# Towards Understanding the Mechanism of Contrastive Learning via Similarity Structure: A Theoretical Analysis

Hiroki Waida[1]
waida.h.aa@m.titech.ac.jp

Yuichiro Wada[2,3]
wada.yuichiro@fujitsu.com

Léo Andéol[4,5,6,7]
leo.andeol@math.univ-toulouse.fr

Takumi Nakagawa[1,3]
nakagawa.t.as@m.titech.ac.jp

Yuhui Zhang[1]
zhang.y.av@m.titech.ac.jp

Takafumi Kanamori[1,3]
kanamori@c.titech.ac.jp

[1]Tokyo Institute of Technology, Japan
[2]Fujitsu, Japan
[3]RIKEN AIP, Japan
[4]Institut de Mathématiques de Toulouse, France
[5]SNCF, France
[6]Université de Toulouse, France
[7]CNRS, France

## Abstract

Contrastive learning is an efficient approach to self-supervised representation learning. Although recent studies have made progress in the theoretical understanding of contrastive learning, the investigation of how to characterize the clusters of the learned representations is still limited. In this paper, we aim to elucidate the characterization from theoretical perspectives. To this end, we consider a kernel-based contrastive learning framework termed Kernel Contrastive Learning (KCL), where kernel functions play an important role when applying our theoretical results to other frameworks. We introduce a formulation of the similarity structure of learned representations by utilizing a statistical dependency viewpoint. We investigate the theoretical properties of the kernel-based contrastive loss via this formulation. We first prove that the formulation characterizes the structure of representations learned with the kernel-based contrastive learning framework. We show a new upper bound of the classification error of a downstream task, which explains that our theory is consistent with the empirical success of contrastive learning. We also establish a generalization error bound of KCL. Finally, we show a guarantee for the generalization ability of KCL to the downstream classification task via a surrogate bound.

## 1 Introduction

Recently, many studies on self-supervised representation learning have been paying much attention to contrastive learning (Caron et al., 2020; Chen et al., 2020a,b; Dwibedi et al., 2021; HaoChen

1

et al., 2021; Li et al., 2021). Through contrastive learning, encoder functions acquire how to encode unlabeled data to good representations by utilizing some information of similarity behind the data, where recent works (Chen et al., 2020a,b; Dwibedi et al., 2021) use several data augmentation techniques to produce pairs of similar data. It is empirically shown by many works (Caron et al., 2020; Chen et al., 2020a,b; Dwibedi et al., 2021; HaoChen et al., 2021) that contrastive learning produces effective representations that are fully adaptable to downstream tasks, such as classification and transfer learning.

Besides the practical development of contrastive learning, the theoretical understanding of contrastive learning is essential to construct more efficient self-supervised learning algorithms. In this paper, we tackle the following fundamental question of contrastive learning from the theoretical side: *How are the clusters of feature vectors output from an encoder model pretrained by contrastive learning characterized?*

Recently, several works have shed light on several theoretical perspectives on this problem to study the generalization guarantees of contrastive learning to downstream classification tasks (Arora et al., 2019; Dufumier et al., 2022; HaoChen and Ma, 2023; HaoChen et al., 2021; Huang et al., 2023; Wang et al., 2022a; Zhao et al., 2023). One of the primary approaches of these works is to introduce some similarity measures in the data. Arora et al. (2019) has introduced the conditional independence assumption, which assumes that data $x$ and its positive data $x^+$ are sampled independently according to the conditional probability distribution $\mathcal{D}_c$, given the *latent class $c$* drawn from the latent class distribution. Although the concepts of latent classes and conditional independence assumption are often utilized to formulate semantic similarity of $x$ and $x^+$ (Arora et al., 2019; Ash et al., 2022; Awasthi et al., 2022; Bao et al., 2022; Zou and Liu, 2023), it is pointed out by several works (HaoChen et al., 2021; Wang et al., 2022a) that this assumption can be violated in practice. Several works (HaoChen et al., 2021; Wang et al., 2022a) have introduced different ideas about the similarity between data to alleviate this assumption. HaoChen et al. (2021) have introduced the notion called *population augmentation graph* to provide a theoretical analysis for Spectral Contrastive Learning (SCL) without the conditional independence assumption on $x$ and $x^+$. Some works also focus on various graph structures (Dufumier et al., 2022; HaoChen and Ma, 2023; Wang et al., 2022a). Although these studies give interesting insights into contrastive learning, the applicable scope of their analyses is limited to specific objective functions. Recently, several works (Huang et al., 2023; Zhao et al., 2023) consider the setup where raw data in the same latent class are aligned well in the sense that the *augmented distance* is small enough. Although their theoretical guarantees can apply to multiple contrastive learning frameworks, their assumptions on the function class of encoders are strong, and it needs to be elucidated whether their guarantees can hold in practice. Therefore, the investigation of the above question from unified viewpoints is ongoing, and more perspectives are required to understand the structure learned by contrastive learning.

## 1.1 Contributions

In this paper, we aim to theoretically investigate the above question from a unified perspective by introducing a formulation based on a statistical similarity between data. The main contributions of this paper are summarized below:

1. Since we aim to elucidate the mechanism of contrastive learning, we need to consider a unified framework that can apply to others, not specific frameworks such as SimCLR (Chen et al., 2020a) and SCL (HaoChen et al., 2021). Li et al. (2021) pointed out that kernel-based self-supervised learning objectives are related to other contrastive losses, such as the InfoNCE loss (Chen et al., 2020a; van den Oord et al., 2018). Therefore, via a kernel-based contrastive learning framework, other frameworks can be investigated through the lens of kernels. Motivated by this, we utilize the framework termed *Kernel Contrastive Learning* (KCL) as a tool for achieving the goal. The loss of KCL, which is called *kernel contrastive loss*, is a contrastive loss that has a simple and general form, where the similarity between two feature vectors is measured by a reproducing kernel (Aronszajn, 1950; Berlinet and Thomas-Agnan, 2004; Steinwart and Christmann, 2008) (Section 3). One of our contributions is employing KCL to study the mechanism of contrastive learning from a new unified theoretical perspective.

2. We introduce a new formulation of similarity between data (Section 4). Our formulation of similarity begins with the following intuition: if raw or augmented data $x$ and $x'$ belong to the same class, then the similarity measured by some function should be higher than a threshold. Following this, we introduce a formulation (Assumption 2) based on the similarity function (2).

3. We present the theoretical analyses towards elucidating the above question (Section 5). We first show that KCL can distinguish the clusters of representations according to this formulation (Section 5.1). This result shows that our formulation is closely connected to the mechanism of contrastive learning. Next, we establish a new upper bound for the classification error of the downstream task (Section 5.2), which indicates that our formulation does not contradict the practical efficiency of contrastive learning shown by a line of work (Chen et al., 2020a,b; Dwibedi et al., 2021; HaoChen et al., 2021). Notably, our upper bound is valid under more realistic assumptions on the encoder functions, compared to the previous works (Huang et al., 2023; Zhao et al., 2023). We also establish the generalization error bound for KCL (Section 5.3). Finally, applying our theoretical results, we show a guarantee for the generalization of KCL to the downstream classification task via a surrogate bound (Section 5.4).

## 1.2 Related Work

Contrastive learning methods have been investigated from the empirical side (Caron et al., 2020; Chen et al., 2020a, 2021, 2020b). Chen et al. (2020a) propose a method called SimCLR, which utilizes a variant of InfoNCE (Chen et al., 2020a; van den Oord et al., 2018). Several works have recently improved contrastive methods from various viewpoints (Caron et al., 2020; Dwibedi et al., 2021; Robinson et al., 2021a,b). Contrastive learning is often utilized in several fundamental tasks, such as clustering (Van Gansbeke et al., 2020) and domain adaptation (Singh, 2021), and applied to some domains such as vision (Chen et al., 2020a), natural language processing (Gao et al., 2021), and speech (Jiang et al., 2021). Besides the contrastive methods, several works (Chen and He, 2021; Grill et al., 2020) also study non-contrastive methods. Investigation toward the theoretical understanding of contrastive learning is also a growing focus. For instance, the generalization ability of contrastive learning to the downstream classification task has been investigated from many kinds

of settings (Arora et al., 2019; Bao et al., 2022; HaoChen and Ma, 2023; HaoChen et al., 2021; Huang et al., 2023; Saunshi et al., 2022; Tosh et al., 2021b; Wang et al., 2022a; Zhao et al., 2023). Several works investigate contrastive learning from various theoretical and empirical viewpoints to elucidate its mechanism, such as the geometric properties of contrastive losses (Huang et al., 2023; Wang and Isola, 2020), formulation of similarity between data (Arora et al., 2019; Dufumier et al., 2022; HaoChen et al., 2021; Huang et al., 2023; von Kügelgen et al., 2021; Wang et al., 2022a; Zhao et al., 2023), inductive bias (HaoChen and Ma, 2023; Saunshi et al., 2022), transferablity (HaoChen et al., 2022; Shen et al., 2022; Zhao et al., 2023), feature suppression (Chen et al., 2021; Robinson et al., 2021a), negative sampling methods (Chuang et al., 2020; Robinson et al., 2021b), and optimization viewpoints (Tian, 2022; Wen and Li, 2021).

Several works (Dufumier et al., 2022; Johnson et al., 2023; Kiani et al., 2022; Li et al., 2021; Tsai et al., 2022; Zhang et al., 2022) study the connection between contrastive learning and the theory of kernels. Li et al. (2021) investigate some contrastive losses, such as InfoNCE, from a kernel perspective. Zhang et al. (2022) show a relation between the kernel method and $f$-mutual information and apply their theory to contrastive learning. Tsai et al. (2022) tackle the conditional sampling problem using kernels as similarity measurements. Dufumier et al. (2022) consider incorporating prior information in contrastive learning by using the theory of kernel functions. Kiani et al. (2022) connect several self-supervised learning algorithms to kernel methods through optimization problem viewpoints. Note that different from these works, our work employs kernel functions to investigate a new unified perspective of contrastive learning via the statistical similarity.

Many previous works investigate various interpretations of self-supervised representation learning objectives. For instance, the InfoMax principle (Poole et al., 2019; Tschannen et al., 2020), spectral clustering (HaoChen et al., 2021) (see Ng et al. (2002) for spectral clustering), and Hilbert-Schmidt Independence Criterion (HSIC) (Li et al., 2021) (see Gretton et al. (2005) for HSIC). However, the investigation of contrastive learning from unified perspectives is worth addressing to elucidate its mechanism, as recent works on self-supervised representation learning tackle it from the various standpoints (Dubois et al., 2022; Huang et al., 2023; Johnson et al., 2023; Kiani et al., 2022; Tian, 2022).

## 2 Preliminaries

### 2.1 Problem Setup

We give the standard setup of contrastive learning. Our setup closely follows that of HaoChen et al. (2021), though we also introduce additional technically necessary settings to maintain the mathematical rigorousness. Let $\overline{\mathbb{X}} \subset \mathbb{R}^p$ be a topological space consisting of raw data, and let $P_{\overline{\mathbb{X}}}$ be a Borel probability measure on $\overline{\mathbb{X}}$. A line of work on contrastive learning (Chen et al., 2020a,b; Dwibedi et al., 2021; HaoChen et al., 2021) uses data augmentation techniques to obtain similar augmented data points. Hence, we define a set $\mathcal{T}$ of maps transforming a point $\overline{x} \in \overline{\mathbb{X}}$ into $\mathbb{R}^p$, where we assume that $\mathcal{T}$ includes the identity map on $\mathbb{R}^p$. Then, let us define $\mathbb{X} = \bigcup_{t \in \mathcal{T}} \{t(\overline{x}) : \overline{x} \in \overline{\mathbb{X}}\}$. Every element $t$ in $\mathcal{T}$ can be regarded as a map returning an augmented data $x = t(\overline{x})$ for a raw data point $\overline{x} \in \overline{\mathbb{X}}$. Note that since the identity map belongs to $\mathcal{T}$, $\overline{\mathbb{X}}$ is a subset of $\mathbb{X}$. We endow $\mathbb{X}$ with some topology. Let $\nu_{\mathbb{X}}$ be a $\sigma$-finite and non-negative Borel

measure in $\mathbb{X}$. Following the idea of HaoChen et al. (2021), we denote $a(x|\overline{x})$ as the conditional probability density function of $x$ given $\overline{x} \sim P_{\overline{\mathbb{X}}}$ and define the weight function $w : \mathbb{X} \times \mathbb{X} \to \mathbb{R}$ as $w(x, x') = \mathbb{E}_{\overline{x} \sim P_{\overline{\mathbb{X}}}}[a(x|\overline{x})a(x'|\overline{x})]$. From the definition, $w$ is a joint probability density function on $\mathbb{X} \times \mathbb{X}$. Let us define the marginal $w(\cdot)$ of the weight function to be $w(x) = \int w(x, x')d\nu_{\mathbb{X}}(x')$. The marginal $w(\cdot)$ is also a probability density function on $\mathbb{X}$, and the corresponding probability measure is denoted by $dP_{\mathbb{X}}(x) = w(x)d\nu_{\mathbb{X}}(x)$. Denote by $\mathbb{E}_{x,x^+}[\cdot], \mathbb{E}_{x,x^-}[\cdot]$ respectively, the expectation w.r.t. the probability measure $w(x, x')d\nu_{\mathbb{X}}^{\otimes 2}(x, x'), w(x)w(x')d\nu_{\mathbb{X}}^{\otimes 2}(x, x')$ on $\mathbb{X} \times \mathbb{X}$, where $\nu_{\mathbb{X}}^{\otimes 2} := \nu_{\mathbb{X}} \otimes \nu_{\mathbb{X}}$ is the product measure on $\mathbb{X} \times \mathbb{X}$. To rigorously formulate our framework of contrastive learning, we assume that the marginal $w$ is positive on $\mathbb{X}$. Indeed, a point $x \in \mathbb{X}$ satisfying $w(x) = 0$ is not included in the support of $a(\cdot|\overline{x})$ for $P_{\overline{\mathbb{X}}}$-almost surely $\overline{x} \in \overline{\mathbb{X}}$, which means that such a point $x$ merely appears as augmented data.

Let $f_0 : \mathbb{X} \to \mathbb{R}^d$ be an encoder model mapping augmented data to the feature space, and let $\mathcal{F}_0$ be a class of functions consisting of such encoders. In practice, $f_0$ is defined by a backbone architecture (e.g., ResNet (He et al., 2016); see Chen et al. (2020a)), followed by the additional multi-layer perceptrons called *projection head* (Chen et al., 2020a). We assume that $\mathcal{F}_0$ is uniformly bounded, i.e., there exists a universal constant $c \in \mathbb{R}$ such that $\sup_{f_0 \in \mathcal{F}_0} \sup_{x \in \mathbb{X}} \|f_0(x)\|_2 \leq c$. For instance, a function space of bias-free fully connected neural networks on a bounded domain, where every neural network has the continuous activate function at each layer, satisfies this condition. Since a feature vector output from the encoder model is normalized using the Euclidean norm in many empirical studies (Chen et al., 2020a; Dwibedi et al., 2021) and several theoretical studies (Wang and Isola, 2020; Wang et al., 2022a), we consider the function space of normalized functions $\mathcal{F} = \{f \mid \exists f_0 \in \mathcal{F}_0, \ f(x) = f_0(x)/\|f_0(x)\|_2 \text{ for } \forall x \in \mathbb{X}\}$. Here, to guarantee that every $f \in \mathcal{F}$ is well-defined, suppose that $\mathfrak{m}(\mathcal{F}_0) := \inf_{f \in \mathcal{F}} \inf_{x \in \mathbb{X}} \|f_0(x)\|_2 > 0$ holds.

Finally, we introduce several notations used throughout this paper. Let $\mathbb{M} \subset \mathbb{X}$ be a measurable set, then we write

$$\mathbb{E}[f(x)|\mathbb{M}] := \int_{\mathbb{X}} f(x)P_{\mathbb{X}}(dx|\mathbb{M}) = P_{\mathbb{X}}(\mathbb{M})^{-1}\int_{\mathbb{M}} f(x)w(x)d\nu_{\mathbb{X}}(x).$$

We also use the notation $\mathbb{E}[f(x); \mathbb{M}] := \int_{\mathbb{M}} f(x)w(x)d\nu_{\mathbb{X}}(x)$. Denote by $\mathbb{1}_{\mathbb{M}}(\cdot)$, the indicator function of a set $\mathbb{M}$. We also use $[n] := \{1, \cdots, n\}$ for $n \in \mathbb{N}$.

## 2.2 Reproducing Kernels

We provide several notations of reproducing kenrels (Aronszajn, 1950; Berlinet and Thomas-Agnan, 2004; Steinwart and Christmann, 2008). Let $k : \mathbb{S}^{d-1} \times \mathbb{S}^{d-1} \to \mathbb{R}$ be a real-valued, continuous, symmetric, and positive-definite kernel, where $\mathbb{S}^{d-1}$ denotes the unit hypersphere centered at the origin $\mathbf{0} \in \mathbb{R}^d$, and the positive-definiteness means that for every $\{z_i\}_{i=1}^n \subset \mathbb{S}^{d-1}$ and $\{c_i\}_{i=1}^n \subset \mathbb{R}$, $\sum_{i,j=1}^n c_i c_j k(z_i, z_j) \geq 0$ holds (Berlinet and Thomas-Agnan, 2004). Let $\mathcal{H}_k$ be the Reproducing Kernel Hilbert Space (RKHS) with kernel $k$ (Aronszajn, 1950), which satisfies $\phi(z) = \langle \phi, k(\cdot, z) \rangle_{\mathcal{H}_k}$ for all $\phi \in \mathcal{H}_k$ and $z \in \mathbb{S}^{d-1}$. Denote $h(z) = k(\cdot, z)$ for $z \in \mathbb{S}^{d-1}$, where such a map is often called feature map (Steinwart and Christmann, 2008). Here, we impose the following condition.

**Assumption 1.** *For the kernel function $k$, there exists some $\rho$-Lipschitz function $\psi : [-1, 1] \to \mathbb{R}$ such that for every $z, z' \in \mathbb{S}^{d-1}$, $k(z, z') = \psi(z^\top z')$ holds.*

Several popular kernels in machine learning such as the linear kernel, quadratic kernel, and Gaussian kernel, satisfy Assumption 1. Note that the Lipschitz condition in Assumption 1 is useful to derive the generalization error bound for the kernel contrastive loss (see Section 5.3), and sometimes it can be removed when analyzing for a specific kernel. We use this assumption to present general results. Here, we also use the following notion in this paper:

**Proposition 1.** *Let $\mathbb{M} \subset \mathbb{X}$ be a measurable set and $f \in \mathcal{F}$. Define $\mu_{\mathbb{M}}(f) := \mathbb{E}_{P_{\mathbb{X}}}[h(f(x))|\mathbb{M}]$. Then, $\mu_{\mathbb{M}}(f) \in \mathcal{H}_k$.*

The quantity $\mu_{\mathbb{M}}(f)$ with a measurable subset $\mathbb{M}$ can be regarded as a variant of the kernel mean embedding (Muandet et al., 2017). The proof of Proposition 1 is a slight modification of Muandet et al. (2017); see Appendix A.

# 3 Kernel Contrastive Learning

In this section, we introduce a contrastive learning framework to analyze the mechanism of contrastive learning. In representation learning, the InfoNCE loss (Chen et al., 2020a; van den Oord et al., 2018) are widely used in application domains such as vision (Chen et al., 2020a, 2021; Dwibedi et al., 2021). Following previous works (Bao et al., 2022; Chen et al., 2020a; van den Oord et al., 2018), we define the InfoNCE loss as,

$$L_{\mathrm{NCE}}(f;\tau) = -\mathbb{E}_{\substack{x,x^+ \\ \{x_i^-\} \sim P_{\mathbb{X}}}} \left[ \log \frac{e^{f(x)^\top f(x^+)/\tau}}{e^{f(x)^\top f(x^+)/\tau} + \sum_{i=1}^M e^{f(x)^\top f(x_i^-)/\tau}} \right],$$

where $\{x_i^-\}$ are i.i.d. random variables, $\tau > 0$, and $M$ is the number of negative samples. Wang and Isola (2020) introduce the asymptotic of the InfoNCE loss:

$$L_{\infty\text{-NCE}}(f;\tau) = -\mathbb{E}_{x,x^+}\left[ \frac{f(x)^\top f(x^+)}{\tau} \right] + \mathbb{E}_x\left[ \log \mathbb{E}_{x'}\left[ e^{\frac{f(x)^\top f(x')}{\tau}} \right] \right].$$

According to the theoretical analysis of Wang and Isola (2020), they show that the first term represents the *alignment*, i.e., the averaged closeness of feature vectors of the pair $(x, x^+)$, while the second one indicates the *uniformity*, i.e., how far apart the feature vectors of negative samples $x, x'$ are. Besides, Chen et al. (2021) report the efficiency of the generalized contrastive losses, which have the additional weight hyperparameter. Meanwhile, since we aim to study the mechanism of contrastive learning, a simple and general form of contrastive losses related to other frameworks is required. Here, Li et al. (2021) find the connection between self-supervised learning and kernels by showing that some HSIC criterion is proportional to the objective function $\mathbb{E}_{x,x^+}[k(f(x), f(x^+))] - \mathbb{E}_{x,x^-}[k(f(x), f(x^-))]$ (for more detail, see Appendix F.2). Motivated by this connection, we consider a contrastive learning objective where a kernel function measures the similarity of the feature vectors of augmented data points. More precisely, for the kernel function $k$ introduced in Section 2.2, we define the kernel contrastive loss as,

$$L_{\mathrm{KCL}}(f;\lambda) = -\mathbb{E}_{x,x^+}\left[ k(f(x), f(x^+)) \right] + \lambda \mathbb{E}_{x,x^-}\left[ k(f(x), f(x^-)) \right],$$

where the weight hyperparameter $\lambda$ is inspired by Chen et al. (2021). Here, the kernel contrastive loss $L_{\mathrm{KCL}}$ is minimized during the pretraining stage of contrastive learning. Throughout this paper,

the contrastive learning framework with the kernel contrastive loss is called *Kernel Contrastive Learning* (KCL).

Next, we show the connections to other contrastive learning objectives. First, for the InfoNCE loss, we consider the linear kernel contrastive loss $L_{\mathrm{LinKCL}}(f; \lambda)$ defined by selecting $k(z, z') = z^\top z'$. Note that $L_{\mathrm{LinKCL}}$ and its empirical loss are also discussed in several works (Huang et al., 2023; Wang and Liu, 2021). For $L_{\mathrm{LinKCL}}(f; 1)$, we have,

$$\tau^{-1} L_{\mathrm{LinKCL}}(f; 1) \leq L_{\mathrm{NCE}}(f; \tau) + \log M^{-1}. \tag{1}$$

In Appendix F.3, we show a generalized inequality of (1) for the generalized loss (Chen et al., 2021). Note that similar relations hold when $L_{\mathrm{NCE}}$ is replaced with the asymptotic loss (Wang and Isola, 2020) or decoupled contrastive learning loss (Yeh et al., 2022); see Appendix F.3. Therefore, it is possible to analyze the InfoNCE loss and its variants via $L_{\mathrm{LinKCL}}(f; \lambda)$.

The kernel contrastive loss is also related to other contrastive learning objectives. For instance, the quadratic kernel contrastive loss with the quadratic kernel $k(z, z') = (z^\top z')^2$ becomes a lower bound of the spectral contrastive loss (HaoChen et al., 2021) up to an additive constant (see Appendix F.3). Thus, theoretical analyses of the kernel contrastive loss can apply to other contrastive learning objectives.

Note that we empirically demonstrate that the KCL frameworks with the Gaussian kernel and quadratic kernel work, although simple; see Appendix H.1 for the experimental setup and Appendix H.2 and H.3 for the results in the supplementary material. The experimental results also motivate us to use KCL as a theoretical tool for studying contrastive learning.

# 4 A Formulation Based on Statistical Similarity

## 4.1 Key Ingredient: Similarity Function

To study the mechanism of contrastive learning, we introduce a notion of similarity between two augmented data points, which is a key component in our analysis. Let us define,

$$\mathrm{sim}(x, x'; \lambda) := \frac{w(x, x')}{w(x) w(x')} - \lambda, \tag{2}$$

where $\lambda \geq 0$ is the weight parameter of $L_{\mathrm{KCL}}$, and $w(x, x')$ and $w(x)$ have been introduced in Section 2.1. Note that (2) is well-defined since $w(x) > 0$ holds for every $x \in \mathbb{X}$. The quantity $\mathrm{sim}(x, x'; \lambda)$ represents how much statistical dependency $x$ and $x'$ have. The density ratio $w(x, x')/(w(x) w(x'))$ can be regarded as an instance of *point-wise dependency* introduced by Tsai et al. (2020). The hyperparameter $\lambda$ controls the degree of relevance between two augmented data $x, x'$ via their (in-)dependency. For instance, with the fixed $\lambda = 1$, $\mathrm{sim}(x, x'; 1)$ is positive if $w(x, x') > w(x) w(x')$, i.e., $x$ and $x'$ are correlated.

Here, we note that several theoretical works on representation learning (Johnson et al., 2023; Tosh et al., 2021b; Wang et al., 2022b) use the density ratio to study the optimal representations of several contrastive learning objectives. Tosh et al. (2021b) focus on the fact that the minimizer of a logistic loss can be written in terms of the density ratio and utilize it to study *landmark*

7

*embedding* (Tosh et al., 2021a). Johnson et al. (2023) connect the density ratio to the minimizers of several contrastive learning objectives and investigate the quality of representations related to the minimizers. In addition, Wang et al. (2022b) study the minimizer of the spectral contrastive loss (HaoChen et al., 2021). Meanwhile, we emphasize that the purpose of using the density ratio in (2) is not to study the optimal representation of KCL but to give a formulation based on the statistical similarity between augmented data.

**Remark 1.** *We can show that the kernel contrastive loss can be regarded as a relaxation of the population-level Normalized Cut problem (Terada and Yamamoto, 2019), where the integral kernel is defined with (2). Thus, (2) defines the similarity structure utilized by KCL. Detailed arguments and comparison to related works (HaoChen et al., 2021; Tian, 2022) can be found in Appendix G.*

## 4.2 Formulation and Example

We introduce the following formulation based on our problem setting.

**Assumption 2.** *There exist some $\delta \in \mathbb{R}$, number of clusters $K \in \mathbb{N}$, measurable subsets $\mathbb{M}_1, \cdots, \mathbb{M}_K \subset \mathbb{X}$, and a deterministic labeling function $y : \mathbb{X} \to [K]$ such that the following conditions hold:*

**(A)** $\bigcup_{i=1}^{K} \mathbb{M}_i = \mathbb{X}$ *holds.*

**(B)** *For every $i \in [K]$, any points $x, x' \in \mathbb{M}_i$ satisfy $\text{sim}(x, x'; \lambda) \geq \delta$.*

**(C)** *For every $x \in \mathbb{X}$ and the set of indices $J_x = \{j \in [K] \mid x \in \mathbb{M}_j\}$, $y(x) \in J_x$ holds. Moreover, each set $\{x \in \mathbb{X} \mid y(x) = i\}$ is measurable.*

Assumption 2 does not require that $\mathbb{M}_1, \cdots, \mathbb{M}_K$ are disjoint, which is a realistic setting, as Wang et al. (2022a) show that clusters of augmented data can have inter-cluster connections depending on the strength of data augmentation. The conditions **(A)** and **(B)** in Assumption 2 guarantee that each subset $\mathbb{M}_i$ consists of augmented data that have high similarity. The condition **(C)** enables to incorporate label information in our analysis. Note that several works on contrastive learning (HaoChen and Ma, 2023; HaoChen et al., 2021; Saunshi et al., 2022) also employ deterministic labeling functions.

These conditions are useful to analyze the theory of contrastive learning. To gain more intuition, we provide a simple example that satisfies Assumption 2.

**Example 1** (The proof can be found in Appendix F.1)**.** *Suppose that $\mathbb{X}$ consists of disjoint open balls $\mathbb{B}_1, \cdots, \mathbb{B}_K$ of the same radius in $\mathbb{R}^p$, and $\overline{\mathbb{X}} = \mathbb{X}$. Let $a(x|\overline{x}) = \text{vol}(\mathbb{B}_1)^{-1} \sum_{i=1}^{K} \mathbb{1}_{\mathbb{B}_i \times \mathbb{B}_i}(x, \overline{x})$, and let $p_{\overline{\mathbb{X}}}(\overline{x}) = (K\text{vol}(\mathbb{B}_1))^{-1}$ be the probability density function of $P_{\overline{\mathbb{X}}}$. Then, $w(x) > 0$ and $\text{sim}(x, x'; \lambda) = K\mathbb{1}_{\bigcup_{i \in [K]} \mathbb{B}_i \times \mathbb{B}_i}(x, x') - \lambda$ hold. Hence, for instance let $\lambda = 1$ and $\delta = K - 1$, and take $\mathbb{M}_i := \mathbb{B}_i$. Also, define $y : \mathbb{X} \to [K]$ as $y(x) = i$ if $x \in \mathbb{B}_i$ for some $i \in [K]$. Then, Assumption 2 is satisfied in this setting.*

Here, theoretical formulations of similarity have been investigated by several works on contrastive learning (Arora et al., 2019; Dufumier et al., 2022; HaoChen and Ma, 2023; HaoChen

et al., 2021; Huang et al., 2023; Parulekar et al., 2023; von Kügelgen et al., 2021; Wang et al., 2022a; Zhao et al., 2023). The basic notions introduced by these works are: *latent classes* of unlabeled data and conditional independence assumption (Arora et al., 2019), graph structures of augmented data (Dufumier et al., 2022; HaoChen and Ma, 2023; HaoChen et al., 2021; Wang et al., 2022a), the decomposition of unlabled data into the *content* (i.e., invariant against data augmentation) and *style* (i.e., changeable by data augmentation) variables (Parulekar et al., 2023; von Kügelgen et al., 2021), and the geometric structure based on *augmented distance* and the concentration within class subsets (Huang et al., 2023; Zhao et al., 2023). Meanwhile, our formulation in Assumption 2 uses the similarity function (2), which differs from the previous works. Note that Assumption 2 has some relation to Assumption 3 in HaoChen and Ma (2023); see Appendix F.4. Our formulation gives deeper insights into contrastive learning, as shown in Section 5.

# 5 Theoretical Results

## 5.1 KCL as Representation Learning with Statistical Similarity

First, we connect the kernel contrastive loss to the formulation based on the similarity $\text{sim}(x, x'; \lambda)$. The following theorem indicates that KCL has two effects on the way of representation learning, where the clusters $\mathbb{M}_1, \cdots, \mathbb{M}_K$ involve to explain the mathematical relation.

**Theorem 1.** *Suppose that Assumption 1 and 2 hold. Take $\delta \in \mathbb{R}$, $K \in \mathbb{N}$, and $\mathbb{M}_1, \cdots, \mathbb{M}_K$ such that the conditions* **(A)** *and* **(B)** *in Assumption 2 are satisfied. Then, the following inequality holds for every $f \in \mathcal{F}$:*

$$\frac{\delta}{2} \cdot \mathfrak{a}(f) + \lambda \cdot \mathfrak{c}(f) \leq L_{\text{KCL}}(f; \lambda) + R(\lambda). \tag{3}$$

*where $R(\lambda)$ is a function of $\lambda$, $\mu_i(f) = \mu_{\mathbb{M}_i}(f)$, and*

$$\mathfrak{a}(f) := \sum_{i=1}^{K} \mathbb{E}_{x, x^-} \left[ \|h(f(x)) - h(f(x^-))\|_{\mathcal{H}_k}^2 ; \mathbb{M}_i \times \mathbb{M}_i \right],$$

$$\mathfrak{c}(f) := \sum_{i \neq j} P_{\mathbb{X}}(\mathbb{M}_i) P_{\mathbb{X}}(\mathbb{M}_j) \langle \mu_i(f), \mu_j(f) \rangle_{\mathcal{H}_k}.$$

For the proof of Theorem 1, see Appendix B.2. Note that the key point of the proof is the following usage of Assumption 2: for any $x, x' \in \mathbb{M}_i$, the inequality $\text{sim}(x, x'; \lambda) \geq \delta$ implies the relation $w(x, x') \geq (\lambda + \delta) w(x) w(x')$. Here we briefly explain each symbol in Theorem 1. The value $\mathfrak{a}(f)$ quantifies the concentration within each cluster consisting of the representations of augmented data in $\mathbb{M}_i$. The quantity $\mathfrak{c}(f)$ measures how far the subsets $h(f(\mathbb{M}_1)), \cdots, h(f(\mathbb{M}_K))$ are, since $\langle \mu_i(f), \mu_j(f) \rangle_{\mathcal{H}_k} = \int_{\mathbb{M}_i} \int_{\mathbb{M}_j} k(f(x), f(x')) P_{\mathbb{X}}(dx|\mathbb{M}_i) P_{\mathbb{X}}(dx'|\mathbb{M}_j)$ holds (see Lemma 2 in Appendix B.1). These quantities indicate that representation learning by KCL can distinguish the subsets $\mathbb{M}_1, \cdots, \mathbb{M}_K$ in the RKHS (see Figure 1 for illustration). The function $R(\lambda)$ includes the term that represents the hardness of the pretraining task in the space of augmented data $\mathbb{X}$: if the overlaps between two different subsets $\mathbb{M}_i$ and $\mathbb{M}_j$ expand, then $R(\lambda)$ increases (the precise definition is given in Appendix B.2).

9

Figure 1: An illustration of Theorem 1. The clusters $\mathbb{M}_1, \cdots, \mathbb{M}_K$ in the data space are mapped into the RKHS, where each cluster in the RKHS shrinks or expands (via $\mathfrak{a}(f)$) while maintaining the distance to other clusters (via $\mathfrak{c}(f)$).

A key point of Theorem 1 is that $\delta$ and $\lambda$ can determine how learned representations distribute in the RKHS. If $\delta > 0$, then representations of augmented data in each $\mathbb{M}_i$ tend to align as controlling the trade-off between $\mathfrak{a}(f)$ and $\mathfrak{c}(f)$. For $\delta \leq 0$, not just the means $\mu_1(f), \cdots, \mu_K(f)$ but also the representations tend to scatter. We remark that $\delta$ depends on the fixed weight $\lambda$ due to the condition **(B)** in Assumption 2. Intuitively, larger $\lambda$ makes $\delta$ smaller and vice versa.

Here we should remark that under several assumptions, the equality holds in (3), as shown below:

**Corollary 1.** *Suppose that Assumption 1 and 2 hold. Take $\delta \in \mathbb{R}$, $K \in \mathbb{N}$, and $\mathbb{M}_1, \cdots, \mathbb{M}_K$ such that the conditions **(A)** and **(B)** in Assumption 2 are satisfied. Suppose that $\mathbb{M}_1, \cdots, \mathbb{M}_K$ are disjoint, and for every pair $(i, j) \in [K] \times [K]$ such that $i \neq j$, every $(x, x') \in \mathbb{M}_i \times \mathbb{M}_j$ satisfies $w(x, x') = 0$. Suppose that for every $i \in [K]$, it holds that $\mathrm{sim}(x, x'; \lambda) = \delta$ for any $x, x' \in \mathbb{M}_i$. Then, for every $f \in \mathcal{F}$, the equality holds in (3), i.e.,*

$$L_{\mathrm{KCL}}(f; \lambda) = \frac{\delta}{2} \cdot \mathfrak{a}(f) + \lambda \cdot \mathfrak{c}(f) - R(\lambda).$$

The proof of Corollary 1 can be found in Appendix B.3. The above corollary means that, under these assumptions, the minimization of the kernel contrastive loss is equivalent to that of the objective function $(\delta/2) \cdot \mathfrak{a}(f) + \lambda \cdot \mathfrak{c}(f)$. Note that Example 1 satisfies all the assumptions enumerated in the statement. In summary, Theorem 1 and Corollary 1 imply that contrastive learning by the KCL framework is characterized as representation learning with the similarity structure of augmented data space $\mathbb{X}$.

### 5.1.1 Comparison to Related Work

The quantity $\mathfrak{a}(f)$ is closely related to the property called *alignment* (Wang and Isola, 2020) since the representations of similar data are learned to be close in the RKHS. Also, the quantity $\mathfrak{c}(f)$ has some connection to *divergence property* (Huang et al., 2023) since it measures how far apart the means $\mu_i(f)$ and $\mu_j(f)$ are. Although the relations between these properties and contrastive learning have been pointed out by Huang et al. (2023); Wang and Isola (2020), we emphasize that our result gives a new characterization of the learned clusters. Furthermore, this theorem also implies that the trade-off between $\mathfrak{a}(f)$ and $\mathfrak{c}(f)$ is determined with the threshold $\delta$ and the hyperparameter $\lambda$.

Therefore, Theorem 1 provides deeper insights into understanding the mechanism of contrastive learning.

## 5.2 A New Upper Bound of the Classification Error

Next, we show how minimization of the kernel contrastive loss guarantees good performance in the downstream classification task, according to our formulation of similarity. To this end, we prove that the properties of contrastive learning shown in Theorem 1 yield the linearly well-separable representations in the RKHS. First, we quantify the linear separability as follows: following HaoChen et al. (2021); Saunshi et al. (2022), under Assumption 2, for a model $f \in \mathcal{F}$, a linear weight $W : \mathcal{H}_k \to \mathbb{R}^K$, and a bias $\beta \in \mathbb{R}^K$, we define the downstream classification error as,

$$L_{\mathrm{Err}}(f, W, \beta; y) := P_{\mathbb{X}}\left(g_{f,W,\beta}(x) \neq y(x)\right),$$

where $g_{f,W,\beta}(x) := \arg\max_{i \in [K]}\{\langle W_i, h(f(x))\rangle_{\mathcal{H}_k} + \beta_i\}$ for $W_i \in \mathcal{H}_k$ and $\beta_i \in \mathbb{R}$. Note that we let arg max and arg min break tie arbitrary as well as HaoChen et al. (2021). Note that in our definition, after augmented data $x \in \mathbb{X}$ is encoded to $f(x) \in \mathbb{S}^{d-1}$, $f(x)$ is further mapped to $h(f(x))$ in the RKHS, and then linear classification is performed using $W$ and $\beta$.

To derive a generalization guarantee for KCL, we focus on the 1-Nearest Neighbor (NN) classifier in the RKHS $\mathcal{H}_k$ as a technical tool, which is a generalization of the 1-NN classifiers utilized in Huang et al. (2023); Robinson et al. (2021b).

**Definition 1** (1-NN classifier in $\mathcal{H}_k$). *Suppose $\mathbb{C}_1, \cdots, \mathbb{C}_K$ are subsets of $\mathbb{X}$. For a model $f : \mathbb{X} \to \mathbb{R}^d$, the 1-NN classifier $g_{1\text{-NN}} : \mathbb{X} \to [K]$ associated with the RKHS $\mathcal{H}_k$ is defined as*

$$g_{1\text{-NN}}(x) := \arg\min_{i \in [K]} \|h(f(x)) - \mu_{\mathbb{C}_i}(f)\|_{\mathcal{H}_k}.$$

Huang et al. (2023) show that the 1-NN classifier they consider can be regarded as a mean classifier (Arora et al., 2019; Wang et al., 2022a) (see Appendix E in Huang et al. (2023)). This fact can also apply to our setup: indeed, under Assumption 1, $g_{1\text{-NN}}$ is equal to

$$g_{f,W_\mu,\beta_\mu}(x) = \arg\max_{i \in [K]} \left\{\langle W_{\mu,i}, h(f(x))\rangle_{\mathcal{H}_k} - \beta_{\mu,i}\right\},$$

where $W_\mu : \mathcal{H}_k \to \mathbb{R}^K$ is defined as $W_\mu(\phi)_i := \langle W_{\mu,i}, \phi\rangle_{\mathcal{H}_k} = \langle \mu_{\mathbb{C}_i}(f), \phi\rangle_{\mathcal{H}_k}$ for each coordinate $i \in [K]$, and $\beta_{\mu,i} := (\|\mu_{\mathbb{C}_i}(f)\|_{\mathcal{H}_k}^2 + \psi(1))/2$ for $i \in [K]$.

Before presenting the result, we need the following notion:

**Definition 2** (Meaningful encoder). *An encoder $f \in \mathcal{F}$ is said to be meaningful if $\min_{i \neq j} \|\mu_i(f) - \mu_j(f)\|_{\mathcal{H}_k}^2 > 0$ holds.*

Note that a meaningful encoder $f \in \mathcal{F}$ avoids the *complete collapse* of feature vectors (Hua et al., 2021; Jing et al., 2022), where many works on self-supervised representation learning (Chen et al., 2020a; Chen and He, 2021; Grill et al., 2020; HaoChen et al., 2021; Li et al., 2021) introduce various architectures and algorithms to prevent it. Now, the theoretical guarantee is presented:

**Theorem 2.** *Suppose that Assumption 1 and 2 hold. Take $\delta \in \mathbb{R}$, $K \in \mathbb{N}$, $\mathbb{M}_1, \cdots, \mathbb{M}_K$, and $y$ such that the conditions* **(A)**, **(B)**, *and* **(C)** *in Assumption 2 are satisfied. Then, for each meaningful encoder $f \in \mathcal{F}$, we have*

$$L_{\text{Err}}(f, W_\mu, \beta_\mu; y) \leq \frac{8(K-1)}{\Delta_{\min}(f) \cdot \min_{i \in [K]} P_{\mathbb{X}}(\mathbb{M}_i)} \mathfrak{a}(f)$$

*where $\Delta_{\min}(f) = \min_{i \neq j} \|\mu_i(f) - \mu_j(f)\|_{\mathcal{H}_k}^2$.*

The proof of Theorem 2 can be found in Appendix C. The upper bound in Theorem 2 becomes smaller if the representations of any two points $x, x'$ belonging to $\mathbb{M}_i$ are closer for each $i \in [K]$ and the closest centers $\mu_i(f)$ and $\mu_j(f)$ of different subsets $\mathbb{M}_i$ and $\mathbb{M}_j$ become distant from each other. Since $\|\mu_i(f) - \mu_j(f)\|_{\mathcal{H}_k}^2 = -2\langle \mu_i(f), \mu_j(f) \rangle_{\mathcal{H}_k} + \|\mu_i(f)\|_{\mathcal{H}_k}^2 + \|\mu_j(f)\|_{\mathcal{H}_k}^2$, Theorem 1 and 2 indicate that during the optimization for the kernel contrastive loss, the quantities $\mathfrak{a}(f)$ and $\mathfrak{c}(f)$ can contribute to making the learned representations linearly well-separable. Thus, our theory is consistent with the empirical success of contrastive learning shown by a line of research (Chen et al., 2020a,b; Dwibedi et al., 2021; HaoChen et al., 2021).

### 5.2.1 Comparison to Related Work

We discuss Theorem 2. 1) Several works (Huang et al., 2023; Robinson et al., 2021b) also show that the classification loss or error is upper bounded by the quantity related to the alignment of feature vectors within each cluster. However, their results do not address the following conjecture: *Does the distance between the centers of each cluster consisting of feature vectors affect the linear separability?* Theorem 2 indicates that the answer is yes via the quantity $\Delta_{\min}(f)$. Note that Theorem 3.2 of Zhao et al. (2023) implies a similar answer, but their result is for the squared loss and requires several strong assumptions on the encoder functions. Meanwhile, our Theorem 2 for the classification error requires the meaningfulness of encoder functions (Definition 2), which is more practical than those of Zhao et al. (2023). 2) Furthermore, our result is different from previous works (Huang et al., 2023; Robinson et al., 2021b; Zhao et al., 2023) in the problem setup. Indeed, Robinson et al. (2021b) follow the setup of Arora et al. (2019), and Huang et al. (2023); Zhao et al. (2023) formulate their setup by imposing the $(\sigma, \delta)$-*augmentation* property to given latent class subsets. Meanwhile, our formulation is mainly based on the statistical similarity (2). Furthermore, we note that our Theorem 2 can be extended to the case that $K \in \mathbb{N}$, $\mathbb{M}_1, \cdots, \mathbb{M}_K \subset \mathbb{X}$, and $y$ are taken to satisfy the conditions **(A)** and **(C)** in Assumption 2 (see Theorem 5 in Appendix C), implying that our result can apply to other problem setups of contrastive learning. Due to space limitations, we present more detailed explanations in Appendix F.6 and F.7.

### 5.3 A Generalization Error Bound for KCL

Since in practice we minimize the empirical kernel contrastive loss, we derive a generalization error bound for KCL. The empirical loss is defined as follows: denote $dP_+(x, x') = w(x, x') d\nu_{\mathbb{X}}^{\otimes 2}(x, x')$. Let $(X_1, X_1'), \cdots, (X_n, X_n')$ be pairs of random variables drawn independently from $P_+$, where $X_i$ and $X_j'$ are assumed to be independent for each pair of distinct indices $i, j \in [n]$. Following the standard setup that a pair of two augmented samples is obtained by randomly transforming the

same raw sample, which is considered in many empirical works (Chen et al., 2020a,b; Dwibedi et al., 2021), for each $i \in [n]$, we consider the case that $X_i, X_i'$ are not necessarily independent. The empirical kernel contrastive loss is defined as,

$$\widehat{L}_{\mathrm{KCL}}(f; \lambda) = -\frac{1}{n} \sum_{i=1}^{n} k(f(X_i), f(X_i')) + \frac{\lambda}{n(n-1)} \sum_{i \neq j} k(f(X_i), f(X_j')). \tag{4}$$

In the statement below, denote $\mathcal{Q} = \{f(\cdot)^\top f(\cdot) : \mathbb{X} \times \mathbb{X} \to \mathbb{R} \mid f \in \mathcal{F}\}$. Define the Rademacher complexity (Mohri et al., 2018) as, $\mathfrak{R}_n^+(\mathcal{Q}) := \mathbb{E}_{P_+, \sigma_{1:n}}[\sup_{q \in \mathcal{Q}} n^{-1} \sum_{i=1}^n \sigma_i q(X_i, X_i')]$, where $\sigma_1, \cdots, \sigma_n$ are independent random variables taking $\pm 1$ with probability one half for each. We also define the Rademacher compexltiy $\mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*)$ with the optimal choice $s^*$ from the symmetric group $S_n$ of degree $n$:

$$\mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*) := \max_{s \in S_n} \mathbb{E}_{\substack{X, X' \\ \sigma_{1:(n/2)}}} \left[ \sup_{q \in \mathcal{Q}} \frac{2}{n} \sum_{i=1}^{n/2} \sigma_i q(X_{s(2i-1)}, X'_{s(2i)}) \right].$$

Note that $\mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*)$ is related to the *average of "sums-of-i.i.d." blocks* technique for $U$-statistics explained in Clémençon et al. (2008); for more detail, see also Remark 3 in Appendix D.4. The generalization error bound for KCL is presented below:

**Theorem 3.** *Suppose that Assumption 1 holds, and $n$ is even. Furthermore, suppose that the minimizer $\widehat{f} \in \mathcal{F}$ of $\widehat{L}_{\mathrm{KCL}}(f; \lambda)$ exists. Then, with probability at least $1 - 2\varepsilon$ where $\varepsilon > 0$, we have*

$$L_{\mathrm{KCL}}(\widehat{f}; \lambda) \leq L_{\mathrm{KCL}}(f; \lambda) + 2 \cdot \mathrm{Gen}(n, \lambda, \varepsilon),$$

*where $\mathfrak{R}_n(\mathcal{Q}) := \mathfrak{R}_n^+(\mathcal{Q}) + \lambda \mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*)$, and*

$$\mathrm{Gen}(n, \lambda, \varepsilon) = O\left( \mathfrak{R}_n(\mathcal{Q}) + (1 + \lambda)\sqrt{\frac{\log(2/\varepsilon)}{n}} \right).$$

**Remark 2.** *In Appendix D.2, we show that under some conditions, $\mathrm{Gen}(n, \lambda, \varepsilon) \downarrow 0$ holds as $n \to \infty$.*

The proof of Theorem 3 can be found in Appendix D.1. Since $X_i, X_i'$ are not necessarily independent for each $i \in [n]$, the standard techniques (e.g., Theorem 3.3 in Mohri et al. (2018)) are not applicable in the proof. We instead utilize the results by Zhang et al. (2019) to overcome this difficulty, which is different from the previous bounds for contrastive learning (Arora et al., 2019; Ash et al., 2022; HaoChen et al., 2021; Lei et al., 2023; Nozawa et al., 2020; Wang et al., 2022b; Zhang et al., 2022; Zou and Liu, 2023) (see Appendix F.5 for more detail). Here, if $\mathfrak{R}_n(\mathcal{Q}) \downarrow 0$ as $n \to \infty$, then by using our result, we can prove the consistency of the empirical contrastive loss to the population one for each $f \in \mathcal{F}$.

## 5.4 Application of the Theoretical Results: A New Surrogate Bound

Recent works (Arora et al., 2019; Ash et al., 2022; Awasthi et al., 2022; Bao et al., 2022; Dufumier et al., 2022; HaoChen et al., 2021; Nozawa et al., 2020; Nozawa and Sato, 2021; Saunshi et al., 2022;

Tosh et al., 2021b; Wang et al., 2022a; Zou and Liu, 2023) show that some contrastive learning objectives can be regarded as surrogate losses of the supervised loss or error in downstream tasks. Here, Arora et al. (2019) show that, a contrastive loss $L_{\mathrm{CL}}$ *surrogates* a supervised loss or error $L_{\mathrm{Sup}}$: for every $f \in \mathcal{F}$, $L_{\mathrm{sup}}(W \circ f) \lesssim L_{\mathrm{CL}}(f) + \alpha$ holds for some $\alpha \in \mathbb{R}$ and matrix $W$. This type of inequality is also called *surrogate bound* (Bao et al., 2022). Arora et al. (2019) show that the inequality guarantees that $L_{\mathrm{sup}}(W^* \circ \widehat{f}) \lesssim L_{\mathrm{CL}}(f) + \alpha$ holds with high probability, where $\widehat{f}$ is a minimizer of the empirical loss for $L_{\mathrm{CL}}$, $W^*$ is the optimal weight, and $\alpha$ is some term. Motivated by these works, we show a surrogate bound for KCL.

**Theorem 4.** *Suppose that Assumption 1 and 2 hold, $n$ is even, and there exists a minimizer $\widehat{f}$ of $\widehat{L}_{\mathrm{KCL}}(f; \lambda)$ such that $\widehat{f}$ is meaningful. Take $\delta \in \mathbb{R}$, $K \in \mathbb{N}$, $\mathbb{M}_1, \cdots, \mathbb{M}_K$, and $y$ such that the conditions* **(A)**, **(B)**, *and* **(C)** *in Assumption 2 are satisfied. Then, for any $f \in \mathcal{F}$ and $\varepsilon > 0$, with probability at least $1 - 2\varepsilon$,*

$$L_{\mathrm{Err}}(\widehat{f}, W^*, \beta^*; y) \lesssim L_{\mathrm{KCL}}(f; \lambda) + (1 - \frac{\delta}{2})\mathfrak{a}(\widehat{f}) - \lambda\mathfrak{c}(\widehat{f}) + R(\lambda) + 2\mathrm{Gen}(n, \lambda, \varepsilon),$$

*where $L_{\mathrm{Err}}(\widehat{f}, W^*, \beta^*; y) = \inf_{W,\beta} L_{\mathrm{Err}}(\widehat{f}, W, \beta; y)$, and $\lesssim$ omits the coefficient $8(K-1)/(\Delta_{\min}(\widehat{f}) \cdot \min_{i \in [K]} P_{\mathbb{X}}(\mathbb{M}_i))$.*

The proof of Theorem 4 can be found in Appendix E. This theorem indicates that minimization of the kernel contrastive loss in $\mathcal{F}$ can reduce the infimum of the classification error with high probability. Note that since larger $\lambda$ can make $\delta$ smaller due to the relation in condition **(B)** of Assumption 2, larger $\lambda$ may result in enlarging $(1 - \delta/2)\mathfrak{a}(\widehat{f}) - \lambda\mathfrak{c}(\widehat{f})$ and loosening the upper bound if $\mathfrak{a}(\widehat{f}) > 0$ and $\mathfrak{c}(\widehat{f}) < 0$. We empirically find that the KCL framework with larger $\lambda$ degrades its performance in the downstream classification task; see Appendix H.4.

### 5.4.1 Comparison to Related Work

Several works also establish the surrogate bounds for some contrastive learning objectives (Arora et al., 2019; Ash et al., 2022; Awasthi et al., 2022; Bao et al., 2022; Dufumier et al., 2022; HaoChen et al., 2021; Nozawa et al., 2020; Nozawa and Sato, 2021; Saunshi et al., 2022; Tosh et al., 2021b; Wang et al., 2022a; Zou and Liu, 2023). The main differences between the previous works and Theorem 4 are summarized in three points: 1) Theorem 4 indicates that the kernel contrastive loss is a surrogate loss of the classification error, while the previous works deal with other contrastive learning objectives. 2) Recent works (Bao et al., 2022; Wang et al., 2022a) prove that the InfoNCE loss is a surrogate loss of the cross-entropy loss. However, since the theory of classification calibration losses (see e.g., Zhang (2004)) indicates that the relation between the classification loss and the cross-entropy loss is complicated under the multi-class setting, the relation between the InfoNCE loss and the classification error is non-trivial from the previous results. On the other hand, combining Theorem 4 and (1), we can show that the InfoNCE loss is also a surrogate loss of the classification error. Note that Theorem 4 can apply to other contrastive learning objectives. 3) The bound in Theorem 4 is established by introducing the formulation presented in Section 4. Especially our bound includes the geometric quantity $\delta$ and hyperparameter $\lambda$.

# 6 Conclusion and Discussion

In this paper, we studied the characterization of the structure of the representations learned by contrastive learning. By employing Kernel Contrastive Learning (KCL) as a unified framework, we showed that the formulation based on statistical similarity characterizes the clusters of learned representations and guarantees that the kernel contrastive loss minimization can yield good performance in the downstream classification task. As a limitation of this paper, we point out that in practice, it is challenging to compute the true $\text{sim}(x, x'; \lambda)$ and $\delta$ for datasets. However, we believe that our theory promotes future theoretical and empirical research to investigate the practical success of contrastive learning via the sets of augmented data defined by $\text{sim}(x, x'; \lambda)$ and $\delta$. Note that as recent works (Tsai et al., 2021, 2020) tackle the estimation of the point-wise dependency by using neural networks, the estimation problem is an important future work. As a future work, it is worth studying how the selection of kernels affects the quality of representations via our theory. The investigation of transfer learning perspectives of KCL is also an interesting future work, as recent works (HaoChen et al., 2022; Shen et al., 2022; Zhao et al., 2023) also address the problem for some contrastive learning frameworks.

**Ethical Statement**

Since this paper mainly studies theoretical analysis of contrastive learning, it will not be thought that there is a direct negative social impact. However, revealing detailed properties of contrastive learning could promote an opportunity to misuse the knowledge. We point out that such wrong usage is not straightforward with the proposed method, as the application is not discussed much in the paper.

# References

Arendt, W., Batty, C. J., Hieber, M., and Neubrander, F. (2011). *Vector-valued Laplace transforms and Cauchy problems*. Birkhäuser Basel. 24

Aronszajn, N. (1950). Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68(3):337–404. 3, 5, 51

Arora, S., Khandeparkar, H., Khodak, M., Plevrakis, O., and Saunshi, N. (2019). A theoretical analysis of contrastive unsupervised representation learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5628–5637. PMLR. 2, 4, 8, 9, 11, 12, 13, 14, 40, 47, 48

Ash, J., Goel, S., Krishnamurthy, A., and Misra, D. (2022). Investigating the role of negatives in contrastive representation learning. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 7187–7209. PMLR. 2, 13, 14, 40, 47

Awasthi, P., Dikkala, N., and Kamath, P. (2022). Do more negative samples necessarily hurt in contrastive learning? In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 1101–1116. PMLR. 2, 13, 14

Bao, H., Nagano, Y., and Nozawa, K. (2022). On the surrogate gap between contrastive and supervised losses. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 1585–1606. PMLR. 2, 4, 6, 13, 14

Berlinet, A. and Thomas-Agnan, C. (2004). *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media. 3, 5, 23, 51

Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., and Joulin, A. (2020). Unsupervised learning of visual features by contrasting cluster assignments. In *Advances in Neural Information Processing Systems*, volume 33, pages 9912–9924. Curran Associates, Inc. 1, 2, 3

Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. (2020a). A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607. PMLR. 1, 2, 3, 4, 5, 6, 11, 12, 13, 44, 55, 57, 58

Chen, T., Luo, C., and Li, L. (2021). Intriguing properties of contrastive losses. In *Advances in Neural Information Processing Systems*, volume 34, pages 11834–11845. Curran Associates, Inc. 3, 4, 6, 7, 44, 60

Chen, X., Fan, H., Girshick, R., and He, K. (2020b). Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297v1*. 1, 2, 3, 4, 12, 13, 56

Chen, X. and He, K. (2021). Exploring simple siamese representation learning. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15745–15753. 3, 11, 55, 56, 57, 60

Chuang, C.-Y., Robinson, J., Lin, Y.-C., Torralba, A., and Jegelka, S. (2020). Debiased contrastive learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 8765–8775. Curran Associates, Inc. 4

Clémençon, S., Lugosi, G., and Vayatis, N. (2008). Ranking and Empirical Minimization of U-statistics. *The Annals of Statistics*, 36(2):844 – 874. 13, 39, 40, 48

Coates, A., Ng, A., and Lee, H. (2011). An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 215–223. PMLR. 57

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. Ieee. 57

Dubois, Y., Ermon, S., Hashimoto, T., and Liang, P. (2022). Improving self-supervised learning by characterizing idealized representations. In *Advances in Neural Information Processing Systems*. 4

Dufumier, B., Barbano, C. A., Louiset, R., Duchesnay, E., and Gori, P. (2022). Rethinking positive sampling for contrastive learning with kernel. *arXiv preprint arXiv:2206.01646v1*. 2, 4, 8, 9, 13, 14

Dwibedi, D., Aytar, Y., Tompson, J., Sermanet, P., and Zisserman, A. (2021). With a little help from my friends: Nearest-neighbor contrastive learning of visual representations. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9568–9577. 1, 2, 3, 4, 5, 6, 12, 13, 55

Gao, T., Yao, X., and Chen, D. (2021). SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910. Association for Computational Linguistics. 3

Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. (2017). Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677v2*. 56

Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. (2005). Measuring statistical dependence with hilbert-schmidt norms. In *Algorithmic Learning Theory*, pages 63–77. Springer Berlin Heidelberg. 4

Grill, J.-B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., Piot, B., Kavukcuoglu, K., Munos, R., and Valko, M. (2020). Bootstrap your own latent - a new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 21271–21284. Curran Associates, Inc. 3, 11

Guedj, B. (2019). A primer on pac-bayesian learning. *arXiv preprint arXiv:1901.05353v3*. 47

HaoChen, J. Z. and Ma, T. (2023). A theoretical study of inductive biases in contrastive learning. The Eleventh International Conference on Learning Representations. https://openreview.net/forum?id=AuEgNlEAmed. 2, 4, 8, 9, 45, 46

HaoChen, J. Z., Wei, C., Gaidon, A., and Ma, T. (2021). Provable guarantees for self-supervised deep learning with spectral contrastive loss. In *Advances in Neural Information Processing Systems*, volume 34, pages 5000–5011. Curran Associates, Inc. 1, 2, 3, 4, 5, 7, 8, 9, 11, 12, 13, 14, 45, 47, 54, 56, 57, 58

HaoChen, J. Z., Wei, C., Kumar, A., and Ma, T. (2022). Beyond separability: Analyzing the linear transferability of contrastive representations to related subpopulations. In *Advances in Neural Information Processing Systems*. 4, 15

Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825):357–362. 58

He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778. 5, 56, 57

Hua, T., Wang, W., Xue, Z., Ren, S., Wang, Y., and Zhao, H. (2021). On feature decorrelation in self-supervised learning. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9578–9588. 11

Huang, W., Yi, M., Zhao, X., and Jiang, Z. (2023). Towards the generalization of contrastive self-supervised learning. The Eleventh International Conference on Learning Representations. https://openreview.net/forum?id=XDJwuEYHhme. 2, 3, 4, 7, 9, 10, 11, 12, 48, 49

Hunter, J. D. (2007). Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95. 58

Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456. PMLR. 55

Jiang, D., Li, W., Cao, M., Zou, W., and Li, X. (2021). Speech SimCLR: Combining contrastive and reconstruction objective for self-supervised speech representation learning. In *Proc. Interspeech 2021*, pages 1544–1548. 3

Jing, L., Vincent, P., LeCun, Y., and Tian, Y. (2022). Understanding dimensional collapse in contrastive self-supervised learning. International Conference on Learning Representations. https://openreview.net/forum?id=YevsQ05DEN7. 11

Johnson, D. D., Hanchi, A. E., and Maddison, C. J. (2023). Contrastive learning can find an optimal basis for approximately view-invariant functions. The Eleventh International Conference on Learning Representations. https://openreview.net/forum?id=AjCOKBjiMu. 4, 7, 8

Kiani, B. T., Balestriero, R., Chen, Y., Lloyd, S., and LeCun, Y. (2022). Joint embedding self-supervised learning in the kernel regime. *arXiv preprint arXiv:2209.14884v1*. 4

Krizhevsky, A. (2009). Learning multiple layers of features from tiny images. Technical report. 57

Lei, Y., Yang, T., Ying, Y., and Zhou, D.-X. (2023). Generalization analysis for contrastive representation learning. *arXiv preprint arXiv:2302.12383v2*. 13, 47

Li, Y., Pogodin, R., Sutherland, D. J., and Gretton, A. (2021). Self-supervised learning with kernel dependence maximization. In *Advances in Neural Information Processing Systems*, volume 34, pages 15543–15556. Curran Associates, Inc. 2, 3, 4, 6, 11, 43, 44

Loshchilov, I. and Hutter, F. (2017). SGDR: Stochastic gradient descent with warm restarts. International Conference on Learning Representations. https://openreview.net/forum?id=Skq89Scxx. 56

McDiarmid, C. (1989). *On the method of bounded differences*, page 148–188. London Mathematical Society Lecture Note Series. Cambridge University Press. 37

Mohri, M., Rostamizadeh, A., and Talwalkar, A. (2018). *Foundations of machine learning*. MIT press. 13, 33, 38, 47

Muandet, K., Fukumizu, K., Sriperumbudur, B., and Schölkopf, B. (2017). Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141. 6, 23

Ng, A. Y., Jordan, M. I., and Weiss, Y. (2002). On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems*, pages 849–856. 4

Nozawa, K., Germain, P., and Guedj, B. (2020). Pac-bayesian contrastive unsupervised representation learning. In *Conference on Uncertainty in Artificial Intelligence*, pages 21–30. PMLR. 13, 14, 47

Nozawa, K. and Sato, I. (2021). Understanding negative samples in instance discriminative self-supervised representation learning. In *Advances in Neural Information Processing Systems*, volume 34, pages 5784–5797. Curran Associates, Inc. 13, 14

Parulekar, A., Collins, L., Shanmugam, K., Mokhtari, A., and Shakkottai, S. (2023). Infonce loss provably learns cluster-preserving representations. *arXiv preprint arXiv:2302.07920v1*. 9

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc. 55, 58

Poole, B., Ozair, S., van den Oord, A., Alemi, A., and Tucker, G. (2019). On variational bounds of mutual information. In *International Conference on Machine Learning*, pages 5171–5180. PMLR. 4

Robinson, J., Sun, L., Yu, K., Batmanghelich, K., Jegelka, S., and Sra, S. (2021a). Can contrastive learning avoid shortcut solutions? In *Advances in Neural Information Processing Systems*, volume 34, pages 4974–4986. Curran Associates, Inc. 3, 4

Robinson, J. D., Chuang, C.-Y., Sra, S., and Jegelka, S. (2021b). Contrastive learning with hard negative samples. International Conference on Learning Representations. https://openreview.net/forum?id=CR1XOQ0UTh-. 3, 4, 11, 12, 48

Saunshi, N., Ash, J., Goel, S., Misra, D., Zhang, C., Arora, S., Kakade, S., and Krishnamurthy, A. (2022). Understanding contrastive learning requires incorporating inductive biases. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 19250–19286. PMLR. 4, 8, 11, 13, 14

Shalev-Shwartz, S. and Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge University Press. 33, 39

Shen, K., Jones, R. M., Kumar, A., Xie, S. M., Haochen, J. Z., Ma, T., and Liang, P. (2022). Connect, not collapse: Explaining contrastive learning for unsupervised domain adaptation. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 19847–19878. PMLR. 4, 15

Shi, J. and Malik, J. (2000). Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):888–905. 49, 50

Singh, A. (2021). Clda: Contrastive learning for semi-supervised domain adaptation. In *Advances in Neural Information Processing Systems*, volume 34, pages 5089–5101. Curran Associates, Inc. 3

Smola, A., Gretton, A., Song, L., and Schölkopf, B. (2007). A hilbert space embedding for distributions. In *Algorithmic Learning Theory*, pages 13–31. Springer Berlin Heidelberg. 23, 43

Steinwart, I. and Christmann, A. (2008). *Support vector machines*. Springer Science & Business Media. 3, 5, 51

Susmelj, I., Helle, M., Wirth, P., Prescott, J., and Ebner et al., M. (2020). Lightly. https://github.com/lightly-ai/lightly. 57, 58

Terada, Y. and Yamamoto, M. (2019). Kernel normalized cut: A theoretical revisit. In *International Conference on Machine Learning*, pages 6206–6214. PMLR. 8, 49, 50

Tian, Y. (2022). Understanding deep contrastive learning via coordinate-wise optimization. In *Advances in Neural Information Processing Systems*. 4, 8, 54

Tian, Y., Krishnan, D., and Isola, P. (2020). Contrastive multiview coding. In *Computer Vision – ECCV 2020*, pages 776–794. Springer International Publishing. 57

TorchVision maintainers and contributors (2016). Torchvision: Pytorch's computer vision library. https://github.com/pytorch/vision. 58

Tosh, C., Krishnamurthy, A., and Hsu, D. (2021a). Contrastive estimation reveals topic posterior information to linear models. *Journal of Machine Learning Research*, 22(281):1–31. 8

Tosh, C., Krishnamurthy, A., and Hsu, D. (2021b). Contrastive learning, multi-view redundancy, and linear models. In *Proceedings of the 32nd International Conference on Algorithmic Learning Theory*, volume 132 of *Proceedings of Machine Learning Research*, pages 1179–1206. PMLR. 4, 7, 14

Trillos, N. G., Slepčev, D., Von Brecht, J., Laurent, T., and Bresson, X. (2016). Consistency of cheeger and ratio graph cuts. *The Journal of Machine Learning Research*, 17(1):6268–6313. 50

Tsai, Y.-H. H., Li, T., Ma, M. Q., Zhao, H., Zhang, K., Morency, L.-P., and Salakhutdinov, R. (2022). Conditional contrastive learning with kernel. International Conference on Learning Representations. https://openreview.net/forum?id=AAJLBoGt0XM. 4

Tsai, Y.-H. H., Ma, M. Q., Yang, M., Zhao, H., Morency, L.-P., and Salakhutdinov, R. (2021). Self-supervised representation learning with relative predictive coding. International Conference on Learning Representations. https://openreview.net/forum?id=068E_JSq9O. 15

Tsai, Y.-H. H., Zhao, H., Yamada, M., Morency, L.-P., and Salakhutdinov, R. R. (2020). Neural methods for point-wise dependency estimation. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 62–72. Curran Associates, Inc. 7, 15

Tschannen, M., Djolonga, J., Rubenstein, P. K., Gelly, S., and Lucic, M. (2020). On mutual information maximization for representation learning. International Conference on Learning Representations. https://openreview.net/forum?id=rkxoh24FPH. 4

Tu, Z., Zhang, J., and Tao, D. (2019). Theoretical analysis of adversarial learning: A minimax approach. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc. 35, 36

van den Oord, A., Li, Y., and Vinyals, O. (2018). Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748v2*. 3, 6, 44

Van Gansbeke, W., Vandenhende, S., Georgoulis, S., Proesmans, M., and Van Gool, L. (2020). Scan: Learning to classify images without labels. In *Proceedings of the European Conference on Computer Vision*. 3

von Kügelgen, J., Sharma, Y., Gresele, L., Brendel, W., Schölkopf, B., Besserve, M., and Locatello, F. (2021). Self-supervised learning with data augmentations provably isolates content from style. In *Advances in Neural Information Processing Systems*, volume 34, pages 16451–16467. Curran Associates, Inc. 4, 9

Von Luxburg, U. (2007). A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416. 49, 50, 54

Wainwright, M. J. (2019). *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press. 34, 35, 36, 39

Wang, F. and Liu, H. (2021). Understanding the behaviour of contrastive loss. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2495–2504. 7

Wang, T. and Isola, P. (2020). Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR. 4, 5, 6, 7, 10, 44

Wang, Y., Zhang, Q., Wang, Y., Yang, J., and Lin, Z. (2022a). Chaos is a ladder: A new theoretical understanding of contrastive learning via augmentation overlap. International Conference on Learning Representations. https://openreview.net/forum?id=ECvgmYVyeUz. 2, 4, 5, 8, 9, 11, 14

Wang, Z., Luo, Y., Li, Y., Zhu, J., and Schölkopf, B. (2022b). Spectral representation learning for conditional moment models. *arXiv preprint arXiv:2210.16525v2*. 7, 8, 13, 47, 48

Waskom, M. L. (2021). seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60):3021. 58

Wen, Z. and Li, Y. (2021). Toward understanding the feature learning process of self-supervised contrastive learning. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 11112–11122. PMLR. 4

Yeh, C.-H., Hong, C.-Y., Hsu, Y.-C., Liu, T.-L., Chen, Y., and LeCun, Y. (2022). Decoupled contrastive learning. In *Computer Vision – ECCV 2022*, pages 668–684. Springer Nature Switzerland. 7, 44

Zhang, G., Lu, Y., Sun, S., Guo, H., and Yu, Y. (2022). $f$-mutual information contrastive learning. https://openreview.net/forum?id=3kTt_W1_tgw. 4, 13, 47

Zhang, R. R., Liu, X., Wang, Y., and Wang, L. (2019). Mcdiarmid-type inequalities for graph-dependent variables and stability bounds. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc. 13, 36, 37, 38, 47, 48

Zhang, T. (2004). Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5:1225–1251. 14

Zhao, X., Du, T., Wang, Y., Yao, J., and Huang, W. (2023). ArCL: Enhancing contrastive learning with augmentation-robust representations. The Eleventh International Conference on Learning Representations. https://openreview.net/forum?id=n0Pb9T5kmb. 2, 3, 4, 9, 12, 15, 48, 49

Zou, X. and Liu, W. (2023). Generalization bounds for adversarial contrastive learning. *arXiv preprint arXiv:2302.10633v1*. 2, 13, 14, 47

# A    Proof in Section 2.2

First we prove Proposition 1.

*Proof of Proposition 1.* The proof of the claim closely follows the proof of Lemma 3.1 in Muandet et al. (2017) (see also Smola et al. (2007)), which shows that if $\mathbb{E}_P[\sqrt{k(x,x)}] < +\infty$ where $x \sim P$, then $\mathbb{E}_P[k(\cdot, x)] \in \mathcal{H}_k$. For the sake of completeness, we provide the proof of Proposition 1 by modifying the proof of Muandet et al. (2017) slightly.

Let $\mathbb{M}$ be a measurable set in $\mathbb{X}$. Define $\mu_{\mathbb{M}}(f) := \mathbb{E}[h(f(x))|\mathbb{M}]$. Our goal is to show that $\mu_{\mathbb{M}}(f) \in \mathcal{H}_k$ holds. To this end, for $\phi \in \mathcal{H}_k$, we compute

$$
\begin{aligned}
|\mathbb{E}[\phi(f(x))|\mathbb{M}]| &= |\mathbb{E}[\langle \phi, k(\cdot, f(x)) \rangle_{\mathcal{H}_k}|\mathbb{M}]| \\
&\leq \mathbb{E}[|\langle \phi, k(\cdot, f(x)) \rangle_{\mathcal{H}_k}||\mathbb{M}] \\
&\leq \mathbb{E}[\|\phi\|_{\mathcal{H}_k}\|k(\cdot, f(x))\|_{\mathcal{H}_k}|\mathbb{M}] \qquad \text{(Cauchy-Schwarz ineq.)} \\
&= \|\phi\|_{\mathcal{H}_k}\mathbb{E}\left[\sqrt{k(f(x), f(x))}|\mathbb{M}\right].
\end{aligned}
$$

Since $\sup_{z,z' \in \mathbb{S}^{d-1}} k(z, z') < \infty$ holds, we have $\mathbb{E}[\sqrt{k(f(x), f(x))}|\mathbb{M}] < +\infty$. Hence, the map $\phi \mapsto \mathbb{E}[\phi(f(x))|\mathbb{M}]$ is a bounded linear functional on $\mathcal{H}_k$, and thus from Riesz's representation theorem, there exists some $\xi \in \mathcal{H}_k$ such that $\mathbb{E}[\phi(f(x))|\mathbb{M}] = \langle \xi, \phi \rangle_{\mathcal{H}_k}$. However, let $\phi = k(\cdot, z)$, then $\xi(z) = \langle \xi, k(\cdot, z) \rangle_{\mathcal{H}_k} = \mathbb{E}[k(f(x), z)|\mathbb{M}]$. This implies $\xi = \mathbb{E}[k(f(x), \cdot)|\mathbb{M}] \in \mathcal{H}_k$. Since $k$ is symmetric, we have $\mu_{\mathbb{M}}(f) = \mathbb{E}[k(\cdot, f(x))|\mathbb{M}] \in \mathcal{H}_k$.  □

# B    Proofs in Section 5.1

## B.1    Useful Lemmas for the Proof of Theorem 1

Before showing Theorem 1, we give several basic and useful lemmas that are used in the proof of the theorem. Since the definition of $\mu_{\mathbb{M}}(f)$, where $\mathbb{M}$ is a measurable subset of $\mathbb{X}$ and $f \in \mathcal{F}$, is slightly different from the kernel mean embedding of the usual form (Berlinet and Thomas-Agnan, 2004; Muandet et al., 2017) due to the existence of the encoder function $f$, we provide the proof for each lemma for the sake of completeness.

**Lemma 1.** *Let $\{e_j\}$ be an orthonormal basis of $\mathcal{H}_k$, and let $\mathbb{M}$ be a measurable set. Let $f \in \mathcal{F}$. Then, the following identity holds for each $j$:*

$$
\int_{\mathbb{M}} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}) = \langle \mu_{\mathbb{M}}(f), e_j \rangle_{\mathcal{H}_k}.
$$

*Proof.* We calculate,

$$\langle \mu_{\mathbb{M}}(f), e_j \rangle_{\mathcal{H}_k} = \left\langle \int_{\mathbb{M}} h(f(x)) P_{\mathbb{X}}(dx|\mathbb{M}), e_j \right\rangle_{\mathcal{H}_k}$$

$$= \left\langle \int_{\mathbb{M}} \sum_{j'} \langle h(f(x)), e_{j'} \rangle_{\mathcal{H}_k} e_{j'} P_{\mathbb{X}}(dx|\mathbb{M}), e_j \right\rangle_{\mathcal{H}_k}$$

$$= \left\langle \sum_{j'} e_{j'} \int_{\mathbb{M}} \langle h(f(x)), e_{j'} \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}), e_j \right\rangle_{\mathcal{H}_k}$$

$$= \int_{\mathbb{M}} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}),$$

where in the third line, we use the Dominated Convergence Theorem for the Bochner integral (e.g., see Theorem 1.1.8 in Arendt et al. (2011)). Hence, we obtain the claim. $\square$

**Lemma 2.** *Let $\mathbb{M}, \mathbb{M}'$ be measurable subsets of $\mathbb{X}$. Let $f \in \mathcal{F}$. Then, we have*

$$\int_{\mathbb{M}} \int_{\mathbb{M}'} \langle h(f(x)), h(f(x')) \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}) P_{\mathbb{X}}(dx'|\mathbb{M}') = \langle \mu_{\mathbb{M}}(f), \mu_{\mathbb{M}'}(f) \rangle_{\mathcal{H}_k}.$$

*Proof.* Let $\{e_j\}$ be an orthonormal basis of $\mathcal{H}_k$. Then we have,

$$\int_{\mathbb{M}} \int_{\mathbb{M}'} \langle h(f(x)), h(f(x')) \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}) P_{\mathbb{X}}(dx'|\mathbb{M}')$$

$$= \int_{\mathbb{M}} \int_{\mathbb{M}'} \left\langle \sum_j \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} e_j, \sum_j \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} e_j \right\rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}) P_{\mathbb{X}}(dx'|\mathbb{M}')$$

$$= \int_{\mathbb{M}} \int_{\mathbb{M}'} \sum_j \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}) P_{\mathbb{X}}(dx'|\mathbb{M}')$$

$$= \sum_j \int_{\mathbb{M}} \int_{\mathbb{M}'} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}) P_{\mathbb{X}}(dx'|\mathbb{M}') \qquad (5)$$

$$= \sum_j \left( \int_{\mathbb{M}} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx|\mathbb{M}) \right) \left( \int_{\mathbb{M}'} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} P_{\mathbb{X}}(dx'|\mathbb{M}') \right)$$

$$= \sum_j \langle \mu_{\mathbb{M}}(f), e_j \rangle_{\mathcal{H}_k} \langle \mu_{\mathbb{M}'}(f), e_j \rangle_{\mathcal{H}_k} \qquad \text{(Lemma 1)}$$

$$= \left\langle \sum_j \langle \mu_{\mathbb{M}}(f), e_j \rangle_{\mathcal{H}_k} e_j, \sum_j \langle \mu_{\mathbb{M}'}(f), e_j \rangle_{\mathcal{H}_k} e_j \right\rangle_{\mathcal{H}_k}$$

$$= \langle \mu_{\mathbb{M}}(f), \mu_{\mathbb{M}'}(f) \rangle_{\mathcal{H}_k},$$

where in (5) we use the Dominated Convergence Theorem. Hence we obtain the claim. $\square$

24

## B.2 Proof of Theorem 1

The following notation is used in the proof of Theorem 1.

**Definition 3.** *Denote $M_k = \sup_{z,z' \in \mathbb{S}^{d-1}} \|k(\cdot, z) - k(\cdot, z')\|^2_{\mathcal{H}_k}$. We define,*

$$R(\lambda) := \frac{M_k}{2} \sum_{i \neq j} P_+((\mathbb{M}_i \cap \mathbb{M}_j) \times (\mathbb{M}_i \cap \mathbb{M}_j)) + \lambda\psi(1) \sum_{i=1}^{K} P_{\mathbb{X}}(\mathbb{M}_i)(1 - P_{\mathbb{X}}(\mathbb{M}_i)) + (1-\lambda)\psi(1),$$

*where $dP_+(x, x') = w(x, x')d\nu_{\mathbb{X}}^{\otimes 2}(x, x').$*

Note that under Assumption 1, $k(z) := k(z, z) = \psi(z^\top z) = \psi(1)$ is a constant function on $\mathbb{S}^{d-1}$. We are now ready to present the proof of Theorem 1.

*Proof of Theorem 1.* It is convenient to analyze the following form instead of the kernel contrastive loss:

$$\widetilde{L}_{\mathrm{KCL}}(f; \lambda) := \underbrace{\mathbb{E}_{x,x^+} \left[ \|h(f(x)) - h(f(x^+))\|^2_{\mathcal{H}_k} \right]}_{\text{the positive term}} - \lambda \underbrace{\mathbb{E}_{x,x^-} \left[ \|h(f(x)) - h(f(x^-))\|^2_{\mathcal{H}_k} \right]}_{\text{the negative term}}. \quad (6)$$

Note that, $\widetilde{L}_{\mathrm{KCL}}(f; \lambda) = 2(1-\lambda)\psi(1) + 2L_{\mathrm{KCL}}(f; \lambda)$ holds since $f(x) \in \mathbb{S}^{d-1}$ for all $x \in \mathbb{X}$. For the positive term of $\widetilde{L}_{\mathrm{KCL}}(f)$, we can evaluate that,

$$\mathbb{E}_{x,x^+} \left[ \|h(f(x)) - h(f(x^+))\|^2_{\mathcal{H}_k} \right]$$

$$\geq \int_{\bigcup_{i=1}^{K} \mathbb{M}_i \times \mathbb{M}_i} \|h(f(x)) - h(f(x'))\|^2_{\mathcal{H}_k} w(x, x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x') \quad (7)$$

$$\geq \sum_{i=1}^{K} \int_{\mathbb{M}_i \times \mathbb{M}_i} \|h(f(x)) - h(f(x'))\|^2_{\mathcal{H}_k} w(x, x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x')$$

$$- \sum_{j \neq i} \int_{(\mathbb{M}_i \cap \mathbb{M}_j) \times (\mathbb{M}_i \cap \mathbb{M}_j)} \|h(f(x)) - h(f(x'))\|^2_{\mathcal{H}_k} w(x, x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x') \quad (8)$$

$$\geq \sum_{i=1}^{K} \left( \int_{\mathbb{M}_i \times \mathbb{M}_i} \|h(f(x)) - h(f(x'))\|^2_{\mathcal{H}_k} w(x, x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x') - M_k \sum_{j \neq i} P_+((\mathbb{M}_i \cap \mathbb{M}_j) \times (\mathbb{M}_i \cap \mathbb{M}_j)) \right) \quad (9)$$

where in the second inequality we use the fact that

$$Q(\bigcup_{i=1}^{K} \mathbb{M}_i \times \mathbb{M}_i) \geq \sum_{i=1}^{K} Q(\mathbb{M}_i \times \mathbb{M}_i) - \sum_{i \neq j} Q((\mathbb{M}_i \times \mathbb{M}_i) \cap (\mathbb{M}_j \times \mathbb{M}_j)),$$

for any probability measure $Q$ in $\mathbb{X} \times \mathbb{X}$, and in the last inequality we use the definition $M_k = \sup_{z,z' \in \mathbb{S}^{d-1}} \|k(\cdot, z) - k(\cdot, z')\|_{\mathcal{H}_k}^2$. The first term of the above lower bound can be bounded as

$$\sum_{i=1}^K \int_{\mathbb{M}_i \times \mathbb{M}_i} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x, x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$\geq \sum_{i=1}^K \int_{\mathbb{M}_i \times \mathbb{M}_i} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 \cdot (\lambda + \delta) w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x'), \tag{10}$$

where we utilize the definition of $\mathbb{M}_i$ for each $i \in \{1, \cdots, K\}$; recall that due to the condition **(B)** in Assumption 2, for every $x, x' \in \mathbb{M}_i$ we have $\text{sim}(x, x'; \lambda) \geq \delta$.

On the other hand, for the negative term we can compute as follows:

$$- \mathbb{E}_{x, x^-} \left[ \|h(f(x)) - h(f(x^-))\|_{\mathcal{H}_k}^2 \right]$$

$$= - \int_{\bigcup_{i=1}^K \mathbb{M}_i \times \mathbb{M}_i} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$- \int_{\mathbb{X} \times \mathbb{X} \setminus (\bigcup_{i=1}^K \mathbb{M}_i \times \mathbb{M}_i)} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$\geq - \sum_{i=1}^K \int_{\mathbb{M}_i \times \mathbb{M}_i} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$- \int_{\mathbb{X} \times \mathbb{X} \setminus (\bigcup_{i=1}^K \mathbb{M}_i \times \mathbb{M}_i)} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x'), \tag{11}$$

where the last inequality is due to the union bound. For the second term in the right hand side of the inequality above, we have

$$- \int_{\mathbb{X} \times \mathbb{X} \setminus (\bigcup_{i=1}^K \mathbb{M}_i \times \mathbb{M}_i)} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$\geq - \sum_{i \neq j} \int_{\mathbb{M}_i \times \mathbb{M}_j} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$- \int_{\mathbb{X} \times \mathbb{X} \setminus (\bigcup_{i,j=1}^K \mathbb{M}_i \times \mathbb{M}_j)} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x') \tag{12}$$

$$\geq - \sum_{i \neq j} \int_{\mathbb{M}_i \times \mathbb{M}_j} \|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x) w(x') d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x') - M_k P_{\mathbb{X}}^{\otimes 2} \left( \mathbb{X} \times \mathbb{X} \setminus \left( \bigcup_{i,j=1}^K \mathbb{M}_i \times \mathbb{M}_j \right) \right).$$
$$\tag{13}$$

Here, the second term of (13) vanishes since Assumption 2 implies $\mathbb{X} \times \mathbb{X} = \bigcup_{i,j=1}^K \mathbb{M}_i \times \mathbb{M}_j$. The

first term of (13) is further lower bounded as,

$$-\sum_{i\neq j}\int_{\mathbb{M}_i\times\mathbb{M}_j}\|h(f(x))-h(f(x'))\|_{\mathcal{H}_k}^2 w(x)w(x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x')$$

$$\geq 2\sum_{i\neq j}\left\{-\sup_{z\in\mathbb{S}^{d-1}}k(z,z)P_{\mathbb{X}}(\mathbb{M}_i)P_{\mathbb{X}}(\mathbb{M}_j)\right.$$

$$\left.+\int_{\mathbb{M}_i\times\mathbb{M}_j}\langle h(f(x)),h(f(x'))\rangle_{\mathcal{H}_k}w(x)w(x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x')\right\} \tag{14}$$

$$= 2\sum_{i\neq j}\left\{-\sup_{z\in\mathbb{S}^{d-1}}k(z,z)P_{\mathbb{X}}(\mathbb{M}_i)P_{\mathbb{X}}(\mathbb{M}_j)+P_{\mathbb{X}}(\mathbb{M}_i)P_{\mathbb{X}}(\mathbb{M}_j)\langle\mu_i(f),\mu_j(f)\rangle_{\mathcal{H}_k}\right\} \quad \text{(Lemma 2)}$$

$$= -2\sup_{z\in\mathbb{S}^{d-1}}k(z,z)\sum_{i=1}^K P_{\mathbb{X}}(\mathbb{M}_i)\left(1-P_{\mathbb{X}}(\mathbb{M}_i)\right)+2\mathfrak{c}(f).$$

Thus for the negative term we obtain the inequality,

$$-\mathbb{E}_{x,x^-}\left[\|h(f(x))-h(f(x^-))\|_{\mathcal{H}_k}^2\right]$$

$$\geq -\sum_{i=1}^K\int_{\mathbb{M}_i\times\mathbb{M}_i}\|h(f(x))-h(f(x'))\|_{\mathcal{H}_k}^2 w(x)w(x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x')$$

$$+2\mathfrak{c}(f)-2\sup_{z\in\mathbb{S}^{d-1}}k(z,z)\sum_{i=1}^K P_{\mathbb{X}}(\mathbb{M}_i)\left(1-P_{\mathbb{X}}(\mathbb{M}_i)\right). \tag{15}$$

Combining (6),(9),(10), and (15), we have

$$L_{\text{KCL}}(f;\lambda)+(1-\lambda)\psi(1)$$

$$\geq \frac{\delta}{2}\sum_{i=1}^K\mathbb{E}_{x,x^-}\left[\|h(f(x))-h(f(x^-))\|_{\mathcal{H}_k}^2;\mathbb{M}_i\times\mathbb{M}_i\right]+\lambda\mathfrak{c}(f)$$

$$-\frac{M_k}{2}\sum_{i\neq j}P_+((\mathbb{M}_i\cap\mathbb{M}_j)\times(\mathbb{M}_i\cap\mathbb{M}_j))-\lambda\psi(1)\sum_{i=1}^K P_{\mathbb{X}}(\mathbb{M}_i)\left(1-P_{\mathbb{X}}(\mathbb{M}_i)\right). \tag{16}$$

Therefore, we complete the proof. $\qquad\square$

## B.3  Proof of Corollary 1

*Proof of Corollary 1.* The proof of Corollary 1 is completed by checking whether the equality holds in each inequality that appears in the proof of Theorem 1. We list the detail of the checks below:

(7): Since $w(x,x') = 0$ for any $(x,x') \in \mathbb{M}_i \times \mathbb{M}_j$ $(i \neq j)$, we have $\int_{\mathbb{M}_i\times\mathbb{M}_j}\|h(f(x)) - h(f(x'))\|_{\mathcal{H}_k}^2 w(x,x')d\nu_{\mathbb{X}}(x)d\nu_{\mathbb{X}}(x') = 0$. Here, we have the decomposition $\mathbb{X}\times\mathbb{X} = (\bigcup_{i=1}^K \mathbb{M}_i)\times(\bigcup_{j=1}^K \mathbb{M}_j) = \bigcup_{i,j=1}^K \mathbb{M}_i \times \mathbb{M}_j$, where $(\mathbb{M}_i \times \mathbb{M}_j)\cap(\mathbb{M}_{i'} \times \mathbb{M}_{j'}) = \emptyset$ for any $(i,j,i',j')$ such

that $i \neq i'$ or $j \neq j'$, from the assumption that $\mathbb{M}_1, \cdots, \mathbb{M}_K$ are disjoint. Hence, using the additivity of a probability measure yields the equality.

(8): Since $\mathbb{M}_i \cap \mathbb{M}_j = \emptyset$ for $i, j \in [K]$ such that $i \neq j$, the first term in the right-hand-side of (8) is equal to the first term in the left-hand-side of (8). On the other hand, the second term in the right-hand-side of (8) is equal to 0. Hence, the equality holds.

(9): Since the second term of the right-hand-side of (9) is 0 under the assumption that $\mathbb{M}_i \cap \mathbb{M}_j = \emptyset$ for $i, j$ ($i \neq j$), the equality holds.

(10): The equality holds from the assumption that for any $x, x' \in \mathbb{M}_i$ ($i \in [K]$), $\text{sim}(x, x'; \lambda) = \delta$ holds. Indeed, this assumption implies that $w(x, x') = (\lambda + \delta)w(x)w(x')$ for any $x, x' \in \mathbb{M}_i$.

(11): Since $\mathbb{M}_1 \times \mathbb{M}_1, \cdots, \mathbb{M}_K \times \mathbb{M}_K$ are disjoint, the equality holds.

(12): Since $\bigcup_{i,j=1}^K \mathbb{M}_i \times \mathbb{M}_j = \mathbb{X} \times \mathbb{X}$, the second term of the right-hand-side of (12) is equal to 0. Thus, the equality holds.

(13): The equality holds due to the same reason as (12) above.

(14): Since $\|h(f(x))\|_{\mathcal{H}_k}^2 = k(f(x), f(x)) = \psi(f(x)^\top f(x)) = \psi(1)$ for any $x \in \mathbb{X}$ and $f \in \mathcal{F}$, the equality holds.

(15): Since (15) is the combination of (11), (12), (13), and (14), the equality in (15) holds in this case.

(16): Since (16) is obtained by combining (9), (10), and (15), the equality holds.

Therefore, we obtain the result. $\qquad\square$

# C   Proof in Section 5.2

We present the proof of a generalized version of Theorem 2. The generalized theorem is presented below.

**Theorem 5** (The generalization of Theorem 2). *Suppose that Assumption 1 and 2 hold. Take $K \in \mathbb{N}$ and $\mathbb{M}_1, \cdots, \mathbb{M}_K$ such that the condition (**A**) in Assumption 2 is satisfied. Let $\widetilde{\mathbb{M}}_1, \cdots \widetilde{\mathbb{M}}_K$ be a disjoint partition of $\mathbb{X}$ satisfying $\widetilde{\mathbb{M}}_i \subset \mathbb{M}_i$ for each $i \in [K]$. Define $\widetilde{y} : \mathbb{X} \to [K]$ as $\widetilde{y}(x) = i$ for every $x \in \widetilde{\mathbb{M}}_i$. Then, for each meaningful encoder $f \in \mathcal{F}$, we have*

$$L_{\text{Err}}(f, W_\mu, \beta_\mu; \widetilde{y}) \leq \frac{8(K-1)}{\Delta_{\min}(f) \cdot \min_{i \in [K]} P_\mathbb{X}(\mathbb{M}_i)} \mathfrak{a}(f)$$

*where $\Delta_{\min}(f) = \min_{i \neq j} \|\mu_i(f) - \mu_j(f)\|_{\mathcal{H}_k}^2$.*

*Proof of Theorem 5.* From the definition, we have $\widetilde{\mathbb{M}}_i \subset \mathbb{M}_i$ and $\widetilde{\mathbb{M}}_i \cap \widetilde{\mathbb{M}}_j = \emptyset$ for all the pairs of distinct indices $i, j \in [K]$. Let us recall the definition of $L_{\text{Err}}(f, W_\mu, \beta_\mu; \widetilde{y})$:

$$L_{\text{err}}(f, W_\mu, \beta_\mu; \widetilde{y}) = P_\mathbb{X}\left(g_{f, W_\mu, \beta_\mu}(x) \neq \widetilde{y}(x)\right).$$

Here recall that we let $\arg\max, \arg\min$ also breaks tie arbitrary. For instance, if there are distinct integers $i_1, \cdots, i_j \in [K]$ such that $g_{f,W_\mu,\beta_\mu}(x) = \{i_1, \cdots, i_j\}$, then we define $g_{f,W_\mu,\beta_\mu}(x) = \widetilde{y}(x)$ if $\widetilde{y}(x) \in \{i_1, \cdots, i_j\}$, and $g_{f,W_\mu,\beta_\mu}(x) = i_1$ if $\widetilde{y}(x) \notin \{i_1, \cdots, i_j\}$. The event $\mathbb{A} := \{x \mid g_{f,W_\mu,\beta_\mu}(x) \neq \widetilde{y}(x)\} = \{x \mid g_{f,W_\mu,\beta_\mu}(x) \neq \widetilde{y}(x)\} \cap \bigcup_{i=1}^{K} \widetilde{\mathbb{M}}_i \subset \mathbb{X}$ is a subset of the event $\mathbb{D} := \bigcup_{i=1}^{K} \bigcup_{j \neq i} \{x \mid \|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k} \geq \|h(f(x)) - \mu_j(f)\|_{\mathcal{H}_k}\} \cap \widetilde{\mathbb{M}}_i$, since

$$x \in \mathbb{A}$$

$$\iff \arg\min_{i \in [K]} \|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k} \neq \widetilde{y}(x) \quad \text{and} \quad x \in \bigcup_{i=1}^{K} \widetilde{\mathbb{M}}_i \quad (\text{def. of } g_{f,W_\mu,\beta_\mu} \text{ and } g_{1\text{-NN}})$$

$$\iff x \in \bigcup_{j \neq \widetilde{y}(x)} \{x' \mid \|h(f(x')) - \mu_{\widetilde{y}(x)}(h_f)\|_{\mathcal{H}_k} \geq \|h(f(x')) - \mu_j(f)\|_{\mathcal{H}_k}\}$$

$$\text{and} \quad x \in \bigcup_{i=1}^{K} \widetilde{\mathbb{M}}_i$$

$$\iff x \in \bigcup_{j \neq \widetilde{y}(x)} \{x' \mid \|h(f(x')) - \mu_{\widetilde{y}(x)}(h_f)\|_{\mathcal{H}_k} \geq \|h(f(x')) - \mu_j(f)\|_{\mathcal{H}_k}\} \cap \widetilde{\mathbb{M}}_{\widetilde{y}(x)}$$

$$\implies x \in \bigcup_{i=1}^{K} \bigcup_{j \neq i} \{x' \mid \|h(f(x')) - \mu_i(f)\|_{\mathcal{H}_k} \geq \|h(f(x')) - \mu_j(f)\|_{\mathcal{H}_k}\} \cap \widetilde{\mathbb{M}}_i = \mathbb{D}.$$

Define $\mathbb{L}_{ij} := \{c(\mu_j(f) - \mu_i(f)) \mid c \in \mathbb{R}\} \subset \mathcal{H}_k$ for every $i, j \in [K], i \neq j$. Since each $\mathbb{L}_{ij}$ is a closed subspace of $\mathcal{H}_k$, for every $z \in \mathcal{H}_k$ there exists some $z_1 \in \mathbb{L}_{ij}$ and $z_2 \in \mathbb{L}_{ij}^\perp$ (where $\mathbb{L}_{ij}^\perp$ is the orthogonal complement space of $\mathbb{L}_{ij}$) such that $z$ admits the unique decomposition $z = z_1 + z_2$. Here define the projection $\widetilde{\pi}_{ij} : \mathcal{H}_k \to \mathbb{L}_{ij}$ as $\widetilde{\pi}_{ij}(z) = z_1$, and define the shifted projection $\pi_{ij}$ as $\pi_{ij} : \mathcal{H}_k \to \mathcal{H}_k, \pi_{ij}(z) := \widetilde{\pi}_{ij}(z - \mu_i(f)) + \mu_i(f)$. From the definition, we have that $\|\pi_{ij}(z) - \mu_i(f)\|_{\mathcal{H}_k} \leq \|z - \mu_i(f)\|_{\mathcal{H}_k}$ and $\|\pi_{ij}(z) - \mu_j(f)\|_{\mathcal{H}_k} \leq \|z - \mu_j(f)\|_{\mathcal{H}_k}$.

Hereafter, we use the abbreviation $\Delta_{ij} := \Delta_{ij}(f) = \|\mu_i(f) - \mu_j(f)\|_{\mathcal{H}_k}^2$ for the sake of convenience. Using $\pi_{ij}, i, j \in [K], i \neq j$, the event $\mathbb{D}$ can be decomposed into,

$$\mathbb{D} = \underbrace{\left( \mathbb{D} \cap \left( \bigcup_{i=1}^{K} \bigcup_{j \neq i} \left\{ x \mid \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} \leq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}} \right\} \cap \widetilde{\mathbb{M}}_i \right) \right)}_{= \mathbb{D}_1}$$

$$\cup \underbrace{\left( \mathbb{D} \cap \left( \bigcup_{i=1}^{K} \bigcup_{j \neq i} \left\{ x \mid \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} \leq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}} \right\} \cap \widetilde{\mathbb{M}}_i \right)^c \right)}_{= \mathbb{D}_2}.$$

For $\mathbb{D}_1$, we have

$$P_{\mathbb{X}}(\mathbb{D}_1)$$

$$\leq P_{\mathbb{X}}\left(\bigcup_{i=1}^{K}\bigcup_{j\neq i}\left\{x \mid \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} \leq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\right\} \cap \widetilde{\mathbb{M}}_i\right)$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} P_{\mathbb{X}}\left(\left\{x \mid \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} \leq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\right\} \cap \widetilde{\mathbb{M}}_i\right) \qquad \text{(the union bound)}$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} P_{\mathbb{X}}\left(\left\{x \mid -\|\pi_{ij}(h(f(x))) - \mu_i(f)\|_{\mathcal{H}_k} + \|\mu_i(f) - \mu_j(f)\|_{\mathcal{H}_k} \leq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\right\} \cap \widetilde{\mathbb{M}}_i\right)$$

$$\text{(triangle ineq.)}$$

$$= \sum_{i=1}^{K}\sum_{j\neq i} P_{\mathbb{X}}\left(\left\{x \mid \|\pi_{ij}(h(f(x))) - \mu_i(f)\|_{\mathcal{H}_k} \geq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\right\} \cap \widetilde{\mathbb{M}}_i\right)$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} \frac{4}{\Delta_{ij}}\mathbb{E}\left[\|\pi_{ij}(h(f(x))) - \mu_i(f)\|_{\mathcal{H}_k}^2; \widetilde{\mathbb{M}}_i\right] \qquad \text{(Markov's ineq.)}$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} \frac{4}{\Delta_{ij}}\mathbb{E}\left[\|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k}^2; \widetilde{\mathbb{M}}_i\right] \qquad \text{(def. of } \pi_{ij})$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} \frac{4}{\Delta_{ij}}\mathbb{E}\left[\|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k}^2; \mathbb{M}_i\right]. \qquad \text{(def. of } \widetilde{\mathbb{M}}_i)$$

For $\mathbb{D}_2$, we note that we can rewrite as,

$$P_{\mathbb{X}}(\mathbb{D}_2)$$

$$= P_{\mathbb{X}}\left(\mathbb{D}\cap\left(\bigcup_{i=1}^{K}\bigcup_{j\neq i}\left\{x \mid \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} \leq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\right\} \cap \widetilde{\mathbb{M}}_i\right)^c\right)$$

$$= P_{\mathbb{X}}\left(\left(\bigcup_{i=1}^{K}\bigcup_{j\neq i}\{x \mid \|h(f(x)) - \mu_i(f)\|_2 \geq \|h(f(x)) - \mu_j(f)\|_2\} \cap \widetilde{\mathbb{M}}_i\right) \cap\right.$$

$$\left.\left(\bigcap_{i=1}^{K}\bigcap_{j\neq i}\left\{x \mid \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\right\} \cup \widetilde{\mathbb{M}}_i^c\right)\right)$$

$$= P_{\mathbb{X}}\left(\bigcup_{i=1}^{K}\bigcup_{j\neq i}\bigcap_{i'=1}^{K}\bigcap_{j'\neq i'}\left(\{x \mid \|h(f(x)) - \mu_i(f)\|_2 \geq \|h(f(x)) - \mu_j(f)\|_2\} \cap \widetilde{\mathbb{M}}_i \cap\right.\right.$$

$$\left.\left.\left(\left\{x \mid \|\pi_{i'j'}(h(f(x))) - \mu_{j'}(h_f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{i'j'}^{\frac{1}{2}}\right\} \cup \widetilde{\mathbb{M}}_{i'}^c\right)\right)\right)$$

30

By using above, we have

$$P_{\mathbb{X}}(\mathbb{D}_2)$$

$$\leq P_{\mathbb{X}}\left(\bigcup_{i=1}^{K}\bigcup_{j\neq i}\left(\{x \mid \|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k} \geq \|h(f(x)) - \mu_j(f)\|_{\mathcal{H}_k}\right.\right.$$

$$\left.\left.\text{and } \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\} \cap \widetilde{\mathbb{M}}_i\right)\right) \qquad (17)$$

$$\leq P_{\mathbb{X}}\left(\bigcup_{i=1}^{K}\bigcup_{j\neq i}\left(\{x \mid \|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k} \geq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\} \cap \widetilde{\mathbb{M}}_i\right)\right) \qquad \text{(def. of } \pi_{ij})$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} P_{\mathbb{X}}\left(\{x \mid \|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k} \geq \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\} \cap \widetilde{\mathbb{M}}_i\right) \qquad \text{(the union bound)}$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} \frac{4}{\Delta_{ij}} \mathbb{E}\left[\|\pi_{ij}(h(f(x))) - \mu_i(f)\|_{\mathcal{H}_k}^2 ; \widetilde{\mathbb{M}}_i\right] \qquad \text{(Markov's ineq.)}$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} \frac{4}{\Delta_{ij}} \mathbb{E}\left[\|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k}^2 ; \widetilde{\mathbb{M}}_i\right] \qquad \text{(def. of } \pi_{ij})$$

$$\leq \sum_{i=1}^{K}\sum_{j\neq i} \frac{4}{\Delta_{ij}} \mathbb{E}\left[\|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k}^2 ; \mathbb{M}_i\right]. \qquad \text{(def. of } \widetilde{\mathbb{M}}_i)$$

Here, let us show (17). First let us fix $i, j \in [K]$, where $i \neq j$. For $i', j' \in [K]$ satisfying $i' \neq j'$, we consider the following two cases.

- If $i' = i$ and $j' = j$, then $\widetilde{\mathbb{M}}_i \cap \widetilde{\mathbb{M}}_i^c = \emptyset$, which implies

$$\{x \mid \|h(f(x)) - \mu_i(f)\|_2 \geq \|h(f(x)) - \mu_j(f)\|_2\} \cap \widetilde{\mathbb{M}}_i \cap \left(\left\{x \mid \|\pi_{i'j'}(h(f(x))) - \mu_{j'}(h_f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{i'j'}^{\frac{1}{2}}\right\} \cup \widetilde{\mathbb{M}}_{i'}^c\right)$$

$$= \left\{x \mid \|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k} \geq \|h(f(x)) - \mu_j(f)\|_{\mathcal{H}_k} \text{ and } \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\right\} \cap \widetilde{\mathbb{M}}_i.$$

- if $i' \neq i$ or $j' \neq j$, then

$$\{x \mid \|h(f(x)) - \mu_i(f)\|_2 \geq \|h(f(x)) - \mu_j(f)\|_2\} \cap \widetilde{\mathbb{M}}_i \cap \left(\left\{x \mid \|\pi_{i'j'}(h(f(x))) - \mu_{j'}(h_f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{i'j'}^{\frac{1}{2}}\right\} \cup \widetilde{\mathbb{M}}_{i'}^c\right)$$

$$\subset \widetilde{\mathbb{M}}_i.$$

Thus,

$$\bigcap_{i'=1}^{K} \bigcap_{j' \neq i'} \left( \{x \mid \|h(f(x)) - \mu_i(f)\|_2 \geq \|h(f(x)) - \mu_j(f)\|_2\} \right.$$

$$\left. \cap \widetilde{\mathbb{M}}_i \cap \left( \left\{ x \mid \|\pi_{i'j'}(h(f(x))) - \mu_{j'}(h_f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{i'j'}^{\frac{1}{2}} \right\} \cup \widetilde{\mathbb{M}}_{i'}^c \right) \right)$$

$$\subset \{x \mid \|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k} \geq \|h(f(x)) - \mu_j(f)\|_{\mathcal{H}_k} \text{ and } \|\pi_{ij}(h(f(x))) - \mu_j(f)\|_{\mathcal{H}_k} > \frac{1}{2}\Delta_{ij}^{\frac{1}{2}}\} \cap \widetilde{\mathbb{M}}_i.$$

By combining all the results, we obtain

$$P_{\mathbb{X}}(\mathbb{A}) \leq P_{\mathbb{X}}(\mathbb{D})$$

$$\leq P_{\mathbb{X}}(\mathbb{D}_1) + P_{\mathbb{X}}(\mathbb{D}_2)$$

$$\leq \sum_{i=1}^{K} \sum_{j \neq i} \frac{8}{\Delta_{ij}} \mathbb{E}\left[\|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k}^2; \mathbb{M}_i\right]$$

$$\leq \frac{8(K-1)}{\Delta_{\min}(f)} \sum_{i=1}^{K} \mathbb{E}[\|h(f(x)) - \mu_i(f)\|_{\mathcal{H}_k}^2; \mathbb{M}_i]$$

$$\leq \frac{8(K-1)}{\Delta_{\min}(f)} \sum_{i=1}^{K} \frac{1}{P_{\mathbb{X}}(\mathbb{M}_i)} \mathbb{E}_{x,x^-}[\|h(f(x)) - h(f(x^-))\|_{\mathcal{H}_k}^2; \mathbb{M}_i \times \mathbb{M}_i] \quad \text{(Jensen's inequality)}$$

$$\leq \frac{8(K-1)}{\Delta_{\min}(f) \cdot \min_{i \in [K]} P_{\mathbb{X}}(\mathbb{M}_i)} \mathfrak{a}(f),$$

and we complete the proof. $\qquad \square$

*Proof of Theorem 2.* From the definition of $y$, it is guaranteed that the sets $\{x \in \mathbb{X} \mid y(x) = i\}$ for $i = 1, \cdots, K$ are disjoint and satisfy the relation $\{x \in \mathbb{X} \mid y(x) = i\} \subseteq \mathbb{M}_i$ for every $i \in [K]$. Thus, Theorem 5 can apply to this case, and we obtain the result. $\qquad \square$

# D  Proofs in Section 5.3

## D.1  Proof of Theorem 3

First, we prove Theorem 3. Before that, we present the following theorem, which is a part of the proof of Theorem 3.

**Theorem 6.** *Let $(X_1, X_1'), \cdots, (X_n, X_n')$ be random variables introduced in Section 5.3. Suppose that Assumption 1 holds, and suppose that $n$ is even. Then, with probability at least $1 - \varepsilon$, the following inequality holds:*

$$\sup_{f \in \mathcal{F}} \left( -\frac{1}{n(n-1)} \sum_{i \neq j} k(f(X_i), f(X_j')) + \mathbb{E}_{X,X^-}\left[k(f(X), f(X^-))\right] \right) \leq 2\rho \mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*) + \sqrt{\frac{10b^2 \log{(1/\varepsilon)}}{n}},$$

*where we define $\mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*)$ with the symmetric group $S_n$ of degree $n$:*

$$\mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*) := \max_{s \in S_n} \mathop{\mathbb{E}}_{\substack{X,X' \\ \sigma_{1:(n/2)}}} \left[ \sup_{f \in \mathcal{F}} \frac{2}{n} \sum_{i=1}^{n/2} \sigma_i f(X_{s(2i-1)})^\top f(X'_{s(2i)}) \right].$$

We remark that in Theorem 6, we need to deal with more delicate technical matters compared to the typical generalization error bounds (e.g., Theorem 3.3 of Mohri et al. (2018)), since in our setup $X_1, X_1', \cdots, X_n, X_n'$ are not necessarily independent to each other. We give the proof of Theorem 6 in Appendix D.4.

Now, we can show Theorem 3.

*Proof of Theorem 3.* First observe that,

$$\sup_{f \in \mathcal{F}} \left( -\widehat{L}_{\mathrm{KCL}}(f; \lambda) + L_{\mathrm{KCL}}(f; \lambda) \right)$$

$$= \sup_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n k(f(X_i), f(X_i')) - \frac{\lambda}{n(n-1)} \sum_{i \neq j} k(f(X_i), f(X_j')) \right.$$

$$\left. - \mathbb{E}_{X,X^+} \left[ k(f(X), f(X^+)) \right] + \lambda \mathbb{E}_{X,X^-} \left[ k(f(X), f(X^-)) \right] \right)$$

$$\leq \underbrace{\sup_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n k(f(X_i), f(X_i')) - \mathbb{E}_{X,X^+} \left[ k(f(X), f(X^+)) \right] \right)}_{\text{(i)}}$$

$$+ \lambda \underbrace{\sup_{f \in \mathcal{F}} \left( -\frac{1}{n(n-1)} \sum_{i \neq j} k(f(X_i), f(X_j')) + \mathbb{E}_{X,X^-} \left[ k(f(X), f(X^-)) \right] \right)}_{\text{(ii)}}.$$

Let us define the function space $\mathcal{K} := \{k(f(\cdot), f(\cdot)) : \mathbb{X} \times \mathbb{X} \to \mathbb{R} \mid f \in \mathcal{F}\}$. Then $\mathcal{K}$ is uniformly bounded with constant $b = \sup_{z,z' \in \mathbb{S}^{d-1}} |k(z, z')|$. Here we note that $b < +\infty$ holds since $k$ is continuous and $\mathbb{S}^{d-1}$ is compact; see Section 2.1. From the ULLNs (Theorem 3.3 in Mohri et al. (2018)), with probability at least $1 - \varepsilon/2$, we have

$$\text{(i)} \leq 2\mathfrak{R}_n^+(\mathcal{K}) + \sqrt{\frac{2b^2 \log(2/\varepsilon)}{n}}.$$

Since $k$ is represented by $k(x, x') = \psi(x^\top x')$ for some $\rho$-Lipshitz function $\psi$ from Assumption 1, by applying Talagrand's lemma (Lemma 26.9 in Shalev-Shwartz and Ben-David (2014)) we have $\mathfrak{R}_n^+(\mathcal{K}) \leq \rho \mathfrak{R}_n^+(\mathcal{Q})$. Hence, with probability at least $1 - \varepsilon/2$, we have

$$\text{(i)} \leq 2\rho \mathfrak{R}_n^+(\mathcal{Q}) + \sqrt{\frac{2b^2 \log(2/\varepsilon)}{n}}.$$

For (ii), from Theorem 6, with probability at least $1 - \varepsilon/2$ we have

$$\text{(ii)} \leq 2\rho \mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*) + \sqrt{\frac{10b^2 \log(2/\varepsilon)}{n}}.$$

Therefore, with probability at least $1 - \varepsilon$ we have,

$$\sup_{f \in \mathcal{F}} \left( -\widehat{L}_{\mathrm{KCL}}(f; \lambda) + L_{\mathrm{KCL}}(f; \lambda) \right) \leq 2\rho \mathfrak{R}_n(\mathcal{Q}) + \sqrt{\frac{2b^2 \log(2/\varepsilon)}{n}} + \lambda \sqrt{\frac{10b^2 \log(2/\varepsilon)}{n}}, \qquad (18)$$

where $\mathfrak{R}_n(\mathcal{Q}) := \mathfrak{R}_n^+(\mathcal{Q}) + \lambda \mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*)$.

Note that in the same way as the proof of the above probability bound, we have the following inequality: with probability at least $1 - \varepsilon$,

$$\sup_{f \in \mathcal{F}} \left( \widehat{L}_{\mathrm{KCL}}(f; \lambda) - L_{\mathrm{KCL}}(f; \lambda) \right) \leq 2\rho \mathfrak{R}_n(\mathcal{Q}) + \sqrt{\frac{2b^2 \log(2/\varepsilon)}{n}} + \lambda \sqrt{\frac{10b^2 \log(2/\varepsilon)}{n}}. \qquad (19)$$

Hence, let $\widehat{f}$ be the minimizer of $\widehat{L}_{\mathrm{KCL}}(f; \lambda)$, then from (18) and (19), with probability at least $1 - 2\varepsilon$ we have

$$L_{\mathrm{KCL}}(\widehat{f}; \lambda) \leq L_{\mathrm{KCL}}(f; \lambda) + 4\rho \mathfrak{R}_n(\mathcal{Q}) + 2\sqrt{\frac{2b^2 \log(2/\varepsilon)}{n}} + 2\lambda \sqrt{\frac{10b^2 \log(2/\varepsilon)}{n}},$$

where we note that $\widehat{L}_{\mathrm{KCL}}(\widehat{f}; \lambda) \leq \widehat{L}_{\mathrm{KCL}}(f; \lambda)$ from the definition of $\widehat{f}$. Therefore, we complete the proof. $\qquad \square$

## D.2 An Upper Bound of the Rademacher Complexity

In this section for the sake of simplicity, we consider the case in which for every $f \in \mathcal{F}$, there exists the unique function $f_0 \in \mathcal{F}_0$ such that $f(x) = f_0(x)/\|f_0(x)\|_2$ for every $x \in \mathbb{X}$. First let us recall the definition of a sub-Gaussian process:

**Definition 4** (Quoted from Definition 5.16 in Wainwright (2019))**.** *A collection of zero-mean random variables $\{X_\theta, \theta \in \mathbb{T}\}$ is a sub-Gaussian process with respect to a metric $\rho_X$ on $\mathbb{T}$ if*

$$\mathbb{E}\left[ e^{\lambda(X_\theta - X_{\widetilde{\theta}})} \right] \leq e^{\frac{\lambda^2 \rho_X^2(\theta, \widetilde{\theta})}{2}} \quad \text{for all } \theta, \widetilde{\theta} \in \mathbb{T}, \text{ and } \lambda \in \mathbb{R}.$$

We next upper bound the Rademacher complexity via the chaining technique (Theorem 5.22 in Wainwright (2019)).

**Proposition 2.** *Suppose $n$ is even. For $\mathfrak{R}_n(\mathcal{Q})$, we have the upper bound,*

$$\mathfrak{R}_n(\mathcal{Q}) \leq \frac{64(1 + \sqrt{2}\lambda)}{\mathfrak{m}(\mathcal{F}_0)\sqrt{n}} \int_0^{Cd} \sqrt{\log \mathfrak{C}(u; \mathcal{F}_0, \|\cdot\|_\infty)} du,$$

*where $\|f_0\|_\infty := \sup_{x \in \mathbb{X}} \|f_0(x)\|_2$ for $f_0 \in \mathcal{F}_0$, $\mathfrak{m}(\mathcal{F}_0)$ is defined in Section 2.1, $C$ is a constant independent of $d, n, \lambda$, and $\mathfrak{C}(u; \mathcal{F}_0, \|\cdot\|_\infty)$ is the $u$-covering number of $(\mathcal{F}_0, \|\cdot\|_\infty)$ (for the definition of covering number, see e.g., Definition 5.1 in Wainwright (2019)).*

*Proof of Proposition 2.* In this proof, we follow the proof idea of Tu et al. (2019) (see Lemma 5 in Tu et al. (2019)). Since our setup is different from Tu et al. (2019), we need to modify the proof and add several new techniques. Define

$$Z_{f_0} := \frac{\mathfrak{m}(\mathcal{F}_0)}{2\sqrt{n}} \sum_{i=1}^{n} \sigma_i q(f_0(X_i), f_0(X_i')),$$

where $\sigma_1, \cdots, \sigma_n$ are Rademacher random variables that are independent to each other and to each $(X_i, X_i')$, $i \in [n]$, $f_0 \in \mathcal{F}_0$, $(X_1, X_1), \cdots, (X_n, X_n')$ are the random vectors defined in Section 5.3, and

$$q(z, z') = \frac{z^\top z'}{\|z\|_2 \cdot \|z'\|_2} \quad z, z' \in \mathbb{S}^{d-1}.$$

Also, let us recall the assumption for $\mathcal{F}_0$ introduced in Section 2.1: for every $f \in \mathcal{F}$, there exists the unique function $f_0 \in \mathcal{F}_0$ such that $f(x) = f_0(x)/\|f_0(x)\|_2$ for every $x \in \mathbb{X}$. We show that $\{Z_{f_0}\}_{f_0 \in \mathcal{F}_0}$ is a sub-Gaussian process as follows: note that, for every $f_{1,0}, f_{2,0} \in \mathcal{F}_0$,

$$\frac{\mathfrak{m}(\mathcal{F}_0)}{2\sqrt{n}} \left| \sigma_i(q(f_{1,0}(X_i), f_{1,0}(X_i')) - q(f_{2,0}(X_i), f_{2,0}(X_i'))) \right|$$

$$\leq \frac{\mathfrak{m}(\mathcal{F}_0)}{2\sqrt{n}} \left| f_1(X_i)^\top f_1(X_i') - f_2(X_i)^\top f_2(X_i') \right| \qquad \text{(def. of } \sigma_i\text{)}$$

$$\leq \frac{\mathfrak{m}(\mathcal{F}_0)}{2\sqrt{n}} \left( \left| f_1(X_i)^\top f_1(X_i') - f_1(X_i)^\top f_2(X_i') \right| + \left| f_1(X_i)^\top f_2(X_i') - f_2(X_i)^\top f_2(X_i') \right| \right)$$

$$\text{(triangle ineq.)}$$

$$\leq \frac{\mathfrak{m}(\mathcal{F}_0)}{2\sqrt{n}} \left( \|f_1(X_i)\|_2 \|f_1(X_i') - f_2(X_i')\|_2 + \|f_1(X_i) - f_2(X_i)\|_2 \|f_2(X_i')\|_2 \right)$$

$$\text{(Cauchy-Schwarz ineq.)}$$

$$\leq \frac{\mathfrak{m}(\mathcal{F}_0)}{2\sqrt{n}} \left( \|f_1(X_i') - f_2(X_i')\|_2 + \|f_1(X_i) - f_2(X_i)\|_2 \right) \qquad \text{(def. of } f_1, f_2\text{)}$$

$$\leq \frac{\mathfrak{m}(\mathcal{F}_0)}{\sqrt{n}} \sup_{x \in \mathbb{X}} \|f_1(x) - f_2(x)\|_2$$

$$= \frac{\mathfrak{m}(\mathcal{F}_0)}{\sqrt{n}} \sup_{x \in \mathbb{X}} \left\| \frac{f_{1,0}(x)}{\|f_{1,0}(x)\|_2} - \frac{f_{2,0}(x)}{\|f_{2,0}(x)\|_2} \right\|_2 \qquad \text{(the uniqueness of } f_{1,0}, f_{2,0}\text{)}$$

$$\leq \frac{1}{\sqrt{n}} \|f_{1,0} - f_{2,0}\|_\infty. \qquad \text{(def. of } \mathfrak{m}(\mathcal{F}_0)\text{)}$$

Hence, we have

$$\mathbb{E}_{X_{1:n}, X_{1:n}', \sigma_{1:n}} \left[ \exp\left( t(Z_{f_{1,0}} - Z_{f_{2,0}}) \right) \right] \leq \exp\left( \frac{t^2}{2n} \|f_{1,0} - f_{2,0}\|_\infty^2 \right)^n = \exp\left( \frac{t^2}{2} \|f_{1,0} - f_{2,0}\|_\infty^2 \right).$$

This indicates that $\{Z_{f_0}\}_{f \in \mathcal{F}_0}$ is a sub-Gaussian process with the norm $\| \cdot \|_\infty$. Here note that $\sup_{f_{1,0}, f_{2,0} \in \mathcal{F}_0} \|f_{1,0} - f_{2,0}\|_\infty \leq C\sqrt{d}$ for some constant $C \in \mathbb{R}$ that is independent of $d$, since $\mathcal{F}_0$ is uniformly bounded. By using the chaining theorem (Theorem 5.22 in Wainwright (2019)), we have

$$\mathfrak{R}_n^+(\mathcal{Q}) \leq \frac{64}{\mathfrak{m}(\mathcal{F}_0)\sqrt{n}} \int_0^{C\sqrt{d}} \sqrt{\log \mathfrak{C}(u; \mathcal{F}_0, \| \cdot \|_\infty)} du.$$

For $\mathfrak{R}_{n/2}^-(\mathcal{Q})$, in a similar way we obtain,

$$\mathfrak{R}_{n/2}^-(\mathcal{Q}) \leq \frac{64\sqrt{2}}{\mathfrak{m}(\mathcal{F}_0)\sqrt{n}} \int_0^{C\sqrt{d}} \sqrt{\log \mathfrak{C}(u; \mathcal{F}_0, \|\cdot\|_\infty)} du.$$

Thus, we have

$$\mathfrak{R}_n(\mathcal{Q}) \leq \frac{64(1 + \sqrt{2}\lambda)}{\mathfrak{m}(\mathcal{F}_0)\sqrt{n}} \int_0^{C\sqrt{d}} \sqrt{\log \mathfrak{C}(u; \mathcal{F}_0, \|\cdot\|_\infty)} du,$$

and complete the proof. □

The integral in the above upper bound is often called Dudley entropy integral (Wainwright, 2019). Proposition 2 makes it easier to derive a generalization bound via chaining, since it is enough to evaluate the Dudley entropy integral for the function space $\mathcal{F}_0$ instead of the space of critic functions $\mathcal{Q}$.

Here, denote by $\mathfrak{D}(\mathcal{F}_0, \|\cdot\|_\infty)$, the Dudley entropy integral w.r.t. $(\mathcal{F}_0, \|\cdot\|_\infty)$, i.e.,

$$\mathfrak{D}(\mathcal{F}_0, \|\cdot\|_\infty) = \int_0^{C\sqrt{d}} \sqrt{\log \mathfrak{C}(u; \mathcal{F}_0, \|\cdot\|_\infty)} du.$$

It is shown by Tu et al. (2019) that if $\mathcal{F}_0$ is a function space of feedforward (deep) neural networks, where each neural networks have weight matrices whose norms are bounded by some universal constant, and Lipschitz activation functions that vanish at the origin, then $\mathfrak{D}(\mathcal{F}_0, \|\cdot\|_\infty) < +\infty$ holds. Based on this fact, we introduce:

**Assumption 3.** *The Dudley entropy integral $\mathfrak{D}(\mathcal{F}_0, \|\cdot\|_\infty)$ is finite, and $\mathfrak{R}_n(\mathcal{Q}) \leq O((1 + \lambda)/\sqrt{n})$ holds.*

Consequently, we obtain the generalization error bound.

**Corollary 2.** *Suppose that Assumption 1, 3 hold, and $n$ is even. Then, with probability at least $1 - \varepsilon$ where $\varepsilon > 0$, we have*

$$L_{\mathrm{KCL}}(f; \lambda) \leq \widehat{L}_{\mathrm{KCL}}(f; \lambda) + O\left(\frac{(1 + \lambda)\left(1 + \sqrt{\log(2/\varepsilon)}\right)}{\sqrt{n}}\right).$$

*Proof.* Due to Theorem 3 and Assumption 3. □

## D.3 Useful Results on McDiarmid's Inequality for Dependent Random Variables

Before showing Theorem 6, we need to prepare several definitions and an existing result. The following three definitions are quoted from Zhang et al. (2019).

**Definition 5** (Dependency Graph, quoted from Definition 3.1 in Zhang et al. (2019))**.** *An undirected graph $G$ is called a dependency graph of a random vector $\mathbf{X} = (X_1, \cdots, X_n)$ if*

1. $V(G) = [n]$

2. if $I, J \subset [n]$ are non-adjacent in $G$, then $\{X_i\}_{i \in I}$ and $\{X_j\}_{j \in J}$ are independent.

**Definition 6** (Forest Approximation, quoted from Definition 3.4 in Zhang et al. (2019)). *Given a graph $G$, a forest $F$, and a mapping $\phi : V(G) \to V(F)$, if $\phi(u) = \phi(v)$ or $\langle \phi(u), \phi(v) \rangle \in E(F)$ for any $\langle u, v \rangle \in E(G)$, we say that $(\phi, F)$ is a forest approximation of $G$. Let $\Phi(G)$ denote the set of forest approximations of $G$.*

**Definition 7** (Forest Complexity, quoted from Definition 3.5 in Zhang et al. (2019)). *Given a graph $G$ and any forest approximation $(\phi, F) \in \Phi(G)$ with $F$ consisting of trees $\{T_i\}_{i \in [k]}$, let*

$$\lambda_{(\phi,F)} = \sum_{\langle u,v \rangle \in E(F)} \left( |\phi^{-1}(u)| + |\phi^{-1}(v)| \right)^2 + \sum_{i=1}^{k} \min_{u \in V(T_i)} |\phi^{-1}(u)|^2.$$

*We call*

$$\Lambda(G) = \min_{(\phi,F) \in \Phi(G)} \lambda_{(\phi,F)}$$

*the forest complexity of the graph $G$.*

Zhang et al. (2019) have shown the following result, which is an extension of McDiarmid's inequality (McDiarmid, 1989) for dependent random variables.

**Theorem 7** (Quoted from Theorem 3.6 in Zhang et al. (2019)). *Suppose that $f : \mathbf{\Omega} \to \mathbb{R}$ is a $\mathbf{c}$-Lipschitz function and $G$ is a dependency graph of a random vector $\mathbf{X}$ that takes values in $\mathbf{\Omega}$. For any $t > 0$, the following inequality holds:*

$$\Pr(f(\mathbf{X}) - \mathbf{E}[f(\mathbf{X})] \geq t) \leq \exp\left( -\frac{2t^2}{\Lambda(G)\|\mathbf{c}\|_\infty^2} \right).$$

Note that, in the above theorem $f : \mathbf{\Omega} \to \mathbb{R}$ is said to be $\mathbf{c}$-Lipschitz if $|f(\mathbf{x}) - f(\mathbf{x}')| \leq \sum_{i=1}^{p} \mathbf{c}_i \mathbb{1}_{\{\mathbf{x}_i \neq \mathbf{x}_i'\}}$ for every $\mathbf{x}, \mathbf{x}' \in \mathbf{\Omega}$, where $\mathbf{\Omega} \subset \mathbb{R}^p$ for some $p \in \mathbb{N}$.

## D.4 Proof of Theorem 6

We show Theorem 6 by utilizing the contents in Appendix D.3. Recall the definition of the random variables introduced in Section 5.3: $(X_1, X_1'), \cdots, (X_n, X_n')$ are pairs of random variables sampled independently according to the joint probability distribution with density $w(x, x')$, where $X_i$ and $X_j'$ are independent for each pair of distinct indices $i, j \in [K]$. From the definition, the following claim holds.

**Lemma 3.** *Let $G_n$ be a dependency graph that is defined with a random vector $(X_1, X_1', \cdots, X_n, X_n')$, where the edges in $G_n$ are defined as follows: for any $i, j \in [n]$, $X_i$ and $X_j$ are not connected, and $X_i$ and $X_j'$ are connected by an edge if and only if $i = j$. Then, we have $\Lambda(G_n) \leq 5n$.*

*Proof.* Let $\phi : G_n \to G_n$ be the identity map. From the definition, $G_n$ can be decomposed into trees $\{T_i\}_{i \in [n]}$ where $V(T_i) = \{X_i, X_i'\}$ for each $i \in [n]$. Let $F$ be the forest consisting of the trees $\{T_i\}_{i \in [n]}$. Then, we have $\lambda_{(\phi,F)} = 5n$, which implies $\Lambda(G_n) \leq \lambda_{(\phi,F)} \leq 5n$. □

*Proof of Theorem 6.* The goal of this proof is to upper bound the following quantity with high probability:

$$\sup_{f \in \mathcal{F}} \left( -\frac{1}{n(n-1)} \sum_{i \neq j} k(f(X_i), f(X_j')) + \mathbb{E}_{X,X^-} \left[ k(f(X_i), f(X^-)) \right] \right).$$

However, as explained before, the standard argument (see e.g., Theorem 3.3 in Mohri et al. (2018)) cannot apply to this case since $k(f(X_i), f(X_j'))$, $i, j \in [n], i \neq j$ are not necessarily independent to each other from our problem setup. We instead utilize the McDiarmid's inequality for dependent random variables, which is shown by Zhang et al. (2019), to avoid this problem. Our proof below is mainly based on Theorem 3.3 in Mohri et al. (2018), but it includes some modification due to the application of the results by Zhang et al. (2019). Let $\widetilde{X}_1, \widetilde{X}_1', \cdots, \widetilde{X}_n, \widetilde{X}_n'$ be i.i.d. random variables to the original random variables $X_1, X_1', \cdots, X_n, X_n'$. Define the measurable function $F(f) := F(f)(x_1, x_1', \cdots, x_n, x_n')$ on $\mathbb{X}^{2n}$ as

$$F(f) := \frac{1}{n(n-1)} \sum_{i \neq j} k(f(x_i), f(x_j')) - \mathbb{E}_{X,X^-} \left[ k(f(X), f(X^-)) \right].$$

For simplicity, denote

$$F(f)_{x_\ell} := \frac{1}{n(n-1)} \left( \sum_{j:j \neq \ell} k(f(\widetilde{x}_\ell), f(x_j')) + \sum_{\substack{i,j:i \neq j \\ i \neq \ell}} k(f(x_i), f(x_j')) \right) - \mathbb{E}_{X,X^-} \left[ k(f(X), f(X^-)) \right].$$

In a similar way, we also use the notation $F(f)_{x_\ell'}$. Let $j \in [n]$. Then, for every $f \in \mathcal{F}$, we have

$$F(f) - \sup_{f \in \mathcal{F}} F(f)_{X_j} \leq F(f) - F(f)_{X_j} \leq \left| F(f) - F(f)_{X_j} \right|$$

$$\leq \left| \frac{1}{n(n-1)} \sum_{i \in [n], i \neq j} \left( k(f(X_j), f(X_i')) - k(f(\widetilde{X}_j), f(X_i')) \right) \right|$$

$$\leq \frac{1}{n(n-1)} \cdot 2(n-1)b = \frac{2b}{n},$$

where $b := \sup_{z,z' \in \mathbb{S}^{d-1}} |k(z, z')|$. Hence, $\sup_{f \in \mathcal{F}} F(f) - \sup_{f \in \mathcal{F}} F(f)_{X_j} \leq \frac{2b}{n}$. By applying the same argument several times, $\sup_{f \in \mathcal{F}} F(f)$ satisfies the assumption of Theorem 7. Therefore, from Theorem 7 (i.e., Theorem 3.6 in Zhang et al. (2019)) and Lemma 3, with probability at least $1 - \varepsilon$ we have

$$\sup_{f \in \mathcal{F}} F(f) \leq \mathbb{E} \left[ \sup_{f \in \mathcal{F}} F(f) \right] + \sqrt{\frac{10b^2 \log(1/\varepsilon)}{n}}. \tag{20}$$

Let $\sigma_{1:n/2} := (\sigma_1, \cdots, \sigma_{n/2})$ be a random vector that consists of a Rademacher random variable (i.e., a random variable taking $\pm 1$ with probability $1/2$ each) for each entry, and let $\widetilde{X}_{1:n}, \widetilde{X}'_{1:n}$ be i.i.d. copies of the random vectors $X_{1:n}, X'_{1:n}$, respectively. Denote $m = n/2 \in \mathbb{N}$. Then,

$$
\mathbb{E}\left[\sup_{f \in \mathcal{F}} F(f)\right]
$$

$$
= \mathbb{E}\left[\sup_{f \in \mathcal{F}}\left(\frac{1}{n(n-1)}\sum_{i \neq j} k(f(X_i), f(X'_j)) - \mathbb{E}_{X, X^-}\left[k(f(X), f(X^-))\right]\right)\right]
$$

$$
= \mathbb{E}\left[\sup_{f \in \mathcal{F}}\left(\frac{1}{n(n-1)}\sum_{i \neq j} k(f(X_i), f(X'_j)) - \mathbb{E}_{\widetilde{X}_{1:n}, \widetilde{X}'_{1:n}}\left[\frac{1}{n(n-1)}\sum_{i \neq j} k(f(\widetilde{X}_i), f(\widetilde{X}'_j))\right]\right)\right]
$$

$$
= \mathbb{E}\left[\sup_{f \in \mathcal{F}}\left(\frac{1}{n!m}\sum_{s \in S_n}\left(\sum_{i=1}^{m} k(f(X_{s(2i-1)}), f(X'_{s(2i)})) - \mathbb{E}_{\widetilde{X}'_{1:n}, \widetilde{X}'_{1:n}}\left[\sum_{i=1}^{m} k(f(\widetilde{X}_{s(2i-1)}), f(\widetilde{X}'_{s(2i)}))\right]\right)\right)\right] \tag{21}
$$

$$
\leq \frac{1}{n!}\sum_{s \in S_n} \mathbb{E}\left[\sup_{f \in \mathcal{F}}\left(\frac{1}{m}\sum_{i=1}^{m} k(f(X_{s(2i-1)}), f(X'_{s(2i)})) - \mathbb{E}_{\widetilde{X}_{1:n}, \widetilde{X}'_{1:n}}\left[\frac{1}{m}\sum_{i=1}^{m} k(f(\widetilde{X}_{s(2i-1)}), f(\widetilde{X}'_{s(2i)}))\right]\right)\right]
$$

$$
\leq \frac{1}{n!}\sum_{s \in S_n} \mathbb{E}_{\substack{X_{1:n}, X'_{1:n} \\ \widetilde{X}_{1:n}, \widetilde{X}'_{1:n}}}\left[\sup_{f \in \mathcal{F}}\left(\frac{1}{m}\sum_{i=1}^{m}\left(k(f(X_{s(2i-1)}), f(X'_{s(2i)})) - k(f(\widetilde{X}_{s(2i-1)}), f(\widetilde{X}'_{s(2i)}))\right)\right)\right]
$$

$$
= \frac{1}{n!}\sum_{s \in S_n} \mathbb{E}_{\substack{X_{1:n}, X'_{1:n} \\ \widetilde{X}_{1:n}, \widetilde{X}'_{1:n} \\ \sigma_{1:m}}}\left[\sup_{f \in \mathcal{F}}\left(\frac{1}{m}\sum_{i=1}^{m}\sigma_i\left(k(f(X_{s(2i-1)}), f(X'_{s(2i)})) - k(f(\widetilde{X}_{s(2i-1)}), f(\widetilde{X}'_{s(2i)}))\right)\right)\right] \tag{22}
$$

$$
\leq \frac{2}{n!}\sum_{s \in S_n} \mathbb{E}_{\substack{X_{1:n}, X'_{1:n} \\ \sigma_{1:m}}}\left[\sup_{f \in \mathcal{F}}\frac{1}{m}\sum_{i=1}^{m}\sigma_i k(f(X_{s(2i-1)}), f(X'_{s(2i)}))\right]
$$

$$
\leq \frac{2\rho}{n!}\sum_{s \in S_n} \mathbb{E}_{\substack{X_{1:n}, X'_{1:n} \\ \sigma_{1:m}}}\left[\sup_{f \in \mathcal{F}}\frac{1}{m}\sum_{i=1}^{m}\sigma_i f(X_{s(2i-1)})^\top f(X'_{s(2i)})\right] \tag{23}
$$

$$
\leq 2\rho \max_{s \in S_n} \mathbb{E}_{\substack{X_{1:n}, X'_{1:n} \\ \sigma_{1:m}}}\left[\sup_{f \in \mathcal{F}}\frac{1}{m}\sum_{i=1}^{m}\sigma_i f(X_{s(2i-1)})^\top f(X'_{s(2i)})\right]
$$

$$
= 2\rho\mathfrak{R}_m^-(\mathcal{Q}; s^*),
$$

where in (21) we define $S_n$ as the symmetric group of degree $n$ (see Remark 3 for the relation to the *average of "sums-of-i.i.d." blocks* technique for $U$-statistics which is explained in Clémençon et al. (2008)). Besides in (22), for every $s \in S_n$ the random vectors $(X_{s(2i-1)}, X_{s(2i)})$, $(\widetilde{X}_{s(2i-1)}, \widetilde{X}_{s(2i)})$ for $i = 1, \cdots, m$ are independent and identically distributed, which implies that the standard symmetrization argument (Theorem 4.10 of Wainwright (2019)) is applicable. Finally, in (23), under Assumption 1, we apply Talagrand's lemma (Lemma 26.9 in Shalev-Shwartz and Ben-David

(2014)). Therefore, we obtain with probability at least $1 - \varepsilon$,

$$\sup_{f \in \mathcal{F}} F(f) \leq 2\rho \mathfrak{R}_{n/2}^-(\mathcal{Q}; s^*) + \sqrt{\frac{10b^2 \log{(1/\varepsilon)}}{n}}.$$

Thus, we obtain the claim. □

**Remark 3.** *In* (21) *of the proof of Theorem 6, we use the identity,*

$$\frac{1}{n(n-1)} \sum_{i \neq j} k(f(X_i), f(X_j')) = \frac{1}{n! m} \sum_{s \in S_n} \sum_{i=1}^{m} k(f(X_{s(2i-1)}), f(X_{s(2i)}')).$$

*We notice that the above identity is closely related to the average of "sum-of-i.i.d." blocks technique explained in Appendix A of Clémençon et al. (2008). As well as the technique presented in Clémençon et al. (2008), in* (21) *of our paper we also decompose the sum $\sum_{i \neq j} k(f(X_i), f(X_j'))$ into the sums of the i.i.d. random variables. However, we remark that the definition of the sum $\sum_{i \neq j} k(f(X_i), f(X_j'))$ is different from that presented in Clémençon et al. (2008): indeed, in our case, the random variables $f(X_1), f(X_1'), \cdots, f(X_n), f(X_n')$ are not necessarily independent of each other. To address this problem, we decompose our sum in* (21) *as follows: for $2n$ random variables $X_1, X_1', \cdots, X_n, X_n'$, we create the tuples $(X_{s(1)}, X_{s(2)}', \cdots, X_{s(n-1)}, X_{s(n)}')$ where $s \in S_n$, then sum up all the components $\{\sum_{i=1}^{n} k(f(X_{s(2i-1)}), f(X_{s(2i)}'))\}_{s \in S_n}$.*

# E    Proof in Section 5.4

*Proof of Theorem 4.* First applying Theorem 2 to the empirical loss minimizer $\widehat{f}$, we have

$$L_{\text{Err}}(\widehat{f}, W_\mu, \beta_\mu; y) \leq \frac{8(K-1)}{\Delta_{\min}(\widehat{f}) \cdot \min_{i \in [K]} P_{\mathbb{X}}(\mathbb{M}_i)} \mathfrak{a}(\widehat{f}). \tag{24}$$

Using Theorem 1, we have the inequality,

$$\mathfrak{a}(\widehat{f}) \leq L_{\text{KCL}}(\widehat{f}; \lambda) + (1 - \frac{\delta}{2})\mathfrak{a}(\widehat{f}) - \lambda \mathfrak{c}(\widehat{f}) + R(\lambda). \tag{25}$$

Combining (24) and (25), we obtain

$$L_{\text{Err}}(\widehat{f}, W_\mu, \beta_\mu; y) \leq \frac{8(K-1)}{\Delta_{\min}(\widehat{f}) \cdot \min_{i \in [K]} P_{\mathbb{X}}(\mathbb{M}_i)} \left( L_{\text{KCL}}(\widehat{f}; \lambda) + (1 - \frac{\delta}{2})\mathfrak{a}(\widehat{f}) - \lambda \mathfrak{c}(\widehat{f}) + R(\lambda) \right). \tag{26}$$

Here, using the standard technique for upper bounding the optimal classification loss or error (Arora et al., 2019; Ash et al., 2022), the classification error $L_{\text{Err}}(\widehat{f}, W_\mu, \beta_\mu; y)$ is lower bounded as

$$L_{\text{Err}}(\widehat{f}, W^* \beta^*; y) = \inf_{W, \beta} L_{\text{Err}}(\widehat{f}, W, \beta; y) \leq L_{\text{Err}}(\widehat{f}, W_\mu, \beta_\mu; y). \tag{27}$$

From (26) and (27),

$$L_{\mathrm{Err}}(\widehat{f}, W^*, \beta^*; y) \leq \frac{8(K-1)}{\Delta_{\min}(\widehat{f}) \cdot \min_{i \in [K]} P_{\mathbb{X}}(\mathbb{M}_i)} \left( L_{\mathrm{KCL}}(\widehat{f}; \lambda) + (1 - \frac{\delta}{2})\mathfrak{a}(\widehat{f}) - \lambda\mathfrak{c}(\widehat{f}) + R(\lambda) \right).$$

(28)

Applying Theorem 3 to (28), we obtain: with probability at least $1 - 2\varepsilon$,

$$L_{\mathrm{Err}}(\widehat{f}, W_\mu, \beta_\mu; y) \lesssim L_{\mathrm{KCL}}(f; \lambda) + (1 - \frac{\delta}{2})\mathfrak{a}(\widehat{f}) - \lambda\mathfrak{c}(\widehat{f}) + R(\lambda) + 2\mathrm{Gen}(n, \lambda, \varepsilon),$$

where $\lesssim$ omits the coefficient $\frac{8(K-1)}{\Delta_{\min}(\widehat{f}) \cdot \min_{i \in [K]} P_{\mathbb{X}}(\mathbb{M}_i)}$. Therefore, we obtain the result. $\qquad\square$

# F    Additional Information, Results, and Discussion

## F.1    Examples Satisfying Assumption 2

### F.1.1    Proofs in Example 1

We show the several claims that appear in Example 1 as a proposition.

**Proposition 3.** *Let $r > 0$, $K \in \mathbb{N}$, and $v_1, \cdots, v_K \in \mathbb{R}^p$. For each $i \in [K]$, let $\mathbb{B}_i \subset \mathbb{R}^p$ be the open ball of radius $r$ centered at a point $v_i$. Suppose $\mathbb{B}_1, \cdots, \mathbb{B}_K$ are disjoint to each other. Define $\overline{\mathbb{X}} = \bigcup_{i=1}^K \mathbb{B}_i$, $\mathbb{X} = \overline{\mathbb{X}}$, and the conditional probability $a(x|\overline{x}) = \mathrm{vol}(\mathbb{B}_1)^{-1} \sum_{i=1}^K \mathbb{1}_{\mathbb{B}_i \times \mathbb{B}_i}(x, \overline{x})$, where $\mathrm{vol}(\mathbb{B}_1)$ be the volume of $\mathbb{B}_i$ in $\mathbb{R}^p$. Let $p_{\overline{\mathbb{X}}}(\overline{x}) := (K\mathrm{vol}(\mathbb{B}_1))^{-1}$ be a probability density function of $P_{\overline{\mathbb{X}}}$. Define $y : \mathbb{X} \to [K]$ as $y(x) = i$ if $x \in \mathbb{B}_i$. Then, we have the following properties:*

1. *$w(x) > 0$ for every $x \in \mathbb{X}$.*

2. *$\mathrm{sim}(x, x'; \lambda) = K\mathbb{1}_{\bigcup_{i \in [K]} \mathbb{B}_i \times \mathbb{B}_i}(x, x') - \lambda$ for every $x, x' \in \mathbb{X}$.*

3. *Let $\delta \in (-\lambda, K - \lambda]$. Then, $\delta$, $K$, $\mathbb{B}_1, \cdots, \mathbb{B}_K$, and $y$ satisfy Assumption 2.*

*Proof.* We first show the claim 1. From the definition of $w(x)$, for every $x \in \mathbb{B}_1$ we have

$$w(x) = \int_{\overline{\mathbb{X}}} a(x|\overline{x})p_{\overline{\mathbb{X}}}(\overline{x})d\overline{x} = \int_{\mathbb{B}_1} \frac{1}{K(\mathrm{vol}(\mathbb{B}_1))^2}d\overline{x} = \frac{1}{K\mathrm{vol}(\mathbb{B}_1)}.$$

Similarly, for each $i \in [K]$ we obtain $w(x) = (K\mathrm{vol}(\mathbb{B}_1))^{-1}$ for every $x \in \mathbb{B}_i$. Since $\mathbb{X} = \overline{\mathbb{X}} = \bigcup_{i=1}^K \mathbb{B}_i$, we have that $w(x) > 0$ for every $x \in \mathbb{X}$.

Next, let us show the claim 2. From the claim 1, the function $\mathrm{sim}(x, x'; \lambda)$ is well-defined. To compute $\mathrm{sim}(x, x'; \lambda)$, we need to know the function $w(x, x')$. The computation of $w(x, x')$ is done

as follows:

$$w(x, x') = \int_{\overline{\mathbb{X}}} a(x|\overline{x}) a(x'|\overline{x}) p_{\overline{\mathbb{X}}}(\overline{x}) d\overline{x}$$

$$= \begin{cases} \int_{\mathbb{B}_i} \frac{1}{K(\text{vol}(\mathbb{B}_1))^3} d\overline{x} & \text{if } x, x' \in \mathbb{B}_i \text{ for some } i \in [K] \\ 0 & \text{if } x \in \mathbb{B}_i \text{ and } x' \in \mathbb{B}_j \text{ for some } i \neq j \end{cases}$$

$$= \begin{cases} \frac{1}{K(\text{vol}(\mathbb{B}_1))^2} & \text{if } x, x' \in \mathbb{B}_i \text{ for some } i \in [K] \\ 0 & \text{if } x \in \mathbb{B}_i \text{ and } x' \in \mathbb{B}_j \text{ for some } i \neq j. \end{cases}$$

Hence, it is obvious that the claim 2 holds.

Finally, let us prove the claim 3. However, from the claim 2 we see that $\text{sim}(x, x'; \lambda) \geq \delta$ if and only if $x, x' \in \mathbb{B}_i$ for some $i \in [K]$. Furthermore, $y$ is well-defined and the set $\{x \in \mathbb{X} \mid y(x) = i\} = \mathbb{B}_i$ is measurable for every $i \in [K]$. Thus, the claim 3 is also true, and we end the proof. $\qquad\square$

### F.1.2 An Example When Clusters Overlap

Here, we also deal with an example where the clusters in $\mathbb{X}$ have some overlap. In the following proposition, for the sake of simplicity, we consider the case that there are two clusters in $\mathbb{X}$.

**Proposition 4.** *Let $r > 0$, and $v_1, v_2 \in \mathbb{R}^p$. For each $i \in \{1, 2\}$, let $\mathbb{B}(v_i; r) \subset \mathbb{R}^p$ be the open ball of radius $r$ centered at point $v_i$. Suppose that $\|v_1 - v_2\|_2 = 3r$. Define $\overline{\mathbb{X}} = \mathbb{B}(v_1; r) \cup \mathbb{B}(v_2; r)$, $\mathbb{X} = \mathbb{B}(v_1; 2r) \cup \mathbb{B}(v_2; 2r)$, and $a(x|\overline{x}) = \text{vol}(\mathbb{B}(v_1; 2r))^{-1} \sum_{i=1}^2 \mathbb{1}_{\mathbb{B}(v_i; 2r) \times \mathbb{B}(v_i; r)}(x, \overline{x})$. Let $p_{\overline{\mathbb{X}}}(\overline{x}) := (2 \cdot \text{vol}(\mathbb{B}(v_1; r)))^{-1}$ be a probability density function of $P_{\overline{\mathbb{X}}}$. Define $y : \mathbb{X} \to \{1, 2\}$ as $y(x) = 1$ if $x \in \mathbb{B}(v_1; 2r)$ and $y(x) = 2$ if $x \in \mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r)$. Then, we have the following results:*

1. *$w(x) > 0$ for every $x \in \mathbb{X}$.*

2. *$\text{sim}(x, x'; \lambda) = 2 - \lambda$ if $x, x' \in \mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r)$ or $x, x' \in \mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r)$, $\text{sim}(x, x'; \lambda) = -\lambda$ if $(x, x') \in (\mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r)) \times (\mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r))$ or $(x, x') \in (\mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r)) \times (\mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r))$, and $\text{sim}(x, x'; \lambda) = 1 - \lambda$ otherwise.*

3. *Let $\delta \in (-\lambda, 1 - \lambda]$. Then, $\delta$, $K$, $\mathbb{B}_1, \cdots, \mathbb{B}_K$, and $y$ satisfy Assumption 2.*

*Proof.* Let $\overline{\mathbb{X}}_1$ (resp. $\overline{\mathbb{X}}_2$) denote $\mathbb{B}(v_1; r)$, (resp. $\mathbb{B}(v_2; r)$). Then,

$$w(x) = \mathbb{E}\left[a(x|\overline{x})\right]$$

$$= \int_{\overline{\mathbb{X}}} a(x|\overline{x}) p_{\overline{\mathbb{X}}}(\overline{x}) d\overline{x}$$

$$= \int_{\overline{\mathbb{X}}_1} a(x|\overline{x}) p_{\overline{\mathbb{X}}}(\overline{x}) d\overline{x} + \int_{\overline{\mathbb{X}}_2} a(x|\overline{x}) p_{\overline{\mathbb{X}}}(\overline{x}) d\overline{x}$$

$$= p_{\overline{\mathbb{X}}}(\overline{x}) \left\{ \int_{\overline{\mathbb{X}}_1} a(x|\overline{x}) d\overline{x} + \int_{\overline{\mathbb{X}}_2} a(x|\overline{x}) d\overline{x} \right\} \qquad (p_{\overline{\mathbb{X}}} \text{ is a constant function})$$

Here, we consider Case 1 and Case 2. Firstly, Case 1 is when either $x \in \mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r)$ or $x \in \mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r)$ holds. Since in this case, it is sufficient to prove for the case that $x \in \mathbb{B}(v_1; 2r) \setminus$

$\mathbb{B}(v_2; 2r)$ holds, we may assume this condition. Then, $\int_{\overline{\mathbb{X}}_1} a(x|\overline{x})d\overline{x} = \mathrm{vol}(\mathbb{B}(v_1; r))\mathrm{vol}(\mathbb{B}(v_1; 2r))^{-1}$ and $\int_{\overline{\mathbb{X}}_2} a(x|\overline{x})d\overline{x} = 0$. Thus, $w(x) = 1/(2\mathrm{vol}(\mathbb{B}(v_1; 2r)))$. Secondly, Case 2 is when $x \in \mathbb{B}(v_1; 2r) \cap \mathbb{B}(v_2; 2r)$. Then, $\int_{\overline{\mathbb{X}}_1} a(x|\overline{x})d\overline{x} = \int_{\overline{\mathbb{X}}_2} a(x|\overline{x})d\overline{x} = \mathrm{vol}(\mathbb{B}(v_1; r))\mathrm{vol}(\mathbb{B}(v_1; 2r))^{-1}$. Thus, $w(x) = 1/\mathrm{vol}(\mathbb{B}(v_1; 2r))$. Since $r > 0$, it implies $\mathrm{vol}(\mathbb{B}(v_1; 2r)) > 0$. Thus $w(x) > 0$ for both cases.

Next, we compute

$$\begin{aligned} w(x, x') &= \mathbb{E}_{\overline{x}}\left[a(x|\overline{x})a(x'|\overline{x})\right] \\ &= \int_{\overline{\mathbb{X}}} a(x|\overline{x})a(x'|\overline{x})p_{\overline{\mathbb{X}}}(\overline{x})d\overline{x} \\ &= \int_{\overline{\mathbb{X}}_1} a(x|\overline{x})a(x'|\overline{x})p_{\overline{\mathbb{X}}}(\overline{x})d\overline{x} + \int_{\overline{\mathbb{X}}_2} a(x|\overline{x})a(x'|\overline{x})p_{\overline{\mathbb{X}}}(\overline{x})d\overline{x} \\ &= p_{\overline{\mathbb{X}}}(\overline{x})\left\{\int_{\overline{\mathbb{X}}_1} a(x|\overline{x})a(x'|\overline{x})d\overline{x} + \int_{\overline{\mathbb{X}}_2} a(x|\overline{x})a(x'|\overline{x})d\overline{x}\right\}. \quad (p_{\overline{\mathbb{X}}} \text{ is a constant function}) \end{aligned}$$

Here, we consider Case A, Case B, Case C, and Case D. Firstly Case A is that both $x$ and $x'$ belong to $\mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r)$ (note that the computation for the case that both $x$ and $x'$ belong to $\mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r)$ is the same). Then, $\int_{\overline{\mathbb{X}}_1} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = \mathrm{vol}(\mathbb{B}(v_1; r))/\mathrm{vol}(\mathbb{B}(v_1; 2r))^2$ and $\int_{\overline{\mathbb{X}}_2} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = 0$. Hence, $w(x, x') = \{2(\mathrm{vol}(\mathbb{B}(v_1; 2r)))^2\}^{-1}$. Here recall that $w(x) = w(x') = 1/(2\mathrm{vol}(\mathbb{B}(v_1; 2r)))$, then we have $\mathrm{sim}(x, x'; \lambda) = 2 - \lambda$. Secondly Case B is that $x \in \mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r)$ and $x' \in \mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r)$ (the calculation for the case that $x \in \mathbb{B}(v_2; 2r) \setminus \mathbb{B}(v_1; 2r)$ and $x' \in \mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r)$ is the same). Then, $\int_{\overline{\mathbb{X}}_1} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = \int_{\overline{\mathbb{X}}_2} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = 0$. Therefore, $\mathrm{sim}(x, x'; \lambda) = 0 - \lambda = -\lambda$. Thirdly Case C is that both $x$ and $x'$ belong to $\mathbb{B}(v_1; 2r) \cap \mathbb{B}(v_2; 2r)$. Then, $\int_{\overline{\mathbb{X}}_1} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = \int_{\overline{\mathbb{X}}_2} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = \mathrm{vol}(\mathbb{B}(v_1; r))/\mathrm{vol}(\mathbb{B}(v_1; 2r))^2$. Since $w(x) = w(x') = 1/\mathrm{vol}(\mathbb{B}(v_1; 2r))$, $\mathrm{sim}(x, x'; \lambda) = 1 - \lambda$. Finally in Case D, consider the complementary of the union of the other cases. From the setting, we may assume that $x$ belongs to $\mathbb{B}(v_1; 2r) \cap \mathbb{B}(v_2; 2r)$ and $x'$ to $\mathbb{B}(v_1; 2r) \setminus \mathbb{B}(v_2; 2r)$. Then, $\int_{\overline{\mathbb{X}}_1} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = \mathrm{vol}(\mathbb{B}(v_1; r))/\mathrm{vol}(\mathbb{B}(v_1; 2r))^2$ and $\int_{\overline{\mathbb{X}}_2} a(x|\overline{x})a(x'|\overline{x})d\overline{x} = 0$. Since $w(x) = 1/\mathrm{vol}(\mathbb{B}(v_1; 2r))$ and $w(x') = 1/(2\mathrm{vol}(\mathbb{B}(v_1; 2r)))$, we have $\mathrm{sim}(x, x'; \lambda) = 1 - \lambda$. As a result,

$$\mathrm{sim}(x, x'; \lambda) = \begin{cases} 2 - \lambda, & \text{if Case A holds,} \\ -\lambda, & \text{if Case B holds,} \\ 1 - \lambda, & \text{if Case C holds,} \\ 1 - \lambda, & \text{if Case D holds.} \end{cases}$$

Finally, take $\delta \in (-\lambda, 1 - \lambda]$. Then, from the computation for $\mathrm{sim}(x, x'; \lambda)$ above, the conditions in Assumption 2 are satisfied. $\square$

## F.2 SSL-HSIC Revisit

Li et al. (2021) propose the framework termed SSL-HSIC, which is defined using the notion Hilbert-Schmidt Independence Criterion (HSIC, (Smola et al., 2007)). They show that under some conditions, for a random variable $Z$ (resp. $Y$) that represents the feature vector (resp. the label), one obtains

$$\mathrm{HSIC}(Z, Y) = c\left(\mathbb{E}_{x,x^+}\left[k(f(x), f(x^+))\right] - \mathbb{E}_{x,x^-}\left[k(f(x), f(x^-))\right]\right),$$

where $c > 0$. Li et al. (2021) define the loss of SSL-HSIC as,

$$L_{\text{SSL-HSIC}}(f; \kappa) = -\text{HSIC}(Z, Y) + \kappa\sqrt{\text{HSIC}(Z, Z)},$$

where $\kappa \in \mathbb{R}$.

In the case that $\kappa > 0$, we have

$$L_{\text{KCL}}(f; 1) \lesssim L_{\text{SSL-HSIC}}(f; \kappa).$$

## F.3 Supplementary Information of Section 3

### F.3.1 Relations to Variants of InfoNCE

We first define variants of InfoNCE (Chen et al., 2020a; van den Oord et al., 2018):

- Decoupled InfoNCE loss, which is a variant of the decoupled NT-Xent loss of Chen et al. (2021):

$$\widetilde{L}_{\text{NCE}}(f; \tau, \lambda) = -\mathbb{E}_{x,x^+}\left[\frac{f(x)^\top f(x^+)}{\tau}\right] + \lambda \mathbb{E}_{\substack{x,x^+ \\ \{x_i^-\}}}\left[\log\left(e^{\frac{f(x)^\top f(x^+)}{\tau}} + \sum_{i=1}^M e^{\frac{f(x)^\top f(x_i^-)}{\tau}}\right)\right].$$

- Asymptotic of contrastive loss (Wang and Isola, 2020) decoupled by following the way of Chen et al. (2021):

$$\widetilde{L}_{\infty\text{-NCE}}(f; \tau, \lambda) = -\mathbb{E}_{x,x^+}\left[\frac{f(x)^\top f(x^+)}{\tau}\right] + \lambda \mathbb{E}_x\left[\log \mathbb{E}_{x'}\left[e^{\frac{f(x)^\top f(x')}{\tau}}\right]\right],$$

- InfoNCE loss as a variant of decoupled contrastive learning loss (Yeh et al., 2022):

$$\widetilde{L}_{\text{NCE}}(f; \tau, 1) = -\mathbb{E}_{x,x^+}\left[\frac{f(x)^\top f(x^+)}{\tau}\right] + \mathbb{E}_{x,\{x_i^-\}}\left[\log\left(\sum_{i=1}^M e^{\frac{f(x)^\top f(x_i^-)}{\tau}}\right)\right].$$

- InfoNCE loss as a variant of decoupled contrastive learning loss with additional weight parameter, following Chen et al. (2021):

$$\widetilde{L}_{\text{NCE}}(f; \tau, \lambda) = -\mathbb{E}_{x,x^+}\left[\frac{f(x)^\top f(x^+)}{\tau}\right] + \lambda \mathbb{E}_{x,\{x_i^-\}}\left[\log\left(\sum_{i=1}^M e^{\frac{f(x)^\top f(x_i^-)}{\tau}}\right)\right].$$

Note that $L_{\text{NCE}}(f; \tau)$ and $L_{\infty\text{-NCE}}(f; \tau)$ in Section 2 coincide with $\widetilde{L}_{\text{NCE}}(f; \tau, 1)$ and $\widetilde{L}_{\infty\text{-NCE}}(f; \tau, 1)$ in this subsection, respectively. We show the following facts:

**Proposition 5.** *The following relations hold:*

$$\tau^{-1} L_{\text{LinKCL}}(f; \lambda) \leq \widetilde{L}_{\text{NCE}}(f; \tau, \lambda) + \lambda \log M^{-1}, \tag{29}$$

$$\tau^{-1} L_{\text{LinKCL}}(f; \lambda) \leq \widetilde{L}_{\infty\text{-NCE}}(f; \tau, \lambda), \tag{30}$$

$$\tau^{-1} L_{\text{LinKCL}}(f; \lambda) \leq \widetilde{L}_{\text{NCE}}(f; \tau, \lambda) + \lambda \log M^{-1}. \tag{31}$$

*Proof.* From the definition of $L_{\mathrm{NCE}}(f;\tau,\lambda)$, we have

$$\widetilde{L}_{\mathrm{NCE}}(f;\tau,\lambda) + \lambda\log\frac{1}{M}$$

$$= -\tau^{-1}\mathbb{E}_{x,x^+}\left[f(x)^\top f(x^+)\right] + \lambda\mathbb{E}_{x,x^+,\{x_i^-\}}\left[\log\left(\frac{1}{M}e^{f(x)^\top f(x^+)/\tau} + \frac{1}{M}\sum_{i=1}^M e^{f(x)^\top f(x_i^-)/\tau}\right)\right]$$

$$\geq -\tau^{-1}\mathbb{E}_{x,x^+}\left[f(x)^\top f(x^+)\right] + \lambda\mathbb{E}_{x,\{x_i^-\}}\left[\log\left(\frac{1}{M}\sum_{i=1}^M e^{f(x)^\top f(x_i^-)/\tau}\right)\right]$$

$$\geq -\tau^{-1}\mathbb{E}_{x,x^+}\left[f(x)^\top f(x^+)\right] + \tau^{-1}\lambda\mathbb{E}_{x,\{x_i^-\}}\left[\frac{1}{M}\sum_{i=1}^M f(x)^\top f(x_i^-)\right]$$

$$= \tau^{-1}L_{\mathrm{LinKCL}}(f;\lambda),$$

where in the first inequality we use the fact that $M^{-1}e^{f(x)^\top f(x^+)/\tau} \geq 0$ for any $x,x^+ \in \mathbb{X}$, and in the second inequality we use Jensen's inequality. Note that when $\lambda = 1$, we obtain (1).

The proofs of (31) are almost the same as the proof of (29). The equation (30) is obtained by applying Jensen's inequality. $\qquad\square$

### F.3.2 Relations to SCL

Let us define the quadratic kernel contrastive loss as:

$$L_{\mathrm{QKCL}}(f;\lambda) = -\mathbb{E}_{x,x^+}\left[\left(f(x)^\top f(x^+)\right)^2\right] + \lambda\mathbb{E}_{x,x^-}\left[\left(f(x)^\top f(x^-)\right)^2\right].$$

The spectral contrastive loss $L_{\mathrm{SCL}}(f)$ (HaoChen et al., 2021) is defined as,

$$L_{\mathrm{SCL}}(f) = -2\mathbb{E}_{x,x^+}[f(x)^\top f(x^+)] + \mathbb{E}_{x,x^-}[(f(x)^\top f(x^-))^2]. \tag{32}$$

The following proposition is an elementary result.

**Proposition 6.** *We have,*

$$L_{\mathrm{QKCL}}(f;2^{-1}) \leq \frac{1}{2}L_{\mathrm{SCL}}(f) + \frac{1}{4}.$$

*Proof.* Since $t^2 + 1/4 \geq t$ for every $t \in \mathbb{R}$, we obtain the claim. $\qquad\square$

## F.4 Comparison of Assumption 2 of our work to Assumption 3 in HaoChen and Ma (2023)

Let $\mathbb{M}$ be a measurable subset of $\mathbb{X}$, and let $g : \mathbb{X} \to \mathbb{R}$ be a function. HaoChen and Ma (2023) introduce the following notion that quantifies the inner-connectivity of clusters (see (4) in HaoChen and Ma (2023)):

$$Q_{\mathbb{M}}(g) := \frac{\mathbb{E}_{x,x^+}[(g(x)-g(x^+))^2|\mathbb{M}\times\mathbb{M}]}{\mathbb{E}_{x,x^-}[(g(x)-g(x^-))^2|\mathbb{M}\times\mathbb{M}]}.$$

Here, the expectations above are defined as

$$\mathbb{E}_{x,x^+}[(g(x) - g(x^+))^2 | \mathbb{M} \times \mathbb{M}] = \int_{\mathbb{X} \times \mathbb{X}} (g(x) - g(x'))^2 P_+(dx, dx' | \mathbb{M} \times \mathbb{M}),$$

$$\mathbb{E}_{x,x^-}[(g(x) - g(x^-))^2 | \mathbb{M} \times \mathbb{M}] = \int_{\mathbb{X}} \int_{\mathbb{X}} (g(x) - g(x'))^2 P_{\mathbb{X}}(dx | \mathbb{M}) P_{\mathbb{X}}(dx' | \mathbb{M}),$$

where we use the notation $dP_+ = w(x, x') d\nu_{\mathbb{X}}^{\otimes 2}$. We focus on the following notion, where HaoChen and Ma (2023) denote their subsets by $\{S_1, \cdots, S_m\}$.

**Assumption 4** ($\mathcal{F}$-implementable inner-cluster connection larger than $\beta$, quoted from Assumption 3 in HaoChen and Ma (2023)). *For any function $f \in \mathcal{F}$ and any linear head $w \in \mathbb{R}^k$, let function $g(x) = w^\top f(x)$. For any $i \in [m]$ we have that:*

$$Q_{S_i}(g) \geq \beta.$$

In summary, the relation between Assumption 3 of HaoChen and Ma (2023) and Assumption 2 of our work is given below:

**Proposition 7.** *Suppose that Assumption 2 holds. Take $\delta \in \mathbb{R}$, $K \in \mathbb{N}$, and $\mathbb{M}_1, \cdots, \mathbb{M}_K$ such that the conditions (**A**) and (**B**) are satisfied. Suppose also that: 1) there exists some $c > 0$ such that for every $i \in [K]$, $c \cdot P_+(\mathbb{M}_i \times \mathbb{M}_i) \leq P_{\mathbb{X}}(\mathbb{M}_i)^2$ holds; 2) $\delta + \lambda \geq 0$ holds. Then, the function class $\widetilde{\mathcal{F}}$ including all the maps from $\mathbb{X}$ to $\mathbb{S}^{d-1}$ satisfies Assumption 3 in HaoChen and Ma (2023).*

*Proof.* Take an arbitrary $f \in \widetilde{\mathcal{F}}$ and $w \in \mathbb{R}^d$. For each $i \in [K]$, for any $x, x' \in \mathbb{M}_i$ we have that $w(x, x') \geq (\delta + \lambda)w(x)w(x')$. Since $\delta + \lambda \geq 0$,

$$\int_{\mathbb{M}_i \times \mathbb{M}_i} (g(x) - g(x'))^2 w(x, x') \nu_{\mathbb{X}}^{\otimes 2}(dx, dx') \geq (\delta + \lambda) \int_{\mathbb{M}_i \times \mathbb{M}_i} (g(x) - g(x'))^2 w(x) w(x') \nu_{\mathbb{X}}^{\otimes 2}(dx, dx').$$

Here, using $c \cdot P_+(\mathbb{M}_i \times \mathbb{M}_i) \leq P_{\mathbb{X}}(\mathbb{M}_i)^2$, we have

$$\frac{1}{P_+(\mathbb{M}_i \times \mathbb{M}_i)} \int_{\mathbb{M}_i \times \mathbb{M}_i} (g(x) - g(x'))^2 w(x, x') \nu_{\mathbb{X}}^{\otimes 2}(dx, dx')$$

$$\geq c(\delta + \lambda) \frac{1}{P_{\mathbb{X}}(\mathbb{M}_i)^2} \int_{\mathbb{M}_i \times \mathbb{M}_i} (g(x) - g(x'))^2 w(x) w(x') \nu_{\mathbb{X}}^{\otimes 2}(dx, dx').$$

The above inequality means that,

$$\int_{\mathbb{X} \times \mathbb{X}} (g(x) - g(x'))^2 P_+(dx, dx' | \mathbb{M}_i \times \mathbb{M}_i) \geq c(\delta + \lambda) \int_{\mathbb{X}} \int_{\mathbb{X}} (g(x) - g(x'))^2 P_{\mathbb{X}}(dx | \mathbb{M}_i) P_{\mathbb{X}}(dx' | \mathbb{M}_i).$$

Thus, we obtain $Q_{\mathbb{M}_i}(g) \geq c(\delta + \lambda)$, where $g(x) = w^\top f(x)$. $\qquad \square$

The above proposition indicates that the inner-connectivity $Q_{\mathbb{M}_i}(g)$ for any $i \in [K]$ and $g(x) = w^\top f(x)$ with $f \in \widetilde{\mathcal{F}}$ is lower bounded by $c(\delta + \lambda)$ under the assumptions. Therefore, Assumption 2 of our work is a sufficient condition of Assumption 3 in HaoChen and Ma (2023) if $c \cdot P_+(\mathbb{M}_i \times \mathbb{M}_i) \leq P_{\mathbb{X}}(\mathbb{M}_i)^2$ (for every $i \in [K]$) and $\delta + \lambda \geq 0$ hold.

46

**Remark 4.** *We can construct a positive value c in the above statement explicitly. In this remark, we show a simple way to do so. Let $X, Y$ be random variables on a probability space $(\Omega, P_\Omega)$ with the joint probability distribution $P_+$ and the marginal distribution $P_\mathbb{X}$. Denote $p_1 = P_\Omega(X \in \mathbb{M}_i)$ and $p = P_\Omega(X \in \mathbb{M}_i, Y \in \mathbb{M}_i)$. Here, let $V$ be the covariance matrix of the random variables $\mathbb{1}_{\{X \in \mathbb{M}_i\}}$ and $\mathbb{1}_{\{Y \in \mathbb{M}_i\}}$. The positive semi-definiteness of $V$ implies the inequality $(p_1 - p_1^2)^2 - (p - p_1^2)^2 \geq 0$. This inequality is valid when $2p_1^2 - p_1 \leq p \leq p_1$. Combining this fact with the property $p \geq 0$, we obtain*

$$\max\{2p_1^2 - p_1, 0\} \leq p \leq p_1,$$

*which implies,*

$$P_\mathbb{X}(\mathbb{M}_i) \cdot \max\{2P_\mathbb{X}(\mathbb{M}_i) - 1, 0\} \cdot P_+(\mathbb{M}_i \times \mathbb{M}_i) \leq P_\mathbb{X}(\mathbb{M}_i)^2.$$

*Thus, if $P_\mathbb{X}(\mathbb{M}_i) > 1/2$ holds for every $i \in [K]$, then we can take*

$$c = \min_{i \in [K]}\{P_\mathbb{X}(\mathbb{M}_i)(2P_\mathbb{X}(\mathbb{M}_i) - 1)\}.$$

## F.5    More Discussion about Generalization Bounds in Section 5.3

In this section, we discuss the differences between our generalization error bound and the results presented by other works on contrastive learning. We summarize the differences below.

- Arora et al. (2019); Ash et al. (2022); Lei et al. (2023); Zou and Liu (2023) consider the case that for a pair $(x, x^+)$, $M$-tuple samples $(x_1^-, \cdots, x_M^-)$ independent from other random variables are available. Thus, our problem setup is different from them. Especially, in our analysis, it is also necessary to tackle the cases in which $X_i, X_i'$, $i \in [n]$, are not necessarily independent and the standard techniques (e.g., see Mohri et al. (2018)) cannot be applied. We instead utilized the results of McDiarmid's inequality for dependent random variables shown by Zhang et al. (2019).

- The empirical loss considered in Zhang et al. (2022) is defined in a different way from our empirical kernel contrastive loss. Also, our proof technique is different from Zhang et al. (2022).

- HaoChen et al. (2021); Nozawa et al. (2020) consider the case in which the augmented samples are not necessarily independent. Nozawa et al. (2020) utilize the theory on PAC-Bayes bounds (Guedj, 2019), and HaoChen et al. (2021) provide the high probability bound. Our analysis is different from Nozawa et al. (2020) since our analysis is based on several concentration inequalities. HaoChen et al. (2021) consider the empirical spectral contrastive loss that is defined by raw samples and expressed in the expectation w.r.t. the augmented samples that are drawn according to the conditional distribution given the raw samples (see Section 4.1 in HaoChen et al. (2021)). On the other hand, we derive a generalization error bound for the empirical kernel contrastive loss defined by only augmented samples.

- Wang et al. (2022b) establish the generalization error bound for the spectral contrastive loss (HaoChen et al., 2021), where their analysis improves the convergence rate of HaoChen

et al. (2021). In their analysis, they decompose the second term of the spectral contrastive loss in a different way from us (for the detail of their decomposition, see the proof of Proposition D.1 of Wang et al. (2022b)). Also, they utilize the concentration inequality shown by Clémençon et al. (2008) (see equation (51) in Wang et al. (2022b)), while we use the results proved by Zhang et al. (2019). Thus, the techniques we use in the proof of Theorem 6 are different from those of Wang et al. (2022b).

## F.6   Detailed Comparison to Robinson et al. (2021b)

Robinson et al. (2021b) tackle the hard negative sampling problem in contrastive learning from both the theoretical and empirical perspectives. They also establish generalization bounds for their hard negative objectives by introducing the 1-NN classifier (for their definition of the 1-NN classifier, see the statement of Theorem 5 in Robinson et al. (2021b)). We give a detailed comparison between our results and the theoretical analysis by Robinson et al. (2021b). The main differences are listed below:

- The problem setup of the theoretical results by Robinson et al. (2021b) is based on that of Arora et al. (2019), i.e., they rely on the conditional independence assumption. On the other hand, we does not rely on it, where we utilize the similarity function $\text{sim}(\cdot, \cdot; \lambda)$ instead.

- In the proof of Theorem 5 of Robinson et al. (2021b) (see also Theorem 8 and the proof of their work), the supervised loss is upper-bounded by the term $\mathbb{E}_c \mathbb{E}_{x,x^+ \sim_{iid} p(\cdot|c)} \|f(x) - f(x^+)\|^2$. In summary, the differences between Theorem 5 of Robinson et al. (2021b) and Theorem 2 of our work are: (i) In the numerator of the upper bound in the proof of Theorem 8 of Robinson et al. (2021b), the term mentioned above appears. On the other hand, in Theorem 2 of our work, the quantity $\mathfrak{a}(f)$ appears. (ii) Our upper bound includes the quantity $\Delta_{\min}(f)$. (iii) We also note that the proof techniques used in Theorem 2 in our work are different from Robinson et al. (2021b).

- Note that the label employed in the analysis by Robinson et al. (2021b) is a random variable, while our analysis employ the deterministic labeling function.

## F.7   Detailed Comparison to Huang et al. (2023); Zhao et al. (2023)

Huang et al. (2023) present the generalization bounds that utilizes the 1-NN classifier (for the definition of the 1-NN classifier introduced in Huang et al. (2023), see Section 2 in their paper). Besides, Zhao et al. (2023) extend the results of Huang et al. (2023). Thus, it is worth discussing the differences between the results by Huang et al. (2023); Zhao et al. (2023) and our Theorem 2. We summarize the differences below:

- Huang et al. (2023) show that if the centers of clusters in the feature space are sufficiently apart from each other (note that they call it *divergence*), then their supervised error function is upper bounded by the *alignment* term up to several constants and parameters. Hence, their results do not show that the *divergence* relates directly to the supervised error, i.e., the *divergence* term does not appear in their upper bounds of the supervised error. On the other

hand, we show that the quantities related to the *divergence* in the RKHS can also contribute to upper-bounding the supervised error (see Theorem 2).

- In Huang et al. (2023); Zhao et al. (2023), it is little investigated to what range of encoder models their results can apply. On the other hand, our Theorem 2 requires only the meaningfulness (Definition 2) of encoders belonging to $\mathcal{F}$. Especially, suppose $k$ is the linear kernel, then Theorem 2 in our study refines Theorem 1 of Huang et al. (2023) in this sense.

- Huang et al. (2023); Zhao et al. (2023) utilize the notion termed $(\sigma, \delta)$-*augmentation*, while our analysis utilizes Assumption 2 based on the definition of the similarity function $\mathrm{sim}(\cdot, \cdot; \lambda)$.

- Huang et al. (2023); Zhao et al. (2023) often use the assumption that the encoder $f$ is a Lipschitz function. Meanwhile, our main result does not require that $f$ should be a Lipschitz function.

- Zhao et al. (2023) consider the squared loss for the downstream classification task (see Theorem 3.2 in their paper). On the other hand, we consider the classification error.

# G  Connections between KCL and Normalized Cut

In this section, we present supplementary information of Section 4.1. Throughout this section, we assume that $\inf_{x \in \mathbb{X}} w(x) < \infty$ and $\sup_{\overline{x} \in \overline{\mathbb{X}}} \sup_{x \in \mathbb{X}} a(x|\overline{x}) < \infty$ hold. Note that the assumption $\sup_{\overline{x} \in \overline{\mathbb{X}}} \sup_{x \in \mathbb{X}} a(x|\overline{x}) < \infty$ implies $\sup_{x \in \mathbb{X}} w(x) < \infty$.

## G.1  The Problem Setup of Normalized Cut

In this section, we first explain the population-level normalized cut problem based on Shi and Malik (2000); Terada and Yamamoto (2019); Von Luxburg (2007). Suppose that there are total $K$ clusters in $\mathbb{X}$. Following Terada and Yamamoto (2019), the optimization problem of the population-level normalized cut is given as:

$$\min_{\mathbb{V}_1, \cdots, \mathbb{V}_K} \sum_{i=1}^{K} \frac{W(\mathbb{V}_i, \mathbb{V}_i^c)}{\mathrm{vol}(\mathbb{V}_i)} \tag{33}$$

where the minimum in the above problem is taken over all the possible combinations of $K$ disjoint non-empty measurable subsets $\mathbb{V}_1, \cdots, \mathbb{V}_K$ satisfying $\bigcup_{i=1}^{K} \mathbb{V}_i = \mathbb{X}$, and $W$ and $\mathrm{vol}(\cdot)$ are defined as,

$$W(\mathbb{V}_i, \mathbb{V}_i^c) = \int_{(x,x') \in \mathbb{V}_i \times \mathbb{V}_i^c} \mathrm{sim}(x, x'; \lambda) w(x) w(x') d\nu_{\mathbb{X}}^{\otimes 2}(x, x'), \tag{34}$$

$$\mathrm{vol}(\mathbb{V}_i) = \int_{\mathbb{V}_i} w(x) d\nu_{\mathbb{X}}(x), \tag{35}$$

where $\nu_{\mathbb{X}}^{\otimes 2} := \nu_{\mathbb{X}} \otimes \nu_{\mathbb{X}}$ is the product measure. Here also note that $\mathrm{vol}(\cdot)$ is the volume of a set $\mathbb{V}_i$. Terada and Yamamoto (2019) consider the case that a reproducing kernel is used as similarity

measurement: see Theorem 7 in Terada and Yamamoto (2019). Note that some existing work deals with the measurable partition problems such as the ratio cut and Cheeger cut (Trillos et al., 2016).

Denote by $L^2(\mathbb{X}, P_{\mathbb{X}})$, the Hilbert space over the field $\mathbb{R}$ consisting of real-valued and squared-integrable function defined on $\mathbb{X}$ for $P_{\mathbb{X}}$-a.e., with its inner product $\langle f, g \rangle_{L^2(\mathbb{X}, P_{\mathbb{X}})} = \int f(x)g(x)dP_{\mathbb{X}}(x)$. Let $U : \mathbb{R}^K \to L^2(\mathbb{X}, P_{\mathbb{X}})$ be a linear operator defined as,

$$(Uz)(\cdot) = \sum_{i=1}^{K} \frac{\mathbb{1}_{\mathbb{V}_i}(\cdot)}{\sqrt{\mathrm{vol}(\mathbb{V}_i)}} z_i, \quad z = (z_1, \cdots, z_K)^\top, \tag{36}$$

where $\mathbb{1}_{\mathbb{V}_i}(x) = 1$ if $x \in \mathbb{V}_i$ and $0$ if $x \notin \mathbb{V}_i$, and $z_i := \langle z, e_i \rangle_{\mathbb{R}^K}$ for each $i \in [K]$ with an orthonormal basis $\{e_i\}_{i=1}^K$ of $\mathbb{R}^K$. Note that under the setting that $|\mathbb{X}| < \infty$, the linear operator $U$ is equal to $(\mathbb{1}_{\mathbb{V}_j}(x_i)/\sqrt{\mathrm{vol}(\mathbb{V}_j)})_{ij}$. Therefore, the definition of $U$ matches that of the classical theory of normalized cut (Shi and Malik, 2000). Moreover, every augmented data $x \in \mathbb{X}$ belongs to one of the subsets $\mathbb{V}_1, \cdots, \mathbb{V}_K$. Since $\mathbb{V}_1, \cdots, \mathbb{V}_k$ are disjoint, linear operator $U$ is bounded, and the adjoint operator $U^\dagger$ exists uniquely. Here, Von Luxburg (2007) explain that the objective function of the normalized cut problem can be rewritten as a combinatorial optimization problem. Applying the arguments presented by Von Luxburg (2007) to our setup, we have

$$\sum_{i=1}^{K} \frac{W(\mathbb{V}_i, \mathbb{V}_i^c)}{\mathrm{vol}(\mathbb{V}_i)} = -\mathrm{Tr}(U^\dagger A U) + (1 - \lambda)K, \tag{37}$$

where $A : L^2(\mathbb{X}, P_{\mathbb{X}}) \to L^2(\mathbb{X}, P_{\mathbb{X}})$ is a Hilbert-Schmidt integral operator defined as

$$A\psi(\cdot) = \int \mathrm{sim}(\cdot, x; \lambda)\psi(x)w(x)d\nu_{\mathbb{X}}(x) \quad \psi \in L^2(\mathbb{X}, P_{\mathbb{X}}), \tag{38}$$

and $-\mathrm{Tr}(U^\dagger A U) = -\sum_{i=1}^{K}\langle U^\dagger A U e_i, e_i \rangle_{\mathbb{R}^K}$. The proof of (37) closely follows that of Von Luxburg (2007); in Appendix G.3, we present the proof of an extended version. Here the following proposition shows the well-definedness of $A$.

**Proposition 8.** *Suppose the setting described in Section 2.1 holds*, $\inf_{x \in \mathbb{X}} w(x) > 0$, *and* $\sup_{\overline{x} \in \overline{\mathbb{X}}} \sup_{x \in \mathbb{X}} a(x|\overline{x}) < \infty$ *holds. Then, the integral operator $A$ is well-defined.*

*Proof.* We can evaluate

$$\left| \int \mathrm{sim}^2(x, x; \lambda)w(x)d\nu_{\mathbb{X}}(x) \right| = \left| \int \left( \frac{w(x, x)}{w(x)w(x)} - \lambda \right)^2 w(x)d\nu_{\mathbb{X}}(x) \right|$$

$$\leq \int \left( \left( \frac{w(x, x)}{w(x)w(x)} \right)^2 + 2\lambda \frac{w(x, x)}{w(x)w(x)} + \lambda^2 \right) w(x)d\nu_{\mathbb{X}}(x)$$

$$< +\infty,$$

where we use the assumptions that $\inf_{x \in \mathbb{X}} w(x) > 0$, $\sup_{\overline{x} \in \overline{\mathbb{X}}} \sup_{x \in \mathbb{X}} a(x|\overline{x}) < \infty$. $\qquad \square$

Suppose the dimension of the RKHS $\mathcal{H}_k$ associated with the kernel function $k$ is greater than or equal to $K$. In the following section, it is convenient to redefine (33) as,

$$\min_{\mathbb{V}_1, \cdots, \mathbb{V}_K} \sum_{i=1}^{\infty} \frac{W(\mathbb{V}_i, \mathbb{V}_i^c)}{\mathrm{vol}(\mathbb{V}_i)},$$

where we define $\mathbb{V}_j = \emptyset$ for every $j > K$. Also, let us redefine (36) as the linear operator $U : \mathcal{H}_k \to L^2(\mathbb{X}, P_\mathbb{X})$,

$$(U\psi)(\cdot) = \sum_{i=1}^{\infty} \frac{\mathbb{1}_{\mathbb{V}_i}(\cdot)}{\sqrt{\mathrm{vol}(\mathbb{V}_i)}} \langle \psi, e_i \rangle_{\mathcal{H}_k},$$

where $\{e_j\}_{j=1}^{\infty}$ is an orthonormal basis of $\mathcal{H}_k$ (if $\mathcal{H}_k$ is finite dimensional, then we understand that $\{e_j\}_{j=1}^{\infty}$ consists of finitely many non-zero elements), and we define $\mathbb{1}_{\mathbb{V}_i}(\cdot)/\sqrt{\mathrm{vol}(\mathbb{V}_i)} = 0$ for every $i > K$ as notations. Then, we have the following identity that is analogous of (37):

$$\sum_{i=1}^{\infty} \frac{W(\mathbb{V}_i, \mathbb{V}_i^c)}{\mathrm{vol}(\mathbb{V}_i)} = -\mathrm{Tr}(U^\dagger A U) + (1 - \lambda)K. \tag{39}$$

For the sake of completeness, we provide the proof of the identity (39) in Appendix G.3.

## G.2    Connecting KCL and Normalized Cut via RKHS

Let $k : \mathbb{S}^{d-1} \times \mathbb{S}^{d-1}$ be a continuous, symmetric, and positive-definite kernel function whose RKHS $\mathcal{H}_k$ is $K$-dimensional Hilbert space ($K$ is either finite or $\infty$); For the theory of reproducing kernels, see e.g., Aronszajn (1950); Berlinet and Thomas-Agnan (2004); Steinwart and Christmann (2008). Many kernel functions satisfy these conditions, e.g. the Gaussian kernel, the polynomial kernel, and the linear kernel. Since $\mathbb{S}^{d-1}$ is separable, the RKHS $\mathcal{H}_k$ has an orthonormal basis that is at most countable (e.g., see Berlinet and Thomas-Agnan (2004)). Let $\{e_j\}_{j=1}^{\infty}$ be a countable orthonormal basis of $\mathcal{H}_k$, where $\{e_j\}_{j=1}^{\infty}$ includes only finitely many non-zero elements if $\mathcal{H}_k$ is finite-dimensional. Note that our construction is valid regardless of the choice of $\{e_j\}_{j=1}^{\infty}$. Recall the problem setup presented in Section 2.1. Then, the linear operator $H : \mathcal{H}_k \to L^2(\mathbb{X}, P_\mathbb{X})$ is defined as

$$(H\varphi)(\cdot) = \langle h(f(\cdot)), \varphi \rangle_{\mathcal{H}_k}, \tag{40}$$

for $\varphi \in \mathcal{H}_k$. Let $\| \cdot \|_{\mathcal{H}_k}$ be the norm of the RKHS $\mathcal{H}_k$. Then the following holds for the linear operator $H$ defined in (40):

**Proposition 9.** *The linear operator $H$ is well-defined, i.e., $H\varphi \in L^2(\mathbb{X}, P_\mathbb{X})$ for every $\varphi \in \mathcal{H}_k$. Also, $H$ is continuous, i.e., for a sequence $\varphi_j$ converging strongly to $\varphi$ in $\mathcal{H}_k$, we have that $H\varphi_j$ is convergent to $H\varphi$.*

*Proof.* For any $\varphi \in \mathcal{H}_k$, we have

$$\int |(H\varphi)(x)|^2 w(x) d\nu_\mathbb{X}(dx) = \int |\langle h(f(x)), \varphi \rangle_{\mathcal{H}_k}|^2 w(x) d\nu_\mathbb{X}(x)$$

$$\leq \int \|h(f(x))\|_{\mathcal{H}_k}^2 \|\varphi\|_{\mathcal{H}_k}^2 w(x) d\nu_\mathbb{X}(x)$$

$$\leq \|\varphi\|_{\mathcal{H}_k}^2 \mathbb{E}_x [k(f(x), f(x))] < +\infty.$$

Here in the second inequality we use the Cauchy-Schwarz inequality, and in the last equality we use the fact that $\mathbb{S}^{d-1}$ is compact and $k$ is continuous. Furthermore, if $\varphi_j \to \varphi$ in the sense of strongly

convergence in $\mathcal{H}_k$, then we have

$$\|\langle h(f(\cdot)), \varphi_j\rangle_{\mathcal{H}_k} - \langle h(f(\cdot)), \varphi\rangle_{\mathcal{H}_k}\|_{L^2(\mathbb{X}, P_\mathbb{X})}^2 = \int |\langle h(f(x)), \varphi_j - \varphi\rangle_{\mathcal{H}_k}|^2 w(x) d\nu_\mathbb{X}(x)$$
$$\leq \|\varphi_j - \varphi\|_{\mathcal{H}_k}^2 \mathbb{E}_x\left[k(f(x), f(x))\right]$$
$$\longrightarrow 0 \quad (j \to \infty).$$

Thus $H\varphi_j$ converges to $H\varphi$, and we end the proof. $\qquad\square$

Proposition 9 implies that $H$ is bounded. Therefore, the adjoint operator $H^\dagger : L^2(\mathbb{X}, P_\mathbb{X}) \to \mathcal{H}_k$ exists uniquely.

Now let us recall the definition of the similarity function $\text{sim}(\cdot, \cdot; \lambda)$ with the fixed $\lambda$ in (2), and we consider to relax the combinatorial problem (33) using the linear operator $H$ defined in (40) as follows: we replace the linear operator $U$ in (37) with $H$, which results in the objective function $-\text{Tr}(H^\dagger A H)$. Then, the following proposition holds.

**Proposition 10.** *We have*

$$-\text{Tr}(H^\dagger A H) = -\mathbb{E}_{x,x^+}\left[k(f(x), f(x^+))\right] + \lambda\mathbb{E}_{x,x^-}\left[k(f(x), f(x^-))\right].$$

*Proof.* From the definition of $\text{sim}(x, x'; \lambda)$,

$$(A\psi)(x) = \int \text{sim}(x, x'; \lambda)\psi(x')w(x')d\nu_\mathbb{X}(x')$$
$$= \int \left(\frac{w(x, x')}{w(x)w(x')} - \lambda\right)\psi(x')w(x')d\nu_\mathbb{X}(x')$$
$$= \underbrace{\int \frac{w(x, x')}{w(x)w(x')}\psi(x')w(x')d\nu_\mathbb{X}(x')}_{:=(A_{\text{pos}}\psi)(x)} - \lambda\underbrace{\int \psi(x')w(x')d\nu_\mathbb{X}(x')}_{:=(A_{\text{neg}}\psi)(x)}.$$

Firstly, let us proof the identity

$$\text{Tr}(H^\dagger A_{\text{pos}} H) = \mathbb{E}_{x,x^+}\left[k(f(x), f(x^+))\right].$$

The proof is described as follows: From the definition of $H$,

$$(He_i)(x) = \langle h(f(x)), e_i\rangle_{\mathcal{H}_k} \quad x \in \mathbb{X}.$$

Then we have

$$(A_{\text{pos}}He_i)(x) = \int \frac{w(x, x')}{w(x)w(x')}w(x')\langle h(f(x')), e_i\rangle_{\mathcal{H}_k}d\nu_\mathbb{X}(x')$$
$$= \int \frac{w(x, x')}{w(x)}\langle h(f(x')), e_i\rangle_{\mathcal{H}_k}d\nu_\mathbb{X}(x').$$

Here, the adjoint operator $H^\dagger$ satisfies the following identity; For $\psi \in L^2(\mathbb{X}, P_\mathbb{X})$,

$$\langle He_i, \psi\rangle_{L^2(\mathbb{X}, P_\mathbb{X})} = \langle e_i, H^\dagger\psi\rangle_{\mathcal{H}_k}.$$

Utilizing this relation yields the following representation:

$$H^\dagger A_{\text{pos}} H e_j$$

$$= \sum_{i=1}^{\infty} \langle H^\dagger A_{\text{pos}} H e_j, e_i \rangle_{\mathcal{H}_k} e_i$$

$$= \sum_{i=1}^{\infty} \langle A_{\text{pos}} H e_j, H e_i \rangle_{L^2(\mathbb{X}, P_{\mathbb{X}})} e_i$$

$$= \sum_{i=1}^{\infty} \left( \int w(x) \langle h(f(x)), e_i \rangle_{\mathcal{H}_k} \int \frac{w(x, x')}{w(x)} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} d\nu_{\mathbb{X}}(x') d\nu_{\mathbb{X}}(x) \right) e_i$$

$$= \sum_{i=1}^{\infty} \left( \int \int w(x, x') \langle h(f(x)), e_i \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x') \right) e_i.$$

Therefore,

$$\text{Tr}(H^\dagger A_{\text{pos}} H) = \sum_{j=1}^{\infty} \langle H^\dagger A_{\text{pos}} H e_j, e_j \rangle_{\mathcal{H}_k}$$

$$= \sum_{j=1}^{\infty} \int \int w(x, x') \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$= \int \int w(x, x') \sum_{j=1}^{\infty} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$= \int \int w(x, x') \langle h(f(x)), h(f(x')) \rangle_{\mathcal{H}_k} d\nu_{\mathbb{X}}(x) d\nu_{\mathbb{X}}(x')$$

$$= \mathbb{E}_{x, x^+} \left[ k(f(x), f(x^+)) \right].$$

Note that the third equality above is due to the Dominated Convergence Theorem. Indeed, the sum $\sum_{j=1}^{n} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k}$ converges pointwisely to $\langle h(f(x)), h(f(x)) \rangle_{\mathcal{H}_k}$ on $\mathbb{X} \times \mathbb{X}$, and

$$\left| \sum_{j=1}^{n} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} \right| \leq \sum_{j=1}^{n} \left| \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} \right|$$

$$\leq \left( \sum_{j=1}^{n} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k}^2 \right)^{1/2} \left( \sum_{j=1}^{n} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k}^2 \right)^{1/2}$$

$$\leq \left( \sum_{j=1}^{\infty} \langle h(f(x)), e_j \rangle_{\mathcal{H}_k}^2 \right)^{1/2} \left( \sum_{j=1}^{\infty} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k}^2 \right)^{1/2}$$

$$= \|h(f(x))\|_{\mathcal{H}_k} \|h(f(x'))\|_{\mathcal{H}_k}$$

$$\leq \sup_{x \in \mathbb{X}} k(f(x), f(x)) < +\infty.$$

On the other hand, it is obvious that,

$$
\begin{aligned}
\langle H^\dagger A_{\mathrm{neg}} H e_j, e_j \rangle_{\mathcal{H}_k} &= \langle A_{\mathrm{neg}} H e_j, H e_j \rangle_{L^2(\mathbb{X}, P_\mathbb{X})} \\
&= \int w(x) \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \int w(x') \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} d\nu_\mathbb{X}(x') d\nu_\mathbb{X}(x) \\
&= \int \int \langle h(f(x)), e_j \rangle_{\mathcal{H}_k} \langle h(f(x')), e_j \rangle_{\mathcal{H}_k} w(x) w(x') d\nu_\mathbb{X}(x) d\nu_\mathbb{X}(x').
\end{aligned}
$$

Hence we obtain,

$$
\begin{aligned}
\mathrm{Tr}(H^\dagger A_{\mathrm{neg}} H) &= \sum_{j=1}^{\infty} \langle H^\dagger A_{\mathrm{neg}} H e_j, e_j \rangle_{\mathcal{H}_k} \\
&= \int \int \langle h(f(x)), h(f(x')) \rangle_{\mathcal{H}_k} w(x) w(x') d\nu_\mathbb{X}(x) d\nu_\mathbb{X}(x') \\
&= \mathbb{E}_{x, x^-} \left[ k(f(x), f(x^-)) \right].
\end{aligned}
$$

Hence, we obtain the desired results and end the proof. $\qquad\square$

### G.2.1 Comparison with Related Work from the Graph Cut Viewpoint

HaoChen et al. (2021) has already investigated links between the population-level spectral clustering and contrastive learning. However, our integral kernel (2) introduced in Section 4.1 is slightly different from that of HaoChen et al. (2021), since 1) we divide $w(x, x')$ by $w(x)w(x')$ in the first term rather than by $\sqrt{w(x)w(x')}$ (see Appendix F in HaoChen et al. (2021)), 2) we also incorporate the hyperparemeter $\lambda$.

Note that Tian (2022) introduces a unified framework termed $\alpha$-CL, which connects various contrastive losses from the coordinate-wise optimization perspective. In Tian (2022), the contrastive covariance plays a central role in the theoretical analysis. On the other hand, we use the similarity function defined in Section 4, and thus the approach of our analysis is different from the contrastive covariance of Tian (2022).

### G.3 Proof of (39)

*Proof of* (39). For the proof of (39), we closely follow the approaches presented in Section 5 of Von Luxburg (2007). Since we consider the population-level normalized cut, we present the proof of (39) for the sake of completeness.

Let us define the identity operator $D : L^2(\mathbb{X}, P_\mathbb{X}) \to L^2(\mathbb{X}, P_\mathbb{X})$ as $D\psi = \psi$ for $\psi \in L^2(\mathbb{X}, P_\mathbb{X})$. From the definitions of $D$ and $U$, we have

$$
DU e_i = \begin{cases} \dfrac{\mathbb{1}_{\mathbb{V}_i}(\cdot)}{\sqrt{\mathrm{vol}(\mathbb{V}_i)}} & (i \le K), \\ 0 & (i > K). \end{cases}
$$

Hence, we have the following for $i \le K$:

$$
\langle U^\dagger D U e_i, e_i \rangle_{\mathcal{H}_k} = \int \frac{w(x) \mathbb{1}_{\mathbb{V}_i}(x)^2}{\mathrm{vol}(\mathbb{V}_i)} d\nu_\mathbb{X}(x) = 1.
$$

On the other hand, for $i \leq K$ we have

$$\langle U^{\dagger} A U e_i, e_i \rangle_{\mathcal{H}_k}$$
$$= \int \int (w(x,x') - \lambda w(x)w(x')) \frac{\mathbb{1}_{\mathbb{V}_i}(x)}{\sqrt{\mathrm{vol}(\mathbb{V}_i)}} \frac{\mathbb{1}_{\mathbb{V}_i}(x')}{\sqrt{\mathrm{vol}(\mathbb{V}_i)}} d\nu_{\mathbb{X}}(x') d\nu_{\mathbb{X}}(x)$$

Therefore, we obtain the following:

$$\mathrm{Tr}(U^{\dagger}(D-A)U) = \sum_{i=1}^{\infty} \langle U^{\dagger}(D-A)U e_i, e_i \rangle_{\mathcal{H}_k}$$

$$= \sum_{i=1}^{K} \langle U^{\dagger}(D-A)U e_i, e_i \rangle_{\mathcal{H}_k}$$

$$= \lambda K + \frac{1}{2} \sum_{i=1}^{K} \int \int \left( w(x,x') - \lambda w(x)w(x') \right) \left( \frac{\mathbb{1}_{\mathbb{V}_i}(x)}{\sqrt{\mathrm{vol}(\mathbb{V}_i)}} - \frac{\mathbb{1}_{\mathbb{V}_i}(x')}{\sqrt{\mathrm{vol}(\mathbb{V}_i)}} \right)^2 d\nu_{\mathbb{X}}^{\otimes 2}(x,x')$$

$$= \lambda K + \sum_{i=1}^{K} \int_{x \in \mathbb{V}_i} \int_{x' \in \mathbb{V}_i^c} \frac{(w(x,x') - \lambda w(x)w(x'))}{\mathrm{vol}(\mathbb{V}_i)} d\nu_{\mathbb{X}}(x') d\nu_{\mathbb{X}}(x)$$

$$= \lambda K + \sum_{i=1}^{K} \int_{x \in \mathbb{V}_i} \int_{x' \in \mathbb{V}_i^c} \frac{\mathrm{sim}(x,x';\lambda)}{\mathrm{vol}(\mathbb{V}_i)} w(x)w(x') d\nu_{\mathbb{X}}(x') d\nu_{\mathbb{X}}(x)$$

$$= \lambda K + \sum_{i=1}^{K} \frac{W(\mathbb{V}_i, \mathbb{V}_i^c)}{\mathrm{vol}(\mathbb{V}_i)}.$$

Hence we end the proof. □

# H   Experiments

## H.1   Experimental setup

We provide the setting of the experiments presented in this paper. The code used in our experiments is based on the official implementation of SimSiam[1] and written with PyTorch (Paszke et al., 2019). We basically follow the experimental setting of Chen and He (2021). For the sake of completeness, we provide the detail of the setup used in our experiments. During the stage of pretraining, we construct a trainable encoder model as follows: following Chen and He (2021), we use a backbone architecture whose parameters are initialized, followed by the MLP that consists of linear layers, batch normalization (Ioffe and Szegedy, 2015), and the ReLU activation function. Note that this type of MLP is called *projection head* (Chen et al., 2020a). The output of a trainable encoder model is normalized using the Euclidean norm as several works do (Chen et al., 2020a; Dwibedi et al., 2021). On the other hand, during the stage of linear evaluation (Chen et al., 2020a; Chen and He, 2021), the additional MLP is removed from a trained encoder model, and then a linear classification head is added to the encoder model. The parameters of the trained encoder model

---

[1] https://github.com/facebookresearch/simsiam (Last accessed: March 25, 2023)

are frozen in this stage, and only the linear head is trained. In all of the experiments reported in this paper, we use ResNet-18 (He et al., 2016) as the backbone architecture. We use the 2-layer multi-layer perceptron for the projection head, where the first linear layer is bias-free, and the last linear layer has the bias term.

For the pretraining, we use the same data augmentation techniques as Chen et al. (2020b); Chen and He Chen and He (2021). Note that following Chen and He Chen and He (2021), for the CIFAR-10 experiments, we exclude the Gaussian blur augmentation. For the linear evaluation, we also follow the data augmentation techniques of Chen and He (2021). Note that in both the stage of pretraining and linear evaluation, we set drop_last to True in the training data loader.

For optimization during both the stage of pretraining and linear evaluation, following Chen and He (2021), we use the SGD optimizer. Inspired by HaoChen et al. (2021), we use the cosine-decay learning rate scheduler (Loshchilov and Hutter, 2017) with warmup (Goyal et al., 2017). Note that we use the cosine-decay learning rate scheduler in both the pretraining and linear evaluation and apply warmup in the pretraining. Following the implementation of the learning rate scheduler of the official implementation of SCL[2], we also define our learning rate scheduler by the number of iterations.

### H.1.1 Configurations

In all of the experiments reported in this paper, we use the following configurations:

**Pretraining**   For the learning rate, we set the initial learning rate to 0.0005, the base learning rate to 0.05, and the warmup epochs to 10. Following Chen and He (2021), we use the linear scaling (Goyal et al., 2017) for the learning rate. For the setting of the SGD optimizer, we also follow the setting of Chen and He (2021) used for their CIFAR-10 experiments (see Appendix D in their paper): the momentum is set 0.9, and the weight decay is set 0.0005. For the output dimension of encoders, we set 512.

**Linear Evaluation**   We follow the configurations of Chen and He (2021) for linear evaluation: for the SGD optimizer, the momentum is 0.9, the weight decay is 0, the batch size is fixed to 256, and the learning rate is 30.0, where the linear scaling (Goyal et al., 2017) is applied to the learning rate. Note that we train the linear head for 100 epochs.

### H.1.2 Kernel Functions

In the experiments, we use the following kernel functions:

---

[2]https://github.com/jhaochenz/spectral_contrastive_learning/blob/ee431bdba9bb62ad00a7e55792213ee37712784c/optimizers/lr_scheduler.py (Last accessed: March 25, 2023)

**Gaussian Kernel.**  The Gaussian kernel $k_{\text{Gauss}}$ is defined as,

$$k_{\text{Gauss}}(z, z') = \exp\left(-\frac{\|z - z'\|_2^2}{\sigma^2}\right),$$

where $\sigma^2 > 0$ is the bandwidth parameter.

**Quadratic Kernel.**  The Quadratic kernel $k_{\text{Gauss}}$ is defined as,

$$k_{\text{Quad}}(z, z') = \left(z^\top z'\right)^2.$$

### H.1.3   Loss Functions

For the implementation of the kernel contrastive loss, we implement the empirical kernel contrastive loss (4). Note that in our experiments, the KCL frameworks with the Gaussian kernel and quadratic kernel are called Gaussian KCL (GKCL), and Quadratic KCL (QKCL), respectively.

For comparison, we also perform several reproducing experiments for SimCLR (Chen et al., 2020a) and SCL (HaoChen et al., 2021). For the implementation of the objective function of SimCLR, we use lightly.loss.NTXentLoss[3] of Lightly (Susmelj et al., 2020). For the implementation of the spectral contrastive loss of SCL, we adapt the implementation of the official SCL code[4].

### H.1.4   Datasets

In the experiments, we use the following datasets: CIFAR-10 (Krizhevsky, 2009), STL-10 (Coates et al., 2011), and ImageNet-100 (Tian et al., 2020). Note that ImageNet-100 is a subset of the ImageNet-1K dataset (Deng et al., 2009), where the ImageNet-100 dataset contains images categorized in 100 classes (Tian et al., 2020). When extracting images from the original the ImageNet-1K dataset to create the ImageNet-100 dataset, we select the 100 classes used in Tian et al. (2020). We also remark that for the experiments with the STL-10 dataset, we use the mixed dataset that consists of the unlabeled images and the labeled training images for pretraining, the labeled training images for the training of the linear head in the stage of linear evaluation, and the labeled test images for computing the accuracy in linear evaluation. Throughout the experiments, we use the following image size for each dataset: $32 \times 32$ pixels for CIFAR-10, $96 \times 96$ pixels for STL-10, and $224 \times 224$ pixels for ImageNet-100, where the image sizes of CIFAR-10 and STL-10 are the same as the sizes of the original images, respectively, and the image sizes for ImageNet-100 are inspired by those for the ImageNet-1K dataset used in Chen et al. (2020a); Chen and He (2021).

### H.1.5   Detail of Architectures for the CIFAR-10 Experiments

In the experiments with the CIFAR-10 dataset, following the settings of Chen and He (2021); HaoChen et al. (2021); He et al. (2016), we modify the original ResNet-18 (He et al., 2016) as follows:

---

[3]Note that in our experiments, we use Lightly 1.2.25.

[4]https://github.com/jhaochenz/spectral_contrastive_learning/blob/ee431bdba9bb62ad00a7e55792213ee37712784c/models/spectral.py (Last accessed: March 25, 2023)

in the implementation code of ResNet[5] of torchvision (TorchVision maintainers and contributors, 2016), we replace the first convolution layer with that whose kernel size is 3, stride is 1, and padding is 1, and the maxpool layer with that whose kernel size is 1 and stride is 1.

### H.1.6 Supplementary Information of the Implementation

We use the following packages for the experiments: PyTorch (Paszke et al., 2019), torchvision (TorchVision maintainers and contributors, 2016), NumPy (Harris et al., 2020), Lightly (Susmelj et al., 2020), Matplotlib (Hunter, 2007), and seaborn (Waskom, 2021).

## H.2 Results of Linear Evaluation

We perform pretraining and linear evaluation with the CIFAR-10, STL-10, and ImageNet-100 datasets. In the stage of pretraining, we train the encoder models for 800 epochs. For the experiments with the CIFAR-10 and STL-10 datasets, the batch sizes are set 256. Besides, for the experiments with the ImageNet-100 dataset, we set 512 for the batch sizes. We select the following hyperparameters of the KCL frameworks for all the experiments reported in this subsection: $\sigma^2 = 1$ and $\lambda = 8$ for GKCL, and $\lambda = 4$ for QKCL. For SCL, inspired by HaoChen et al. (2021), we select 3 for the radius parameter. For SimCLR, inspired by Chen et al. (2020a), we select 0.1 for the temperature parameter. These hyperparameters are also used for all the experiments in this subsection.

The results are shown in Table 1. In Table 1, the experiments with the CIFAR-10 dataset are performed using one Quadro P6000 GPU. Besides, the experiments with STL-10 and ImageNet-100 are performed using one Tesla V100S GPU.

## H.3 More Experiments

### H.3.1 Ablation Study on the Weight Parameters and Batch Sizes

We investigate how the selection of $\lambda$ and batch sizes affect the quality of learned representations. In the experiments, we use the CIFAR-10 dataset. We select $\{1, 2, 4, 8, 16, 32\}$ for $\lambda$, and $\{64, 128, 256, 512, 1024\}$ for the batch sizes. In each run, we pretrain an encoder model for 200 epochs. We use both GKCL and QKCL and evaluate those results.

---

[5]https://github.com/pytorch/vision/blob/eac3dc7bab436725b0ba65e556d3a6ffd43c24e1/torchvision/models/resnet.py (Last accessed: March 26, 2023)

Table 1: Top-1 and Top-5 accuracy (%) in linear evaluation for each method. For the CIFAR-10 and STL-10 experiments, we perform three trials of "pretraining+linear evaluation," and the results indicate the mean±standard deviation. For the ImageNet-100 experiments, we perform one trial of "pretraining+linear evaluation." All the results reported below are obtained as follows: we trained the linear heads for 100 epochs and evaluated the final classification accuracy of the models with the corresponding validation or test dataset. The word "repro." is the abbreviation for "reproducing," meaning that we performed several reproducing experiments to compare to the performance of KCL.

|  | CIFAR-10 | | STL-10 | | ImageNet-100 | |
|  | Top-1 | Top-5 | Top-1 | Top-5 | Top-1 | Top-5 |
| --- | --- | --- | --- | --- | --- | --- |
| SimCLR (repro.) | 90.09±0.05 | 99.70±0.01 | 87.23±0.35 | 99.56±0.04 | 77.26 | 94.06 |
| SCL (repro.) | 91.53±0.10 | 99.71±0.05 | 86.68±0.12 | 99.49±0.07 | 75.22 | 93.36 |
| GKCL | 90.87±0.08 | 99.63±0.00 | 86.69±0.09 | 99.38±0.01 | 76.40 | 93.20 |
| QKCL | 90.62±0.10 | 99.59±0.05 | 87.07±0.20 | 99.37±0.03 | 77.12 | 93.96 |



Figure 2: The results of ablation study for GKCL. The number in each cell indicates the Top-1 accuracy (%).

The Top-1 accuracy computed at the end of linear evaluation for GKCL and QKCL are shown in Figure 2 and 3, respectively. Note that the experiments reported in Figure 2 and 3 are performed by using one Tesla V100S GPU. The results of the experiments indicate that 1) the selection of the value for $\lambda$ affects the quality of the representations learned by KCL, 2) the small batch sizes (e.g., 128 and 256) are more efficient when pretraining encoders, while the large batch sizes (e.g., 1024)

Figure 3: The results of ablation study for QKCL. The number in each cell indicates the Top-1 accuracy (%).

degrade the performance. Note that Chen et al. (2021) showed similar findings to the first point for the generalized NT-Xent loss. Besides, Chen and He (2021) point out the efficiency of SimSiam with small batch sizes.

### H.3.2 How Does $\lambda$ Influence the Geometry of Representations Learned?

From Theorem 1, minimization of the kernel contrastive loss makes $\lambda \cdot \mathfrak{c}(f)$ smaller, which can imply that the means of the clusters tend to distribute uniformly as $\lambda$ increases. Motivated by this result, in this subsubsection, we simulate how the mean of the feature vectors belonging to each cluster distributes. In the experiments, we use the STL-10 dataset. In the stage of unsupervised pretraining, we use the combination of the unlabeled images and the labeled training images in the STL-10 dataset. We use GKCL for the pretraining. The weights used in the experiments are $\{1, 2, 4, 8, 16, 32, 64\}$. We pretrain the encoder model for 400 epochs in each run. We set the batch sizes to 256. After the stage of pretraining, we compute the mean for each class and calculate the cosine similarities between those means. Since the clusters $\mathbb{M}_1, \cdots, \mathbb{M}_K$ are hard to obtain for the STL-10 dataset, we instead use the labels included in the labeled training images of the STL-10 dataset to compute the mean over the feature vectors of augmented data transformed from raw data in each class. Note that we draw an augmented image from each raw image when computing the means. The experiments in this subsubsection are performed by using one Tesla V100S GPU.

The results are summarized in Table 4 as a box plot. The results indicate that the variation becomes smaller as $\lambda$ increases. Thus, larger $\lambda$ makes the means to distribute more uniformly in this experimental setting. Note that this result may imply that the clusters $\mathbb{M}_1, \cdots, \mathbb{M}_K$ and the

Figure 4: The box plot of the cosine similarities between the means of the different classes for each encoder model pretrained with different $\lambda$. Note that the horizontal lines in each bar represent, in the order from bottom to top, the minimum value, the first quartile, the median, the third quartile, and the maximum value, respectively.

subsets defined with labels have some relation. We leave the investigation of this question as future work.

## H.4  Too Large $\lambda$ Degrades the Performance in Downstream Classification Tasks

In this subsection, we report the results of the experiments with different values for $\lambda$. In the experiments, we use $\{1, 2, 4, 8, 16, 32, 64, 128, 256, 512\}$ for the weight $\lambda$. For the contrastive learning framework, we use GKCL. We use two datasets, CIFAR-10 and STL-10, and pretrain the encoder during 400 epochs in each run. In the stage of pretraining, we set 128 for the batch size. Each experiment reported in this subsection is performed using one Tesla V100S GPU.

Table 2: Top-1 accuracy (%) in the results of linear evaluation, where the encoder is pretrained with different $\lambda$ for each run.

| | Top-1 Accuracy | |
|---|---|---|
| $\lambda$ | CIFAR-10 | STL-10 |
| 1 | 89.06 | 83.20 |
| 2 | 90.29 | 84.54 |
| 4 | 90.77 | 85.36 |
| 8 | 90.66 | 85.33 |
| 16 | 90.52 | 84.19 |
| 32 | 89.78 | 83.39 |
| 64 | 88.85 | 81.50 |
| 128 | 86.93 | 79.71 |
| 256 | 85.24 | 77.89 |
| 512 | 82.43 | 76.26 |

The results on the Top-1 accuracy at the end of the linear evaluation are presented in Table 2. The results indicate that too large $\lambda$, such as 512, degrades the performance in the downstream task.