\sqrt{2} \lfloor \log_2 M \rfloor}{M} \|\mathbf{S}_{\mathbf{X}}\|_{\text{F}} L_{\ell} \sqrt{R_{\text{NET}}} \\
 &\leq \frac{4 \log(1/\alpha)}{M} + \frac{6\sqrt{2} \log M}{\alpha M \log 2} \|\mathbf{S}_{\mathbf{X}}\|_{\text{F}} \sqrt{R_{\text{NET}}}.
 \end{aligned}$$

This concludes the proof. □

A.1.2. SOME OTHERS LEMMAS

Lemma A.4. For $0 \leq \mathbf{P}, \mathbf{Q} \leq 1$ with the same size $N \times N$,

$$\|\mathbf{P} - \mathbf{Q}\|_{\text{F}}^2 \leq 4N^2 D_H^2(\mathbf{P}, \mathbf{Q}) \leq 4N^2 D_{\text{kl}}(\mathbf{Q} \|\mathbf{P}),$$

where

$$\begin{aligned}
 D_H^2(\mathbf{P}, \mathbf{Q}) &\triangleq \frac{1}{N^2} \sum_{\omega \in [N] \times [N]} d_H^2(P_{\omega}, Q_{\omega}), \\
 D_{\text{kl}}(\mathbf{P} \|\mathbf{Q}) &\triangleq \frac{1}{N^2} \sum_{\omega \in [N] \times [N]} d_{\text{kl}}(P_{\omega} \|\mathbf{Q}_{\omega}),
 \end{aligned} \tag{28}$$

and $d_H(p, q)$, $d_{\text{kl}}(p, q)$ are typical Hellinger distance and KL divergence for $0 \leq p, q \leq 1$, resp, i.e.,

$$\begin{aligned}
 d_H^2(p, q) &\triangleq (\sqrt{p} - \sqrt{q})^2 + (\sqrt{1-p} - \sqrt{1-q})^2, \\
 d_{\text{kl}}(p, q) &\triangleq p \log \left(\frac{p}{q} \right) + (1-p) \log \left(\frac{1-p}{1-q} \right).
 \end{aligned} \tag{29}$$

Proof. For the first inequality, we use the following result.

Lemma A.5 (Lemma 2, scalar version, (Davenport et al., 2014)). Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be any differential function, and s, t are two real numbers satisfying $|s|, |t| \leq \alpha$. Then

$$d_H^2(f(s); f(t)) \geq \inf_{|x| \leq \alpha} \frac{f'(x)^2}{8f(x)(1-f(x))} (s-t)^2.$$

As a result, for $f(x) = x$, $0 \leq p, q \leq 1$,

$$d_H^2(p, q) \geq \inf_{|x| \leq 1} \frac{1}{8x(1-x)} (p-q)^2 = \frac{1}{4} (p-q)^2.$$

Summing across all elements of P, Q leads to,

$$\|P - Q\|_F^2 \leq 4N^2 D_H^2(P, Q).$$

The second inequality is simply derived from Jensen's inequality and the fact that $1 - x \leq -\log x, \forall x > 0$. \square

Lemma A.6 (Dudley entropy integral). ((Shalev-Shwartz & Ben-David, 2014, Lemma 27.4)) Let $A \subset \mathbb{R}^m$, $c = \min_{\bar{\mathbf{a}} \in \mathbb{R}^m} \max_{\mathbf{a} \in A} \|\bar{\mathbf{a}} - \mathbf{a}\|$. Then, for any integer $T > 0$,

$$R(A) \leq \frac{c2^{-T}}{\sqrt{m}} + \frac{6c}{m} \sum_{k=1}^T 2^{-k} \sqrt{\log(N(c2^{-k}, A))}.$$

Lemma A.7. For each $i \in [m]$, let $\phi_i : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$ be a ρ -Lipschitz function; namely, for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{d_1}$ we have $\|\phi_i(\mathbf{x}) - \phi_i(\mathbf{y})\| \leq \rho \|\mathbf{x} - \mathbf{y}\|$. For $\mathbf{A} \in \mathbb{R}^{d_1 \times m}$ let $\Phi(\mathbf{A})$ denote the matrix $[\phi_1(\mathbf{a}_1), \dots, \phi_m(\mathbf{a}_m)]$. For $\mathcal{A} \subset \mathbb{R}^{d_1 \times m}$, let $\Phi \circ \mathcal{A} = \{\Phi(\mathbf{A}) \mid \mathbf{A} \in \mathcal{A}\} \subset \mathbb{R}^{d_2 \times m}$. Then,

$$\mathcal{N}(\rho r, \Phi \circ \mathcal{A}, \|\cdot\|_F) \leq \mathcal{N}(r, \mathcal{A}, \|\cdot\|_F). \quad (30)$$

Remark A.8. Firstly, the choice of Frobenius norm in (30) is tied to ℓ_2 -norm defining the ρ -Lipschitz function. Secondly, Lemma A.7 is a natural generalization of (Shalev-Shwartz & Ben-David, 2014, Lemma 27.3).

Proof. Define \mathcal{A}' as one of the smallest r -cover of \mathcal{A} , i.e., for all $\mathbf{A} \in \mathcal{A}$ there exists $\mathbf{A}' \in \mathcal{A}'$ such that $\|\mathbf{A} - \mathbf{A}'\|_F \leq r$. By definition, $|\mathcal{A}'| = \mathcal{N}(r, \mathcal{A}, \|\cdot\|_F)$. Since

$$\|\Phi(\mathbf{A}) - \Phi(\mathbf{A}')\|_F^2 = \sum_{i=1}^m \|\phi_i(\mathbf{a}_i) - \phi_i(\mathbf{a}'_i)\|^2 \leq \sum_{i=1}^m \rho^2 \|\mathbf{a}_i - \mathbf{a}'_i\|^2 = \rho^2 \|\mathbf{A} - \mathbf{A}'\|_F^2 \leq \rho^2 r^2,$$

$\Phi \circ \mathcal{A}'$ is a ρr -cover of $\Phi \circ \mathcal{A}$. That implies $\mathcal{N}(\rho r, \Phi \circ \mathcal{A}, \|\cdot\|_F) \leq |\Phi \circ \mathcal{A}'| = |\mathcal{A}'| = \mathcal{N}(r, \mathcal{A}, \|\cdot\|_F)$. \square

Lemma A.9. Given two sets $\mathcal{A}_1, \mathcal{A}_2 \subseteq \mathbb{R}^{m \times d}$. Define $\mathcal{A}_3 \triangleq \{[\mathbf{A}_1, \mathbf{A}_2] \in \mathbb{R}^{m \times 2d} \mid \mathbf{A}_1 \in \mathcal{A}_1, \mathbf{A}_2 \in \mathcal{A}_2\}$, then

$$\mathcal{N}(\epsilon/\sqrt{2}, \mathcal{A}_3, \|\cdot\|_F) \leq \mathcal{N}(\epsilon, \mathcal{A}_1, \|\cdot\|_F) \mathcal{N}(\epsilon, \mathcal{A}_2, \|\cdot\|_F).$$

Proof. Let $\bar{\mathcal{A}}_1, \bar{\mathcal{A}}_2$ be minimum ϵ -cover sets of $\mathcal{A}_1, \mathcal{A}_2$, resp. By definition, for all $\mathbf{A}_1 \in \mathcal{A}_1, \mathbf{A}_2 \in \mathcal{A}_2$,

$$\begin{aligned} \exists \bar{\mathbf{A}}_1 \in \bar{\mathcal{A}}_1, \quad \|\mathbf{A}_1 - \bar{\mathbf{A}}_1\|_F^2 &\leq \epsilon^2, \text{ and} \\ \exists \bar{\mathbf{A}}_2 \in \bar{\mathcal{A}}_2, \quad \|\mathbf{A}_2 - \bar{\mathbf{A}}_2\|_F^2 &\leq \epsilon^2 \\ \implies \quad \|\mathbf{A}_1, \mathbf{A}_2 - [\bar{\mathbf{A}}_1, \bar{\mathbf{A}}_2]\|_F^2 &\leq 2\epsilon^2. \end{aligned}$$

Therefore, set $\bar{\mathcal{A}}_3 \triangleq \{[\bar{\mathbf{A}}_1, \bar{\mathbf{A}}_2] \in \mathbb{R}^{m \times 2d} \mid \bar{\mathbf{A}}_1 \in \bar{\mathcal{A}}_1, \bar{\mathbf{A}}_2 \in \bar{\mathcal{A}}_2\}$ is a $(\epsilon/\sqrt{2})$ -cover set of \mathcal{A}_3 , which leads to our conclusion,

$$\mathcal{N}(\epsilon/\sqrt{2}, \mathcal{A}_3, \|\cdot\|_F) \leq |\bar{\mathcal{A}}_3| = |\bar{\mathcal{A}}_1| |\bar{\mathcal{A}}_2| = \mathcal{N}(\epsilon, \mathcal{A}_1, \|\cdot\|_F) \mathcal{N}(\epsilon, \mathcal{A}_2, \|\cdot\|_F). \quad \square$$

Lemma A.10. Covering number of neural networks (Bartlett et al., 2017, Theorem 3.3) Let fixed nonlinearities $(\sigma_1, \dots, \sigma_L)$ and reference matrices (M_1, \dots, M_L) be given, where σ_i is ρ_i -Lipschitz and $\sigma_i(0) = 0$. Let spectral norm bounds (s_1, \dots, s_L) , and matrix (2,1) norm bounds b_1, \dots, b_L be given. Let data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ be given, where the n rows correspond to data points. Let $\mathcal{H}_{\mathbf{X}}$ denote the family of matrices obtained by evaluating \mathbf{X} with all choices of network $F_{\mathcal{A}}$:

$$\mathcal{H}_{\mathbf{X}} \triangleq \left\{ F_{\mathcal{A}}(\mathbf{X}^\top) \mid \mathcal{A} = (\mathbf{A}_1, \dots, \mathbf{A}_L), \|\mathbf{A}_i\|_\sigma \leq s_i, \|\mathbf{A}_i^\top - \mathbf{I}\|_{2,1} \leq b_i \right\},$$

where each matrix has dimension at most D along each axis. Then for any $\epsilon > 0$,

$$\log \mathcal{N}(\epsilon, \mathcal{H}_{\mathbf{X}}, \|\cdot\|_F) \leq \frac{\|\mathbf{X}\|_F^2 R_{NET}}{\epsilon^2},$$

where $R_{NET} = \log(2D^2) \left(\prod_{j=1}^L s_j^2 \rho_j^2 \right) \left(\sum_{L}^i (b_i/s_i)^{2/3} \right)^3$ is a constant depending only on the neural network's properties.

B. Proof of Theorem 3.4

Let $\mathbf{M}^{\natural} = [\mathbf{f}^{\natural}(\mathbf{x}_1), \dots, \mathbf{f}^{\natural}(\mathbf{x}_N)]$. As \mathbf{M}^{\natural} satisfies separability condition, we assume that $\mathbf{M}^{\natural} = [\mathbf{I}, \mathbf{V}^{\natural}]$, which is w.o.l.g. Let $\mathbf{M}^{\star} = [\mathbf{f}^{\star}(\mathbf{x}_1), \dots, \mathbf{f}^{\star}(\mathbf{x}_N)]$ denote the predicted cluster memberships, and also partition \mathbf{M}^{\star} as $\mathbf{M}^{\star} = [\mathbf{U}^{\star}, \mathbf{V}^{\star}]$, $\mathbf{U}^{\star} \in \mathbb{R}^{K \times K}$. Then,

$$\begin{aligned} \|\mathbf{P}^{\star} - \mathbf{P}^{\natural}\|_{\text{F}}^2 &= \|(\mathbf{M}^{\star})^{\top} \mathbf{M}^{\star} - (\mathbf{M}^{\natural})^{\top} \mathbf{M}^{\natural}\|_{\text{F}}^2 \\ &= \left\| \begin{bmatrix} (\mathbf{U}^{\star})^{\top} \\ (\mathbf{V}^{\star})^{\top} \end{bmatrix} \begin{bmatrix} \mathbf{U}^{\star} & \mathbf{V}^{\star} \end{bmatrix} - \begin{bmatrix} \mathbf{I}_K \\ (\mathbf{V}^{\natural})^{\top} \end{bmatrix} \begin{bmatrix} \mathbf{I}_K & \mathbf{V}^{\natural} \end{bmatrix} \right\|_{\text{F}}^2 \\ &= \|(\mathbf{U}^{\star})^{\top} \mathbf{U}^{\star} - \mathbf{I}_K\|_{\text{F}}^2 + 2 \|(\mathbf{U}^{\star})^{\top} \mathbf{V}^{\star} - \mathbf{V}^{\natural}\|_{\text{F}}^2 + \|(\mathbf{V}^{\star})^{\top} \mathbf{V}^{\star} - (\mathbf{V}^{\natural})^{\top} \mathbf{V}^{\natural}\|_{\text{F}}^2. \end{aligned} \quad (31)$$

Analysis of the first term in (31). Let $Z = \frac{1}{N^2} \|\mathbf{P}^{\star} - \mathbf{P}^{\natural}\|_{\text{F}}^2$, $0 \leq Z \leq 1$ be a random variable with a PDF $f(z)$. By Lemma 3.3,

$$\Pr(Z \leq \epsilon(M, \delta)^2) \geq 1 - \delta.$$

We have

$$\begin{aligned} \mathbb{E}[Z] &= \int_0^1 z f(z) dz = \int_0^{\epsilon(M, \delta)^2} z f(z) dz + \int_{\epsilon(M, \delta)^2}^1 z f(z) dz \leq \int_0^{\epsilon(M, \delta)^2} \epsilon(M, \delta)^2 f(z) dz + \int_{\epsilon(M, \delta)^2}^1 f(z) dz \\ &\leq \epsilon(M, \delta)^2 + \delta. \end{aligned}$$

Recall $S = \{(i_1, j_1, y_1), \dots, (i_M, j_M, y_M)\}$. According to Lemma B.3,

$$\mathbb{E}_{\mathbf{X}, S} \left[(P_{ij}^{\star} - P_{ij}^{\natural})^2 \right] = \text{const}, \quad \forall (i, j) \in [N]^2.$$

Therefore, for arbitrary $i, j \in [N]^2$,

$$\mathbb{E}_{\mathbf{X}, S} \left[(P_{ij}^{\star} - P_{ij}^{\natural})^2 \right] = \mathbb{E}_{\mathbf{X}, S} \left[\frac{1}{N^2} \|\mathbf{P}^{\star} - \mathbf{P}^{\natural}\|_{\text{F}}^2 \right] \leq \epsilon(M, \delta)^2 + \delta.$$

Meanwhile, by Markov inequality, for $\tau > 0$,

$$\Pr \left((P_{ij}^{\star} - P_{ij}^{\natural})^2 \leq \tau \mathbb{E}_{\mathbf{X}, S} \left[(P_{ij}^{\star} - P_{ij}^{\natural})^2 \right] \right) \geq 1 - \frac{1}{\tau}.$$

Setting $\tau = M^{1/4}$, $\delta = e^{-\sqrt{M}}$ gives

$$\left(P_{ij}^{\star} - P_{ij}^{\natural} \right)^2 \leq M^{1/4} \mathbb{E}_{\mathbf{X}, S} \left[(P_{ij}^{\star} - P_{ij}^{\natural})^2 \right] \leq M^{1/4} \epsilon(M, e^{-\sqrt{M}})^2 + M^{1/4} e^{-\sqrt{M}}$$

holds with probability at least $1 - 1/M^{1/4}$.

By union bound, the following holds with probability at least $1 - K^2/M^{1/4}$,

$$\|(\mathbf{U}^{\star})^{\top} \mathbf{U}^{\star} - \mathbf{I}_K\|_{\text{F}}^2 = \sum_{(i, j) \in [K]^2} (P_{ij}^{\star} - P_{ij}^{\natural})^2 \leq K^2 \underbrace{(M^{1/4} \epsilon(M, e^{-\sqrt{M}})^2 + M^{1/4} e^{-\sqrt{M}})}_{\epsilon'(M)^2} = K^2 \epsilon'(M)^2, \quad (32)$$

where $\epsilon'(M)$ is defined in (9). Inequality (32) implies

$$\begin{cases} 1 - K \epsilon'(M) \leq \|\mathbf{u}_k^{\star}\|^2 \leq 1, & \forall k \in [K] \end{cases} \quad (33)$$

$$\begin{cases} 0 \leq \langle \mathbf{u}_k^{\star}, \mathbf{u}_\ell^{\star} \rangle \leq K \epsilon'(M), & \forall k \neq \ell, k, \ell \in [K]. \end{cases} \quad (34)$$

For any permutation matrix $\mathbf{\Pi} \in \mathbb{R}^{K \times K}$,

$$\begin{aligned} \|\mathbf{\Pi U}^* - \mathbf{I}\|_{\text{F}}^2 &= \sum_{k=1}^K \|\mathbf{u}_k^* - \boldsymbol{\pi}_k\|^2 = K + \sum_{k=1}^K \|\mathbf{u}_k^*\|^2 - 2 \sum_{k=1}^K \langle \mathbf{u}_k^*, \boldsymbol{\pi}_k \rangle \\ &\leq 2K - 2 \sum_{k=1}^K \langle \mathbf{u}_k^*, \boldsymbol{\pi}_k \rangle \quad (\text{Since } \mathbf{u}_k^* \geq 0, \mathbf{1}^\top \mathbf{u}_k^* = 1 \text{ and thus } \|\mathbf{u}_k^*\| \leq 1). \end{aligned}$$

Minimizing over $\mathbf{\Pi}$ on both sides,

$$\min_{\mathbf{\Pi}} \|\mathbf{\Pi U}^* - \mathbf{I}\|_{\text{F}}^2 \leq 2K - 2 \max_{\mathbf{\Pi}} \sum_{k=1}^K \langle \mathbf{u}_k^*, \boldsymbol{\pi}_k \rangle. \quad (35)$$

Denote $k_i^* := \arg \max_k U_{ki}^*$. Invoking Lemma B.2 where $K\epsilon'(M) < 0.381$ (by (9)) and \mathbf{u}_k^* 's satisfy (33) and (34), it holds that k_1^*, \dots, k_K^* are all different integers, and we assume w.o.l.g. that $k_1^* = 1, \dots, k_K^* = K$, i.e., $U_{ii}^* = \arg \max_k U_{ki}^*$. With that in mind, we get a simple lower bound,

$$\max_{\mathbf{\Pi}} \sum_{k=1}^K (\mathbf{u}_k^*)^\top \boldsymbol{\pi}_k \geq \sum_{k=1}^K U_{kk}^*. \quad (36)$$

Since $\|\mathbf{u}_1^*\|^2 \geq 1 - K\epsilon'(M)$ by (33),

$$(U_{11}^*)^2 \geq 1 - K\epsilon'(M) - \sum_{k=2}^K (U_{k1}^*)^2 \geq 1 - K\epsilon'(M) - \left(\sum_{k=2}^K U_{k1}^* \right)^2 = 1 - K\epsilon'(M) - (1 - U_{11}^*)^2,$$

where the second inequality holds since $U_{k1}^* \geq 0, \forall k \in [K]$, and the last equality holds since \mathbf{u}_1 is a probability simplex vector. As $U_{11}^* = \max_k U_{k1}$, $U_{11}^* \geq 1/K$. Solving the above quadratic inequality with respect to $1/K \leq U_{11}^* \leq 1$ and $K\epsilon'(M) \leq 0.381$ (by (9)) leads to

$$U_{11}^* \geq \frac{1 + \sqrt{1 - 2K\epsilon'(M)}}{2} \geq 1 - K\epsilon'(M).$$

Applying the same argument for $U_{22}^*, \dots, U_{KK}^*$,

$$U_{kk}^* \geq 1 - K\epsilon'(M), \quad k \in [K].$$

Therefore,

$$\min_{\mathbf{\Pi}} \|\mathbf{\Pi U}^* - \mathbf{I}\|_{\text{F}}^2 \leq 2K - 2(K - K^2\epsilon'(M)) = 2K^2\epsilon'(M). \quad (37)$$

This implies that \mathbf{U}^* is close to an identity matrix up to some permutation.

Analysis of the second term in (31). For a fixed $0 < \delta < 1$, suppose that the following event happens

$$\frac{1}{N^2} \|\mathbf{P}^* - \mathbf{P}^\natural\|_{\text{F}}^2 \leq \epsilon(M, \delta)^2. \quad (38)$$

Recall that $\mathbf{M}^* = [\mathbf{U}^*, \mathbf{V}^*], \mathbf{M}^\natural = [\mathbf{I}_K, \mathbf{V}^\natural]$, then

$$\|(\mathbf{U}^*)^\top \mathbf{V}^* - \mathbf{V}^\natural\|_{\text{F}} \leq \frac{N\epsilon(M, \delta)}{\sqrt{2}}.$$

Denote $\mathbf{\Pi}^*$ as the optimal permutation matrix in LHS of (37). We have

$$\begin{aligned} \|(\mathbf{U}^*)^\top \mathbf{V}^* - \mathbf{V}^\natural - ((\mathbf{U}^*)^\top (\mathbf{\Pi}^*)^\top - \mathbf{I}) \mathbf{V}^\natural\|_{\text{F}} &\leq \|(\mathbf{U}^*)^\top \mathbf{V}^* - \mathbf{V}^\natural\|_{\text{F}} + \|((\mathbf{U}^*)^\top (\mathbf{\Pi}^*)^\top - \mathbf{I}) \mathbf{V}^\natural\|_{\text{F}} \\ &\leq \frac{N\epsilon(M, \delta)}{\sqrt{2}} + \sigma_{\max}(\mathbf{V}^\natural) \|\mathbf{\Pi}^* \mathbf{U}^* - \mathbf{I}\|_{\text{F}} \\ &\leq \frac{N\epsilon(M, \delta)}{\sqrt{2}} + \sigma_{\max}(\mathbf{V}^\natural) K \sqrt{2\epsilon'(M)} \quad (\text{thanks to (37)}). \end{aligned} \quad (39)$$

On the other hand,

$$\begin{aligned} \left\| (\mathbf{U}^*)^\top \mathbf{V}^* - \mathbf{V}^\natural - ((\mathbf{U}^*)^\top (\mathbf{\Pi}^*)^\top - \mathbf{I}) \mathbf{V}^\natural \right\|_F &= \left\| (\mathbf{U}^*)^\top (\mathbf{V}^* - (\mathbf{\Pi}^*)^\top \mathbf{V}^\natural) \right\|_F \\ &\geq \sigma_{\min}(\mathbf{U}^*) \left\| \mathbf{V}^* - (\mathbf{\Pi}^*)^\top \mathbf{V}^\natural \right\|_F \\ &= \sigma_{\min}(\mathbf{U}^*) \left\| \mathbf{\Pi}^* \mathbf{V}^* - \mathbf{V}^\natural \right\|_F. \end{aligned} \quad (40)$$

Using Lemma B.1, we can upper bound $\sigma_{\min}(\mathbf{U}^*)$ as

$$\left| \sigma_{\min}(\mathbf{I} + \mathbf{\Pi}^* \mathbf{U}^* - \mathbf{I}) - \sigma_{\min}(\mathbf{I}) \right| \leq \left\| \mathbf{\Pi}^* \mathbf{U}^* - \mathbf{I} \right\|_F \leq K \sqrt{2\epsilon'(M)} \quad (\text{thanks to (37)})$$

or,

$$\sigma_{\min}(\mathbf{U}^*) = \sigma_{\min}(\mathbf{\Pi}^* \mathbf{U}^*) \geq 1 - K \sqrt{2\epsilon'(M)}. \quad (41)$$

Combine (39), (40), (41),

$$\begin{aligned} \left\| \mathbf{\Pi}^* \mathbf{V}^* - \mathbf{V}^\natural \right\|_F^2 &\leq \left(\frac{N\epsilon(M, \delta) + 2\sigma_{\max}(\mathbf{V}^\natural)K\sqrt{\epsilon'(M)}}{\sqrt{2} - 2K\sqrt{\epsilon'(M)}} \right)^2 \\ &\leq 2 \left(N\epsilon(M, \delta) + 2\sigma_{\max}(\mathbf{V}^\natural)K\sqrt{\epsilon'(M)} \right)^2 \quad (\text{Since } 8K^2\epsilon'(M) \leq 1 \text{ by (9)}) \\ &\leq 4N^2\epsilon(M, \delta)^2 + 16\sigma_{\max}^2(\mathbf{V}^\natural)K^2\epsilon'(M) \quad (\text{Cauchy-Schwarz inequality}). \end{aligned} \quad (42)$$

Lastly, combine (37) and (42), we finish our proof,

$$\begin{aligned} \left\| \mathbf{\Pi}^* \mathbf{M}^* - \mathbf{M}^\natural \right\|_F^2 &= \left\| \mathbf{\Pi}^* \mathbf{U} - \mathbf{I} \right\|_F^2 + \left\| \mathbf{\Pi}^* \mathbf{V} - \mathbf{V}^\natural \right\|_F^2 \\ &\leq 2K^2\epsilon'(M) + 4N^2\epsilon(M, \delta)^2 + 16\sigma_{\max}^2(\mathbf{V}^\natural)K^2\epsilon'(M) \\ &= 4N^2\epsilon(M, \delta)^2 + 2K^2(1 + 8\sigma_{\max}^2(\mathbf{V}^\natural))\epsilon'(M). \end{aligned}$$

Note that the whole computation is derived based on the realizations of 2 independent events in (32), (38) with probability of happening at least $1 - K^2/M^{1/4}$ and $1 - \delta$, resp,

$$\frac{1}{N} \left\| \mathbf{\Pi}^* \mathbf{M}^* - \mathbf{M}^\natural \right\|_F^2 \leq 4N\epsilon(M, \delta)^2 + \frac{2}{N}K^2(1 + 8\sigma_{\max}^2(\mathbf{V}^\natural))\epsilon'(M).$$

holds with probability at least $1 - (\delta + K^2/M^{1/4})$.

B.1. Supporting Lemmas for proof of Theorem 3.4

Lemma B.1 ((Weyl, 1912)). *Let $\mathbf{X}, \mathbf{\Delta} \in \mathbb{R}^{m \times n}$,*

$$|\sigma_i(\mathbf{X} + \mathbf{\Delta}) - \sigma_i(\mathbf{\Delta})| \leq \left\| \mathbf{\Delta} \right\|_2 \quad (\leq \left\| \mathbf{\Delta} \right\|_F), \quad 1 \leq i \leq \min(m, n).$$

Lemma B.2. *Define $A_\epsilon = \left\{ \mathbf{x} \in \mathbb{R}^K \mid \mathbf{1}^\top \mathbf{x} = 1, \mathbf{x} \geq 0, \left\| \mathbf{x} \right\|_2^2 \geq 1 - \epsilon \right\}$. If $\epsilon \leq 0.381$, then*

$$\max_{\mathbf{x}, \mathbf{y} \in A_\epsilon} \mathbf{x}^\top \mathbf{y} \leq \epsilon \quad \Rightarrow \quad \arg \max_k x_k \neq \arg \max_k y_k.$$

Proof. We use contradiction to prove the claim. Suppose that $\arg \max_k x_k = \arg \min_k y_k = \tilde{k}$. Without loss of generality, let us assume $\tilde{k} = 1$. Using these assumptions, we will show that $\min_{\mathbf{x}, \mathbf{y} \in A_\epsilon} \mathbf{x}^\top \mathbf{y} > \epsilon$. Indeed,

$$1 - \epsilon \leq \left\| \mathbf{x} \right\|_2^2 \leq x_1 \left(\sum_{k=1}^K x_k \right) = x_1.$$

Similarly, we obtain $y_1 \geq 1 - \epsilon$, which leads to

$$\mathbf{x}^\top \mathbf{y} \geq (1 - \epsilon)^2 > \epsilon,$$

where the second inequality holds because $\epsilon \leq 0.381$. □

Lemma B.3. Assume that $\omega = (\omega^1, \omega^2)$ is a uniform RV over $[N] \times [N]$ where $\omega^1 \in [N], \omega^2 \in [N]$, and $\mathbf{x}_1, \dots, \mathbf{x}_N$ are i.i.d. Define the following mappings:

$$\begin{aligned} \text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_N, \omega) &\triangleq (\mathbf{x}_{\omega^1}, \mathbf{x}_{\omega^2}), \\ \text{gen}(\omega_1, \dots, \omega_M, y_1, \dots, y_M, \mathbf{x}_1, \dots, \mathbf{x}_N) &\triangleq \{\psi_1, \dots, \psi_M\}, \quad \psi_i \triangleq (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_N, \omega_i), y_i), \end{aligned}$$

where

$$\begin{aligned} \omega_i &= (k, \ell) \in [N]^2, \quad 1 \leq i \leq M, \\ y_i &\in \{0, 1\}, \quad 1 \leq i \leq M, \\ \mathbf{x}_i &\in \mathcal{X}, \quad 1 \leq i \leq N. \end{aligned}$$

For a real-valued function $h(\mathbf{u}, \mathbf{v}, \text{gen}(\omega_1, \dots, \omega_M, y_1, \dots, y_M, \mathbf{x}_1, \dots, \mathbf{x}_N))$ with $\mathbf{u}, \mathbf{v} \in \mathcal{X}$, we have: $\forall y_1, \dots, y_M$, and $\forall i, j, i', j' \in [N]$,

$$\begin{aligned} \mathbb{E}_{\substack{\omega_i, 1 \leq i \leq M, \\ \mathbf{x}_i, 1 \leq i \leq N}} [h(\mathbf{x}_i, \mathbf{x}_j, \text{gen}(\omega_1, \dots, \omega_M, y_1, \dots, y_M, \mathbf{x}_1, \dots, \mathbf{x}_N))] = \\ \mathbb{E}_{\substack{\omega_i, 1 \leq i \leq M, \\ \mathbf{x}_i, 1 \leq i \leq N}} [h(\mathbf{x}_{i'}, \mathbf{x}_{j'}, \text{gen}(\omega_1, \dots, \omega_M, y_1, \dots, y_M, \mathbf{x}_1, \dots, \mathbf{x}_N))]. \quad (43) \end{aligned}$$

Proof. We have the following property on the mapping gen : $\forall i \neq j, i, j \in [N]$,

$$\begin{aligned} \text{gen}(\omega_1, \dots, \omega_M, y_1, \dots, y_M, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N) = \\ \text{gen}(s_{ij}(\omega_1), \dots, s_{ij}(\omega_M), y_1, \dots, y_M, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N), \quad (44) \end{aligned}$$

where $s_{ij} : [N]^2 \rightarrow [N]^2$,

$$s_{ij}(\omega) = s_{ij}((k, \ell)) \triangleq \begin{cases} (k, \ell) & \text{if } k \neq i \text{ and } \ell \neq j \\ (j, \ell) & \text{if } k = i \text{ and } \ell \neq j \\ (k, i) & \text{if } k \neq i \text{ and } \ell = j \\ (j, i) & \text{if } k = i \text{ and } \ell = j \end{cases}$$

To see this, let

$$\begin{aligned} \mathcal{G}_1 &= \text{gen}(\omega_1, \dots, \omega_M, y_1, \dots, y_M, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N), \\ \mathcal{G}_2 &= \text{gen}(s_{ij}(\omega_1), \dots, s_{ij}(\omega_M), y_1, \dots, y_M, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N). \end{aligned}$$

Pick $\psi = (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N, \omega), y) \in \mathcal{G}_1$. Assume $\omega = (k, \ell)$. Consider the following cases,

- If $k \neq i$ and $\ell \neq j$, then $s_{ij}(\omega) = \omega$. So

$$\psi = (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, s_{ij}(\omega)), y) \in \mathcal{G}_2.$$

- If $k = i$ and $\ell \neq j$, then $s_{ij}(\omega) = s_{ij}((i, \ell)) = (j, \ell)$. So

$$\begin{aligned} \psi &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N, (i, \ell)), y) \\ &= (\mathbf{x}_i, \mathbf{x}_\ell, y) \\ &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, (j, \ell)), y) \\ &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, s_{ij}(\omega)), y) \in \mathcal{G}_2 \quad (\text{by definition of the mapping gen}). \end{aligned}$$

- If $k \neq i$ and $\ell = j$, then $s_{ij}(\omega) = s_{ij}((k, j)) = (k, i)$. So

$$\begin{aligned} \psi &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N, (k, j)), y) \\ &= (\mathbf{x}_k, \mathbf{x}_j, y) \\ &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, (k, i)), y) \\ &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, s_{ij}(\omega)), y) \in \mathcal{G}_2 \quad (\text{by definition of the mapping gen}). \end{aligned}$$

- If $k = i$ and $\ell = j$, then $s_{ij}(\omega) = s_{ij}((i, j)) = (j, i)$. So

$$\begin{aligned}\psi &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_j, \dots, \mathbf{x}_N, (i, j)), y) \\ &= (\mathbf{x}_i, \mathbf{x}_j, y) \\ &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, (j, i)), y) \\ &= (\text{pick}(\mathbf{x}_1, \dots, \mathbf{x}_j, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N, s_{ij}(\omega)), y) \in \mathcal{G}_2 \quad (\text{by definition of the mapping } \text{gen}).\end{aligned}$$

Since $\psi \in \mathcal{G}_1$ is arbitrary, $\mathcal{G}_1 \subseteq \mathcal{G}_2$. Moreover, $|\mathcal{G}_1| = |\mathcal{G}_2|$, so $\mathcal{G}_1 = \mathcal{G}_2$, and hence (44) holds.

Notice that the mapping $s_{ij}(\omega)$ is surjective and onto $[N]^2$, we proceed with the following

$$\begin{aligned}& \mathbb{E}_{\substack{\mathbf{x}_i, \mathbf{x}_{i'}, \\ \mathbf{x}_{k, k=[N] \setminus \{i, i'\}}, \\ \omega_1, \dots, \omega_M}} [h(\mathbf{x}_i, \mathbf{x}_j, \text{gen}(\omega_1, \dots, \omega_M, y_1, \dots, y_M, \dots, \mathbf{x}_i, \dots, \mathbf{x}_{i'}, \dots))] \\ &= \mathbb{E}_{\substack{\mathbf{x}_i, \mathbf{x}_{i'}, \\ \mathbf{x}_{k, k=[N] \setminus \{i, i'\}}, \\ \omega_1, \dots, \omega_M}} [h(\mathbf{x}_i, \mathbf{x}_j, \text{gen}(s_{ii'}(\omega_1), \dots, s_{ii'}(\omega_M), y_1, \dots, y_M, \dots, \mathbf{x}_{i'}, \dots, \mathbf{x}_i, \dots))] \\ &= \mathbb{E}_{\substack{\mathbf{u}, \mathbf{v}, \\ \mathbf{x}_{k, k=[N] \setminus \{i, i'\}}, \\ \omega_1, \dots, \omega_M}} [h(\mathbf{u}, \mathbf{x}_j, \text{gen}(s_{ii'}(\omega_1), \dots, s_{ii'}(\omega_M), y_1, \dots, y_M, \dots, \mathbf{v}, \dots, \mathbf{u}, \dots))] \\ &= \mathbb{E}_{\substack{\mathbf{x}_{i'}, \mathbf{x}_i, \\ \mathbf{x}_{k, k=[N] \setminus \{i, i'\}}, \\ \omega_1, \dots, \omega_M}} [h(\mathbf{x}_{i'}, \mathbf{x}_j, \text{gen}(s_{ii'}(\omega_1), \dots, s_{ii'}(\omega_M), y_1, \dots, y_M, \dots, \mathbf{x}_i, \dots, \mathbf{x}_{i'}, \dots))].\end{aligned}$$

The first equality holds by property (44), the second and the third equality hold as we just rename some random variables under the expectation. Proceed with the same argument to replace j with j' , we obtain

$$\begin{aligned}\mathbb{E}_{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N, \\ \omega_1, \dots, \omega_M}} [h(\mathbf{x}_{i'}, \mathbf{x}_j, \text{gen}(s_{ii'}(\omega_1), \dots, s_{ii'}(\omega_M), y_1, \dots, y_M, \mathbf{x}_1, \dots, \mathbf{x}_N))] \\ = \mathbb{E}_{\substack{\mathbf{x}_1, \dots, \mathbf{x}_N, \\ \omega_1, \dots, \omega_M}} [h(\mathbf{x}_{i'}, \mathbf{x}_{j'}, \text{gen}(s_{jj'}(s_{ii'}(\omega_1)), \dots, s_{jj'}(s_{ii'}(\omega_M)), y_1, \dots, y_M, \mathbf{x}_1, \dots, \mathbf{x}_N))]. \quad (45)\end{aligned}$$

Since $s_{ii'}(\omega)$ is surjective and onto itself, $s_{ii'}(\omega)$ has the same probability distribution as ω . So $s_{jj'}(s_{ii'}(\omega))$ also define a RV with the same PDF as ω , and hence (45) implies (43). □

C. Proof of Theorem 3.5

We start with the following key lemma:

Lemma C.1. *With probability at least $1 - \delta$ over $\mathbf{x}_1, \dots, \mathbf{x}_N$,*

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{P}_{\mathcal{X}}} \left[\|\mathbf{\Pi}^* \mathbf{f}^*(\mathbf{x}) - \mathbf{f}^{\natural}(\mathbf{x})\|^2 \right] \leq \frac{1}{N} \|\mathbf{\Pi}^* \mathbf{M}^* - \mathbf{M}^{\natural}\|_{\text{F}}^2 + \frac{8}{N} + \frac{12\sqrt{2R_{NET}} \|\mathbf{X}\|_{\text{F}} \log N}{N \log 2} + 8\sqrt{\frac{2 \log(4\delta)}{N}},$$

where $\mathbf{\Pi}^*$ is the optimal permutation matrix in LHS of (37) in the proof of Theorem 3.4.

Proof. Let us define the following:

$$\begin{aligned}L_{\mathbf{X}}(\mathbf{f}) &\triangleq \frac{1}{N} \sum_{n=1}^N \text{SE}(\mathbf{f}, (\mathbf{x}_n, \mathbf{m}_n^{\natural})), \\ \text{SE}(\mathbf{f}, (\mathbf{x}, \mathbf{m})) &\triangleq \|\mathbf{\Pi}^* \mathbf{f}(\mathbf{x}) - \mathbf{m}\|^2.\end{aligned}$$

Note that $L_{\mathbf{X}}(\mathbf{f})$ can be viewed as the empirical loss in statistical learning where a dataset of N i.i.d samples $(\mathbf{x}_1, \mathbf{m}_1^{\natural}), \dots, (\mathbf{x}_N, \mathbf{m}_N^{\natural})$ are given. Denote $\mathcal{P}_{\mathcal{X}, \mathbf{f}^{\natural}}$ as the PDF of the joint distribution of $(\mathbf{x}_n, \mathbf{m}_n^{\natural})$, and define

$$L_{\mathcal{P}_{\mathcal{X}, \mathbf{f}^{\natural}}}(\mathbf{f}) \triangleq \mathbb{E}_{(\mathbf{x}, \mathbf{m}^{\natural}) \sim \mathcal{P}_{\mathcal{X}, \mathbf{f}^{\natural}}} [\text{SE}(\mathbf{f}, (\mathbf{x}, \mathbf{m}^{\natural}))],$$

which is the expected loss. Then it is readily to apply a standard generalization bound from (Shalev-Shwartz & Ben-David, 2014, Theorem 26.5) using the fact that $\text{SE}(\mathbf{f}, (\mathbf{x}, \mathbf{m})) \leq 2$. That is, we have that

$$L_{\mathcal{P}_{\mathbf{x}, \mathbf{f}^\dagger}}(\mathbf{f}^*) \leq L_{\mathbf{X}}(\mathbf{f}^*) + 2\mathcal{R}(\text{SE} \circ \mathcal{F} \circ \mathbf{X}) + 8\sqrt{\frac{2\log(4/\delta)}{N}},$$

holds with probability of at least $1 - \delta$. The first term is

$$L_{\mathbf{X}}(\mathbf{f}^*) = \frac{1}{N} \|\mathbf{\Pi}^* \mathbf{M}^* - \mathbf{M}^\dagger\|_{\text{F}}^2.$$

The second term is bounded by following the proof of Lemma A.3,

$$\begin{aligned} \mathcal{N}(\epsilon, \text{SE} \circ \mathcal{F} \circ \mathbf{X}, \|\cdot\|) &\leq \mathcal{N}\left(\frac{\epsilon}{2\sqrt{2}}, \mathcal{F} \circ \mathbf{X}, \|\cdot\|_{\text{F}}\right) \\ &\leq \exp\left(\frac{8R_{\text{NET}} \|\mathbf{X}\|_{\text{F}}^2}{\epsilon^2}\right) \quad (\text{By Lemma A.10}). \end{aligned} \quad (46)$$

Inequality (46) holds by applying Lemma A.7 and the fact that function $\phi_{\mathbf{m}}(\mathbf{x}) \triangleq \|\mathbf{x} - \mathbf{m}\|^2$ is $2\sqrt{2}$ -Lipschitz continuous on the domain $\Delta = \{\mathbf{x} \in \mathbb{R}^K \mid \mathbf{x} \geq 0, \mathbf{1}^\top \mathbf{x} = 1\}$. This is because of the following fact:

$$\|\nabla \phi_{\mathbf{m}}(\mathbf{x})\| = 2\|\mathbf{x} - \mathbf{m}\| \leq 2\sqrt{2}.$$

We proceed by deriving a bound on $\mathcal{R}(\text{SE} \circ \mathcal{F} \circ \mathbf{X})$ with the help of Dudley's entropy integral technique. In particular, applying Lemma A.6 on the set $\text{SE} \circ \mathcal{F} \circ \mathbf{X}$, the following holds: For any integer $T > 0$,

$$\begin{aligned} \mathcal{R}(\text{SE} \circ \mathcal{F} \circ \mathbf{X}) &\leq \frac{c2^{-T}}{\sqrt{N}} + \frac{6c}{N} \sum_{t=1}^T 2^{-t} \sqrt{\log(\mathcal{N}(c2^{-t}, \text{SE} \circ \mathcal{F} \circ \mathbf{X}))} \\ &\leq \frac{c2^{-T}}{\sqrt{N}} + \frac{6c}{N} \sum_{t=1}^T 2^{-t} \sqrt{\frac{8R_{\text{NET}} \|\mathbf{X}\|_{\text{F}}^2}{c^2 2^{-2t}}} \\ &= \frac{c2^{-T}}{\sqrt{N}} + \frac{12T\sqrt{2R_{\text{NET}}} \|\mathbf{X}\|_{\text{F}}}{N}. \end{aligned}$$

Substituting $c = 2$ as the upper bound of the loss SE and setting $T = \lfloor 0.5 \log_2(N) \rfloor$ gives

$$\mathcal{R}(\text{SE} \circ \mathcal{F} \circ \mathbf{X}) \leq \frac{4}{N} + \frac{6\sqrt{2R_{\text{NET}}} \|\mathbf{X}\|_{\text{F}} \log N}{N \log 2}.$$

Therefore,

$$\mathbb{E}_{(\mathbf{x}, \mathbf{m}^\dagger) \sim \mathcal{P}_{\mathbf{x}, \mathbf{f}^\dagger}} \left[\|\mathbf{\Pi}^* \mathbf{f}^*(\mathbf{x}) - \mathbf{m}^\dagger\|^2 \right] \leq \frac{1}{N} \|\mathbf{\Pi}^* \mathbf{M}^* - \mathbf{M}^\dagger\|_{\text{F}}^2 + \frac{8}{N} + \frac{12\sqrt{2R_{\text{NET}}} \|\mathbf{X}\|_{\text{F}} \log N}{N \log 2} + 8\sqrt{\frac{2\log(4/\delta)}{N}}$$

holds with probability of at least $1 - \delta$. \square

By Theorem 3.4,

$$L_{\mathbf{X}}(\mathbf{f}^*) = \frac{1}{N} \min_{\mathbf{H}} \sum_{n=1}^N \|\mathbf{\Pi} \mathbf{M}^* - \mathbf{M}^\dagger\|_{\text{F}}^2 \leq 4N\epsilon(M, \delta)^2 + \frac{2}{N} K^2(1 + 8\sigma_{\max}^2(\mathbf{V}^\dagger))\epsilon'(M)$$

holds with probability of at least $1 - \delta - K^2/M^{0.25}$. Invoking Lemma C.1 gives the conclusion,

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}, \mathbf{m}^\dagger) \sim \mathcal{P}_{\mathbf{x}, \mathbf{f}^\dagger}} \left[\|\mathbf{\Pi}^* \mathbf{f}^*(\mathbf{x}) - \mathbf{m}^\dagger\|^2 \right] &\leq 4N\epsilon(M, \delta)^2 + \frac{2}{N} K^2(1 + 8\sigma_{\max}^2(\mathbf{V}^\dagger))\epsilon'(M) \\ &\quad + \frac{8}{N} + \frac{12\sqrt{2R_{\text{NET}}} \|\mathbf{X}\|_{\text{F}} \log N}{N \log 2} + 8\sqrt{\frac{2\log(4/\delta)}{N}}, \end{aligned}$$

hold with probability of at least $1 - \delta - K^2/M^{0.25}$.

D. Proof of Theorem 3.7

As $\max(\log(1/\alpha), \log(M)\sqrt{R_{\text{net}}}/\alpha)/\sqrt{M} \rightarrow 0$, $\|\mathbf{P}^* - \mathbf{P}^\natural\|_F^2 \rightarrow 0$ accords to Lemma 3.3, which means

$$(\mathbf{M}^*)^\top \mathbf{M}^* = \mathbf{P}^\natural.$$

This gives a guarantee that there exists permutation matrix $\mathbf{\Pi}$ such that $\mathbf{\Pi}\mathbf{M}^* = \mathbf{M}^\natural$ by invoking (Huang et al., 2014, Theorem 4). Lastly, applying Lemma C.1 concludes the proof.

E. Proof of Theorem 4.1

Invoking Lemma 3.3 and considering $\max(\log(1/\alpha), \log(M)\sqrt{R_{\text{net}}}/\alpha)/\sqrt{M} \rightarrow 0$ gives

$$\mathbf{P}^* = (\mathbf{M}^*)^\top (\mathbf{A}^*)^\top \mathbf{A}^* \mathbf{M}^* = \mathbf{P}^\natural$$

at the limit. As assumed in Theorem 4.1, $\text{rank}((\mathbf{M}^\natural)^\top (\mathbf{A})^\top \mathbf{A}) = K$, thus $\text{rank}(\mathbf{M}^*) = K$. As proved in (Huang et al., 2016), \mathbf{M}^\natural satisfying SSC means $\text{rank}(\mathbf{M}^\natural) = K$, which also means $\text{rank}((\mathbf{M}^\natural)^\top (\mathbf{A})^\top \mathbf{A}) = K$ if $\text{rank}(\mathbf{A}) = K$. It implies that there exists an invertible matrix \mathbf{Q} such that $\mathbf{M}^* = \mathbf{Q}^{-1}\mathbf{M}^\natural$. To proceed, we use the following proposition which is a part of the proof of Theorem 1 in (Fu et al., 2015b).

Proposition E.1. *For $K \leq N$, let $\mathbf{M} \in \mathbb{R}^{K \times N}$ be a matrix of rank K and $\mathbf{M} \geq 0$, $\mathbf{1}^\top \mathbf{M} = \mathbf{1}$. For any invertible matrix \mathbf{Q} such that $\widehat{\mathbf{M}} \triangleq \mathbf{Q}^{-1}\mathbf{M}$ satisfies $\widehat{\mathbf{M}} \geq 0$, $\mathbf{1}^\top \widehat{\mathbf{M}} = \mathbf{1}^\top$, then $|\det \mathbf{Q}| \geq 1$. In addition, the equality holds if and only if \mathbf{Q} is a permutation matrix.*

By Proposition E.1, $|\det(\mathbf{Q})| \geq 1$, and hence

$$\det(\mathbf{M}^* (\mathbf{M}^*)^\top) = \det(\mathbf{Q}^{-1} \mathbf{M}^\natural (\mathbf{M}^\natural)^\top \mathbf{Q}^{-\top}) = \det(\mathbf{Q})^{-2} \det(\mathbf{M}^\natural (\mathbf{M}^\natural)^\top) \leq \det(\mathbf{M}^\natural (\mathbf{M}^\natural)^\top). \quad (47)$$

As we take the solution of (12) with the \mathbf{M}^* that has the largest volume, the equality in (47) is attained. Therefore, \mathbf{Q} can only be a permutation matrix. The rest follows by applying Lemma C.1.

F. Additional Experiments

F.1. Effect of Neural Network Complexity

In this section, we present some additional experiments. To be specific, we repeat experiments in Tables 1 and 2 using a more complex neural network for all the DCC methods. The purpose is to observe the effect of R_{NET} , i.e., the complexity of the neural network. Here, all the DCC methods use a three-hidden-layer fully connected neural networks, where each hidden layer has 512 ReLU activation functions. In the main text, the DCC methods used a two-hidden layer network, also with 512 activation functions in each layer.

Tables 5 and 6 show the new results. Overall, similar results as in the main text can be seen: VolMaxDCC outperforms all other baselines, especially when noise level is getting larger. In addition, VolMaxDCC enjoys a better clustering performance when using a more complex (and thus more expressive) neural network on STL10. When, the noise level is 15%, the new ACC of VolMaxDCC is (0.85,0.86), which is much higher than the previous case that used a two-hidden layer neural network (0.79,0.81). On CIFAR-10, some slight decrease of clustering accuracy occasionally appears. This may also suggest that using a deeper neural network may not always be the best choice. As implied in our theorems, a more complex neural network increases R_{NET} , which may increase the risk of overfitting—especially when N is not very large. In practice, the neural network structure may be selected following standard procedures, e.g., using validation sets, if available.

Table 5: Clustering performance of (seen data, unseen data) on STL10; $N_{\text{unseen}} = 2000$.

Noise level \ Methods		Methods			
		DC-GMM	C-IDEC	VanillaDCC	VolMaxDCC
0.0%	ACC	0.89, 0.87	0.88, 0.87	0.88, 0.86	0.92, 0.91
	NMI	0.83, 0.80	0.79, 0.78	0.79, 0.76	0.84, 0.82
	ARI	0.80, 0.78	0.77, 0.76	0.82, 0.80	0.84, 0.83
8.3%	ACC	0.78, 0.79	0.77, 0.79	0.76, 0.77	0.85, 0.86
	NMI	0.68, 0.70	0.67, 0.69	0.56, 0.59	0.73, 0.75
	ARI	0.59, 0.61	0.59, 0.61	0.66, 0.68	0.77, 0.78
10.3%	ACC	0.72, 0.74	0.69, 0.70	0.69, 0.71	0.84, 0.85
	NMI	0.63, 0.65	0.58, 0.60	0.47, 0.49	0.72, 0.73
	ARI	0.51, 0.53	0.49, 0.51	0.60, 0.61	0.77, 0.78
15.0%	ACC	0.59, 0.59	0.56, 0.57	0.54, 0.55	0.85, 0.86
	NMI	0.52, 0.53	0.49, 0.50	0.33, 0.34	0.72, 0.74
	ARI	0.36, 0.37	0.37, 0.38	0.48, 0.49	0.77, 0.78

 Table 6: Clustering performance of (seen data, unseen data) on CIFAR10; $N_{\text{unseen}} = 45000$.

Noise level \ Methods		Methods			
		DC-GMM	C-IDEC	VanillaDCC	VolMaxDCC
0.0%	ACC	0.90, 0.89	0.88, 0.87	0.89, 0.87	0.91, 0.90
	NMI	0.83, 0.80	0.81, 0.79	0.80, 0.77	0.83, 0.80
	ARI	0.82, 0.79	0.79, 0.76	0.82, 0.79	0.84, 0.82
4.9%	ACC	0.86, 0.86	0.85, 0.86	0.85, 0.86	0.86, 0.86
	NMI	0.77, 0.77	0.77, 0.77	0.73, 0.73	0.73, 0.74
	ARI	0.73, 0.73	0.73, 0.74	0.77, 0.77	0.77, 0.77
8.7%	ACC	0.78, 0.78	0.76, 0.77	0.76, 0.77	0.81, 0.81
	NMI	0.72, 0.72	0.70, 0.70	0.58, 0.58	0.65, 0.66
	ARI	0.59, 0.60	0.59, 0.59	0.70, 0.70	0.73, 0.73
10.9%	ACC	0.70, 0.71	0.74, 0.75	0.67, 0.68	0.81, 0.81
	NMI	0.64, 0.65	0.66, 0.66	0.47, 0.48	0.65, 0.66
	ARI	0.50, 0.51	0.56, 0.57	0.61, 0.61	0.74, 0.74

F.2. Effect of Pretrained Features

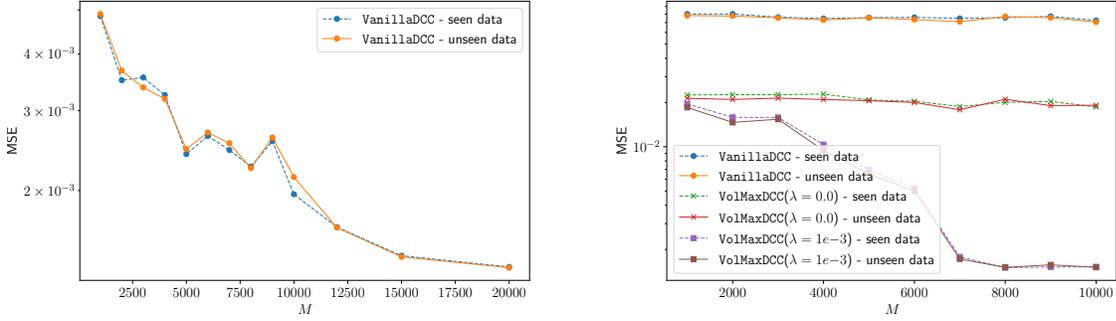
To see the effect of pre-training, we re-run the methods with feature vectors output by less well-trained feature extractors by (Li et al., 2021b). To be specific, here, the extractor is trained using 600 epochs, instead of 1000 epochs as in the main text. Nonetheless, similar results are observed in Table 7.

 Table 7: Clustering performance of (seen data, unseen data) on ImageNet10; $N_{\text{unseen}} = 2000$; embedding vectors are extracted by training (Li et al., 2021b) for 600 epochs.

Noise level \ Methods		Methods						
		K-means	COP-Kmeans	PCKmeans	DC-GMM	C-IDEC	VanillaDCC	VolMaxDCC
0.0%	ACC	0.83, —	0.79±0.06, —	0.74±0.07, —	0.97±0.00, 0.96±0.00	0.96±0.00, 0.96±0.00	0.96±0.00, 0.96±0.00	0.97±0.00, 0.96±0.00
	NMI	0.78, —	0.76±0.03, —	0.76±0.03, —	0.92±0.00, 0.90±0.01	0.91±0.00, 0.91±0.00	0.93±0.00, 0.91±0.00	0.93±0.00, 0.91±0.00
	ARI	0.66, —	0.66±0.06, —	0.59±0.10, —	0.93±0.00, 0.90±0.01	0.92±0.00, 0.91±0.00	0.92±0.00, 0.91±0.00	0.92±0.01, 0.91±0.00
4.4%	ACC	—	0.81±0.03, —	0.66±0.05, —	0.91±0.00, 0.91±0.00	0.94±0.01, 0.94±0.01	0.89±0.01, 0.90±0.01	0.92±0.02, 0.92±0.02
	NMI	—	0.75±0.01, —	0.73±0.02, —	0.85±0.01, 0.86±0.01	0.89±0.01, 0.89±0.01	0.80±0.00, 0.80±0.01	0.85±0.02, 0.85±0.02
	ARI	—	0.66±0.03, —	0.54±0.07, —	0.82±0.01, 0.82±0.01	0.88±0.02, 0.87±0.02	0.83±0.01, 0.84±0.01	0.87±0.01, 0.87±0.01
5.8%	ACC	—	0.80±0.03, —	0.73±0.07, —	0.85±0.03, 0.85±0.04	0.92±0.01, 0.92±0.01	0.81±0.00, 0.81±0.00	0.92±0.01, 0.92±0.01
	NMI	—	0.74±0.01, —	0.76±0.04, —	0.81±0.01, 0.82±0.01	0.86±0.01, 0.87±0.00	0.72±0.01, 0.72±0.01	0.84±0.01, 0.85±0.02
	ARI	—	0.65±0.03, —	0.60±0.06, —	0.75±0.03, 0.75±0.03	0.84±0.01, 0.84±0.01	0.79±0.01, 0.79±0.01	0.86±0.01, 0.87±0.01

F.3. Validating Identifiability Claims using Synthetic Data

We generate synthetic data with $N = 2000$, $K = 3$. The first 1000 data points act as the seen samples, and the rest act as unseen data. In order to generate valid membership \mathcal{M}^{\natural} , we sample N random unit vectors, followed by adding i.i.d element-wise Gaussian noise with mean of 0 and variance of 0.1. The resulting matrix is then truncated to ensure that its elements are nonnegative. Lastly, the columns of \mathcal{M}^{\natural} are normalized using their respective ℓ_1 -norms to become valid probability mass functions. For constructing feature vectors \mathbf{x}^{\natural} corresponding to membership \mathbf{m}^{\natural} , we choose a simple nonlinear function \mathbf{f}



(a) Average MSE of VanillaDCC of 5 random trials against different M 's. (b) Median MSE of 10 random trials against different M 's.

who is defined through its inverse, i.e., \mathbf{f}^{-1} defined as: $\mathbf{f}^{-1}(\mathbf{m}^{\mathfrak{h}}) \triangleq [2m_1^{\mathfrak{h}}, 3m_2^{\mathfrak{h}} + 1, m_1^{\mathfrak{h}}m_2^{\mathfrak{h}} + m_3^{\mathfrak{h}} - 2]^{\top}$. Such construction guarantees the existence of a function mapping from feature vector $\mathbf{x} \in \mathcal{X}$ to its corresponding membership vector \mathbf{m} .

To acquire the pairwise annotations, we sample M pairs of indices uniformly. The pair similarity labels are then drawn from the Bernoulli distribution as described in Section 3.1. The performance of estimated membership $\widehat{\mathbf{M}}$ is measured using the *mean squared error* (MSE), which is defined as

$$\text{MSE}(\widehat{\mathbf{M}}, \mathbf{M}^{\mathfrak{h}}) = \frac{1}{K} \min_{\Pi} \sum_{k=1}^K \left\| \frac{\Pi \mathbf{M}^{\mathfrak{h}}(k, :)}{\|\Pi \mathbf{M}^{\mathfrak{h}}(k, :)\|} - \frac{\widehat{\mathbf{M}}(k, :)}{\|\widehat{\mathbf{M}}(k, :)\|} \right\|^2, \quad \Pi \text{ is permutation matrix.}$$

Fig. 3a shows the median of the MSEs on the training and test sets of the VanillaDCC over 5 random trials. One can see the MSE is fairly low and decreases as M increases—which is consistent with our theory.

To simulate noisy pairwise constraints, we keep the above setting and include a confusion matrix \mathbf{A} as follows

$$\mathbf{A} = \begin{bmatrix} 1.0 & 0.2 & 0.3 \\ 0.0 & 0.8 & 0.3 \\ 0.0 & 0.0 & 0.4 \end{bmatrix}.$$

Fig. 3b shows median MSE over 10 trials on seen and unseen data of VanillaDCC and VolMaxDCC using different number of pairwise constraints. One can see that using $\lambda > 0$, i.e., by promoting maximal volume of solutions, the MSE performance is much better, in the presence of noisy annotations. This shows the usefulness of the volume regularization in enhancing membership identifiability.

F.4. Results with Standard Deviation

In Tables 1, 2, 3, we only reported the average to avoid too-dense looking table. We report those results here with standard deviation in Table 8,9,10, respectively.

Table 8: Clustering performance of (seen data, unseen data) on STL10; $N_{\text{unseen}} = 2000$.

Methods		Kmeans	COP-Kmeans	PCKmeans	DC-GMM	C-IDEC	VanillaDCC	VolMaxDCC
0.0%	ACC	0.71±0.03, —	0.66±0.04, —	0.70±0.08, —	0.88 ± 0.04, 0.87±0.04	0.89±0.01, 0.88±0.02	0.93±0.00, 0.91±0.00	0.91±0.03, 0.89±0.03
	NMI	0.75±0.01, —	0.67±0.02, —	0.72±0.02, —	0.82±0.02, 0.80±0.02	0.81±0.01, 0.80±0.02	0.84±0.01, 0.81±0.01	0.82±0.03, 0.80±0.03
	ARI	0.54±0.03, —	0.52±0.03, —	0.55±0.06, —	0.80±0.04, 0.78±0.03	0.79±0.02, 0.77±0.02	0.85±0.00, 0.83±0.00	0.84±0.01, 0.82±0.01
8.3%	ACC	—	0.70±0.04, —	0.70±0.04, —	0.75±0.03, 0.76±0.03	0.77±0.01, 0.79±0.01	0.78±0.00, 0.80±0.01	0.80±0.01, 0.81±0.00
	NMI	—	0.64±0.03, —	0.71±0.02, —	0.67±0.01, 0.69±0.01	0.67±0.01, 0.69±0.01	0.59±0.00, 0.62±0.01	0.64±0.02, 0.65±0.02
	ARI	—	0.51±0.03, —	0.56±0.04, —	0.57±0.02, 0.59±0.02	0.59±0.02, 0.61±0.01	0.69±0.00, 0.71±0.01	0.73±0.02, 0.74±0.02
10.3%	ACC	—	0.62±0.05, —	0.69±0.04, —	0.70±0.02, 0.72±0.03	0.70±0.03, 0.71±0.03	0.72±0.01, 0.73±0.01	0.79±0.00, 0.81±0.00
	NMI	—	0.59±0.02, —	0.73±0.01, —	0.62±0.02, 0.64±0.02	0.60±0.02, 0.62±0.02	0.50±0.01, 0.51±0.01	0.68±0.00, 0.70±0.00
	ARI	—	0.44±0.03, —	0.55±0.02, —	0.51±0.01, 0.52±0.02	0.50±0.03, 0.52±0.03	0.62±0.00, 0.64±0.01	0.77±0.01, 0.78±0.00
15.0%	ACC	—	0.62±0.05, —	0.64±0.07, —	0.60±0.02, 0.61±0.02	0.57±0.01, 0.57±0.01	0.56±0.03, 0.58±0.03	0.79±0.00, 0.81±0.00
	NMI	—	0.54±0.01, —	0.72±0.01, —	0.54±0.01, 0.55±0.02	0.50±0.02, 0.50±0.02	0.33±0.01, 0.35±0.01	0.68±0.01, 0.69±0.01
	ARI	—	0.41±0.02, —	0.52±0.05, —	0.38±0.02, 0.39±0.01	0.38±0.01, 0.39±0.02	0.50±0.01, 0.51±0.01	0.76±0.02, 0.77±0.02

Deep Clustering With Incomplete Noisy Pairwise Annotations

Table 9: Clustering performance of (seen data, unseen data) on CIFAR10; $N_{\text{unseen}} = 45000$.

Methods		Kmeans	COP-Kmeans	PCKmeans	DC-GMM	C-IDEA	VanillaDCC	VolMaxDCC
0.0%	ACC	0.78±0.00, —	0.67±0.06, —	0.67±0.06, —	0.91±0.01, 0.89±0.01	0.90±0.01, 0.89±0.01	0.92±0.00, 0.90±0.00	0.91±0.01, 0.90±0.00
	NMI	0.71±0.00, —	0.66±0.02, —	0.71±0.02, —	0.83±0.01, 0.81±0.01	0.83±0.01, 0.81±0.01	0.84±0.00, 0.80±0.00	0.83±0.01, 0.80±0.01
	ARI	0.62±0.00, —	0.54±0.05, —	0.55±0.05, —	0.82±0.01, 0.79±0.01	0.81±0.01, 0.79±0.01	0.85±0.00, 0.81±0.00	0.84±0.00, 0.82±0.00
4.9%	ACC	—	0.75±0.04, —	0.70±0.06, —	0.86±0.00, 0.86±0.00	0.86±0.00, 0.86±0.00	0.86±0.00, 0.86±0.00	0.86±0.00, 0.86±0.00
	NMI	—	0.69±0.01, —	0.69±0.02, —	0.77±0.00, 0.77±0.00	0.77±0.00, 0.77±0.00	0.73±0.00, 0.73±0.00	0.73±0.00, 0.74±0.00
	ARI	—	0.60±0.03, —	0.57±0.04, —	0.73±0.00, 0.74±0.00	0.73±0.00, 0.74±0.00	0.77±0.00, 0.77±0.00	0.77±0.00, 0.77±0.00
8.7%	ACC	—	0.64±0.06, —	0.72±0.06, —	0.76±0.04, 0.76±0.04	0.76±0.02, 0.76±0.02	0.76±0.02, 0.77±0.02	0.83±0.01, 0.83±0.01
	NMI	—	0.63±0.02, —	0.69±0.02, —	0.71±0.02, 0.71±0.02	0.70±0.01, 0.70±0.01	0.58±0.01, 0.59±0.01	0.70±0.01, 0.70±0.01
	ARI	—	0.49±0.04, —	0.59±0.03, —	0.58±0.03, 0.59±0.03	0.59±0.01, 0.60±0.01	0.70±0.01, 0.70±0.01	0.75±0.01, 0.75±0.01
10.9%	ACC	—	0.68±0.06, —	0.70±0.04, —	0.74±0.03, 0.74±0.02	0.73±0.02, 0.73±0.02	0.68±0.00, 0.69±0.01	0.82±0.01, 0.82±0.01
	NMI	—	0.61±0.03, —	0.68±0.01, —	0.67±0.02, 0.68±0.02	0.65±0.01, 0.66±0.01	0.48±0.01, 0.49±0.01	0.68±0.02, 0.68±0.02
	ARI	—	0.50±0.05, —	0.57±0.03, —	0.55±0.03, 0.55±0.03	0.56±0.03, 0.57±0.03	0.62±0.01, 0.63±0.01	0.74±0.01, 0.74±0.01

Table 10: Clustering performance of (seen data, unseen data) on ImageNet10; $N_{\text{unseen}} = 2000$.

Methods		Kmeans	COP-Kmeans	PCKmeans	DC-GMM	C-IDEA	VanillaDCC	VolMaxDCC
0.0%	ACC	0.85±0.00, —	0.79±0.06, —	0.70±0.06, —	0.97±0.00, 0.96±0.00	0.97±0.00, 0.96±0.00	0.97±0.00, 0.96±0.00	0.97±0.00, 0.96±0.00
	NMI	0.80±0.00, —	0.77±0.03, —	0.75±0.04, —	0.93±0.00, 0.91±0.00	0.93±0.00, 0.92±0.00	0.94±0.00, 0.91±0.00	0.94±0.00, 0.91±0.00
	ARI	0.68±0.00, —	0.66±0.05, —	0.55±0.07, —	0.94±0.00, 0.92±0.00	0.94±0.00, 0.92±0.00	0.93±0.00, 0.91±0.00	0.93±0.00, 0.91±0.00
3.4%	ACC	—	0.79±0.09, —	0.66±0.12, —	0.93±0.01, 0.92±0.01	0.93±0.00, 0.93±0.00	0.92±0.01, 0.91±0.01	0.94±0.01, 0.94±0.01
	NMI	—	0.75±0.04, —	0.73±0.05, —	0.86±0.01, 0.85±0.01	0.87±0.01, 0.86±0.01	0.83±0.02, 0.82±0.02	0.88±0.01, 0.87±0.02
	ARI	—	0.65±0.07, —	0.52±0.12, —	0.84±0.01, 0.83±0.01	0.86±0.01, 0.85±0.01	0.84±0.01, 0.84±0.01	0.89±0.01, 0.88±0.01
6.9%	ACC	—	0.74±0.07, —	0.72±0.08, —	0.84±0.01, 0.84±0.01	0.88±0.01, 0.88±0.01	0.84±0.00, 0.84±0.00	0.92±0.00, 0.91±0.00
	NMI	—	0.70±0.03, —	0.76±0.04, —	0.79±0.01, 0.79±0.00	0.82±0.01, 0.82±0.01	0.70±0.01, 0.70±0.00	0.84±0.00, 0.83±0.00
	ARI	—	0.58±0.05, —	0.59±0.09, —	0.71±0.01, 0.71±0.01	0.77±0.02, 0.77±0.02	0.77±0.01, 0.77±0.01	0.88±0.00, 0.87±0.00
11.2%	ACC	—	0.72±0.05, —	0.62±0.07, —	0.71±0.03, 0.72±0.03	0.80±0.02, 0.81±0.02	0.65±0.01, 0.66±0.01	0.91±0.01, 0.90±0.01
	NMI	—	0.64±0.02, —	0.73±0.03, —	0.68±0.01, 0.70±0.01	0.74±0.02, 0.76±0.02	0.49±0.03, 0.52±0.02	0.83±0.01, 0.82±0.01
	ARI	—	0.54±0.03, —	0.51±0.09, —	0.56±0.02, 0.58±0.02	0.66±0.03, 0.68±0.03	0.62±0.03, 0.63±0.02	0.87±0.01, 0.86±0.00