

Provable Guarantees for Self-Supervised Deep Learning with Spectral Contrastive Loss

Jeff Z. HaoChen¹ Colin Wei¹ Adrien Gaidon² Tengyu Ma¹

¹ Stanford University ² Toyota Research Institute

{jhaochen, colinwei, tengyuma}@stanford.edu adrien.gaidon@tri.global

Abstract

Recent works in self-supervised learning have advanced the state-of-the-art by relying on the *contrastive learning* paradigm, which learns representations by pushing positive pairs, or similar examples from the same class, closer together while keeping negative pairs far apart. Despite the empirical successes, theoretical foundations are limited – prior analyses assume conditional independence of the positive pairs given the same class label, but recent empirical applications use heavily correlated positive pairs (i.e., data augmentations of the same image). Our work analyzes contrastive learning without assuming conditional independence of positive pairs using a novel concept of the *augmentation graph* on data. Edges in this graph connect augmentations of the same datapoint, and ground-truth classes naturally form connected sub-graphs. We propose a loss that performs spectral decomposition on the population augmentation graph and can be succinctly written as a contrastive learning objective on neural net representations. Minimizing this objective leads to features with provable accuracy guarantees under linear probe evaluation. By standard generalization bounds, these accuracy guarantees also hold when minimizing the training contrastive loss. Empirically, the features learned by our objective can match or outperform several strong baselines on benchmark vision datasets. In all, this work provides the first provable analysis for contrastive learning where guarantees for linear probe evaluation can apply to realistic empirical settings.

1 Introduction

Recent empirical breakthroughs have demonstrated the effectiveness of self-supervised learning, which trains representations on unlabeled data with surrogate losses and self-defined supervision signals (Bachman et al., 2019, Bardes et al., 2021, Caron et al., 2020, Chen and He, 2020, Henaff, 2020, Hjelm et al., 2018, Misra and Maaten, 2020, Oord et al., 2018, Tian et al., 2019, 2020a, Wu et al., 2018, Ye et al., 2019, Zbontar et al., 2021). Self-supervision signals in computer vision are often defined by using data augmentation to produce multiple views of the same image. For example, the recent contrastive learning objectives (Arora et al., 2019, Chen et al., 2020a,b,c, He et al., 2020) encourage closer representations for augmentations/views of the same natural datapoint than for randomly sampled pairs of data.

Despite the empirical successes, there is a limited theoretical understanding of why self-supervised losses learn representations that can be adapted to downstream tasks, for example, using linear heads. Recent mathematical analyses for contrastive learning by Arora et al. (2019), Tosh et al. (2020, 2021) provide guarantees under the assumption that two views are somewhat conditionally independent given the label or a hidden variable. However, in practical algorithms

for computer vision applications, the two views are augmentations of a natural image and usually exhibit a strong correlation that is difficult to be de-correlated by conditioning. They are not independent conditioned on the label, and we are only aware that they are conditionally independent given the natural image, which is too complex to serve as a hidden variable with which prior works can be meaningfully applied. Thus the existing theory does not appear to explain the practical success of self-supervised learning.

This paper presents a theoretical framework for self-supervised learning without requiring conditional independence. We design a principled, practical loss function for learning neural net representations that resembles state-of-the-art contrastive learning methods. We prove that, under a simple and realistic data assumption, linear classification using representations learned on a polynomial number of unlabeled data samples can recover the ground-truth labels of the data with high accuracy.

The fundamental data property that we leverage is a notion of continuity of the *population* data within the same class. Though a random pair of images from the same class can be far apart, the pair is often connected by (many) sequences of natural images, where consecutive images in the sequences are close neighbors within the same class. As shown in Figure 1 (images on the left top part), two very different French bulldogs can be connected by a sequence of French bulldogs (which may not be in the training set but are in the support of the population distribution). Prior work (Wei et al., 2020) empirically demonstrates this type of connectivity property and uses it in the analysis of pseudolabeling algorithms. This property is more salient when the neighborhood of an example includes many different types of augmentations.

More formally, we define the *population augmentation graph*, whose vertices are all the augmented data in the *population* distribution, which can be an exponentially large or infinite set. Two vertices are connected with an edge if they are augmentations of the same natural example. Our main assumption is that for some proper $m \in \mathbb{Z}^+$, we cannot partition the graph into $m + 1$ sub-graphs between which there are few connections (Assumption 3.5). In other words, this intuitively states that there are at most m clusters in the population augmentation graph. This assumption can be seen as a graph-theoretic version of the continuity assumption on the *population* distribution. We also assume that there are very few edges across different ground-truth classes (Assumption 3.6). Figure 1 (left) illustrates a realistic scenario where dog and cat are the ground-truth categories, between which edges are very rare. Each breed forms a sub-graph that has sufficient inner connectivity and thus cannot be further partitioned.

Our assumption fundamentally does not require independence of the two views (the positive pairs) conditioned on the class and can allow disconnected sub-graphs within a class. The classes in the downstream task can be also somewhat flexible as long as they are disconnected in the augmentation graph. For example, when the augmentation graph consists of m disconnected sub-graphs corresponding to fine-grained classes, our assumptions allow the downstream task to have any $r \leq m$ coarse-grained classes containing these fine-grained classes as a sub-partition. Prior work (Wei et al., 2020) on pseudolabeling algorithms essentially requires an exact alignment between sub-graphs and downstream classes (i.e., $r = m$). They face this limitation because their analysis requires fitting discrete pseudolabels on the unlabeled data. We avoid this difficulty because we consider directly learning continuous representations on the unlabeled data.

The main insight of the paper is that contrastive learning can be viewed as a parametric form of spectral clustering (Ng et al., 2001, Shi and Malik, 2000) on the augmentation graph. Concretely, suppose we apply spectral decomposition or spectral clustering—a classical approach for graph partitioning—to the adjacency matrix defined on the population augmentation graph. We form a

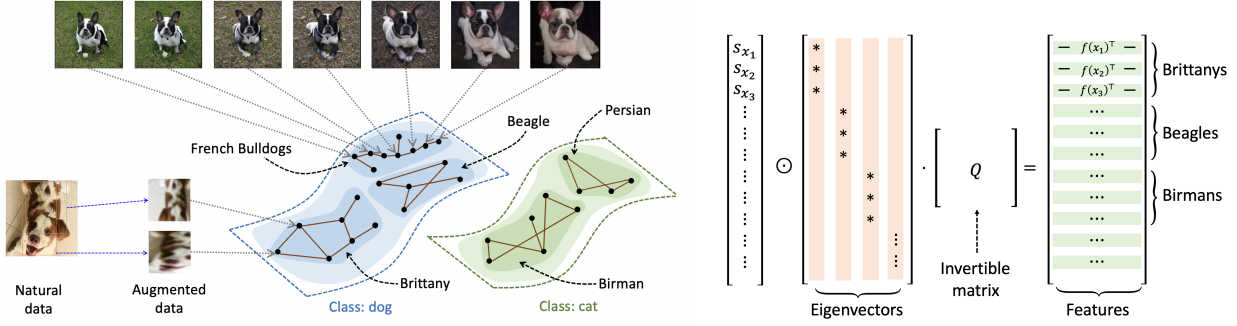


Figure 1: **Left: demonstration of the population augmentation graph.** Two augmented data are connected if they are views of the same natural datapoint. Augmentations of data from different classes in the downstream tasks are assumed to be nearly disconnected, whereas there are more connections within the same class. We allow the existence of disconnected sub-graphs within a class corresponding to potential sub-classes. **Right: decomposition of the learned representations.** The representations (rows in the RHS) learned by minimizing the population spectral contrastive loss can be decomposed as the LHS. The scalar s_{x_i} is positive for every augmented datapoint x_i . Columns of the matrix labeled “eigenvectors” are the top eigenvectors of the normalized adjacency matrix of the augmentation graph defined in Section 3.1. The operator \odot multiplies row-wise each s_{x_i} with the x_i -th row of the eigenvector matrix. When classes (or sub-classes) are exactly disconnected in the augmentation graph, the eigenvectors are sparse and align with the sub-class structure. The invertible Q matrix does not affect the performance of the rows under the linear probe.

matrix where the top- k eigenvectors are the columns and interpret each row of the matrix as the representation (in \mathbb{R}^k) of an example. Somewhat surprisingly, we prove that this feature extractor can be also recovered (up to some linear transformation) by minimizing the following population objective which is a variant of the standard contrastive loss:

$$\mathcal{L}(f) = -2 \cdot \mathbb{E}_{x, x^+} \left[f(x)^\top f(x^+) \right] + \mathbb{E}_{x, x^-} \left[\left(f(x)^\top f(x^-) \right)^2 \right],$$

where (x, x^+) is a pair of augmentations of the same datapoint, (x, x^-) is a pair of independently random augmented data, and f is a parameterized function from augmented data to \mathbb{R}^k . Figure 1 (right) illustrates the relationship between the eigenvector matrix and the learned representations. We call this loss the *population spectral contrastive loss*.

We analyze the linear classification performance of the representations learned by minimizing the population spectral contrastive loss. Our main result (Theorem 3.8) shows that when the representation dimension exceeds the maximum number of disconnected sub-graphs, linear classification with learned representations is guaranteed to have a small error. Our theorem reveals a trend that a larger representation dimension is needed when there are a larger number of disconnected sub-graphs. Our analysis relies on novel techniques tailored to linear probe performance, which have not been studied in the spectral graph theory community to the best of our knowledge.

The spectral contrastive loss also works on empirical data. Since our approach optimizes parametric loss functions, guarantees involving the population loss can be converted to finite sample results using off-the-shelf generalization bounds. The end-to-end result (Theorem 4.3) shows that the number of unlabeled examples required is polynomial in the Rademacher complexity of the

model family and other relevant parameters, whereas the number of downstream labeled examples only needs to be linear in the representation dimension (which needs to be linear in the number of clusters in the graph). This demonstrates that contrastive learning reduces the amount of labeled examples needed.

In summary, our main theoretical contributions are: 1) we propose a simple contrastive loss motivated by spectral decomposition of the population data graph, 2) under simple and realistic assumptions, we provide downstream classification guarantees for the representation learned by minimizing this loss on population data, and 3) our analysis is easily applicable to deep networks with polynomial unlabeled samples via off-the-shelf generalization bounds. Our theoretical framework can be viewed as containing two stages: we first analyze the population loss and the representation that minimizes it (Section 3), then study the empirical loss where the representation is learned with a neural network with bounded capacity (Section 4).

In addition, we implement and test the proposed spectral contrastive loss on standard vision benchmark datasets. Our algorithm is simple and doesn't rely on tricks such as stop-gradient which is essential to SimSiam (Chen and He, 2020). We demonstrate that the features learned by our algorithm can match or outperform several strong baselines (Chen et al., 2020a, Chen and He, 2020, Chen et al., 2020c, Grill et al., 2020) when evaluated using a linear probe.

2 Additional related works

Empirical works on self-supervised learning. Self-supervised learning algorithms have been shown to successfully learn representations that benefit downstream tasks (Bachman et al., 2019, Bardes et al., 2021, Caron et al., 2020, Chen et al., 2020a,b,c, He et al., 2020, Henaff, 2020, Hjelm et al., 2018, Misra and Maaten, 2020, Oord et al., 2018, Tian et al., 2019, 2020a, Wu et al., 2018, Xie et al., 2019, Ye et al., 2019, Zbontar et al., 2021). Many recent self-supervised learning algorithms learn features with siamese networks (Bromley et al., 1993), where two neural networks of shared weights are applied to pairs of augmented data. Introducing asymmetry to siamese networks either with a momentum encoder like BYOL (Grill et al., 2020) or by stopping gradient propagation for one branch of the siamese network like SimSiam (Chen and He, 2020) has been shown to effectively avoid collapsing. Contrastive methods (Chen et al., 2020a,c, He et al., 2020) minimize the InfoNCE loss (Oord et al., 2018), where two views of the same data are attracted while views from different data are repulsed.

Theoretical works on self-supervised learning. As briefly discussed in the introduction, several theoretical works have studied self-supervised learning. Arora et al. (2019) provide guarantees for representations learned by contrastive learning on downstream linear classification tasks under the assumption that the positive pairs are conditionally independent given the class label. Theorem 3.3 and Theorem 3.7 of the work of Lee et al. (2020) show that, under conditional independence given the label and/or additional latent variables, representations learned by reconstruction-based self-supervised learning algorithms can achieve small errors in the downstream linear classification task. Lee et al. (2020, Theorem 4.1) generalizes it to approximate conditional independence for Gaussian data and Theorem 4.5 further weakens the assumptions significantly. Tosh et al. (2020) show that contrastive learning representations can linearly recover any continuous functions of the underlying topic posterior under a topic modeling assumption (which also requires conditional independence of the positive pair given the hidden variable). More recently, Theorem 11 of the work of Tosh et al. (2021) provide novel guarantees for contrastive learning under the assumption that there exists a hidden variable h such that the positive pair (x, x^+) are conditionally independent given h and the random variable $p(x|h)p(x^+|h)/p(x)p(x^+)$ has a small variance. However, in prac-

tical algorithms for computer vision applications, the two views are two augmentations and thus they are highly correlated. They might be only independent when conditioned on very complex hidden variables such as the original natural image, which might be too complex for the previous results to be meaningfully applied.

We can also compare the assumptions and results on a concrete generative model for the data, our Example 3.10 in Section 3.4, where the data are generated by a mixture of Gaussian or a mixture of manifolds, the label is the index of the mixture, and the augmentations are small Gaussian blurring (i.e., adding Gaussian noise). In this case, the positive pairs (x, x^+) are two points that are very close to each other. To the best of our knowledge, applying Theorem 11 of Tosh et al. (2021) to this case with $h = \bar{x}$ (the natural datapoint) would result in requiring a large (if not infinite) representation dimension. Because x^+ and x are very close, the reconstruction-based algorithms in Lee et al. (2020), when used to predict x^+ from x , will not be able to produce good representations as well.¹

On a technical level, to relate prior works’ assumptions to ours, we can consider an almost equivalent version of our assumption (although our proofs do not directly rely on or relate to the discussion below). Let (x, x^+) be a positive pair and let $p(\cdot|x)$ be the conditional distribution of x^+ given x . Starting from x_0 , let us consider a hypothetical Markov chain x_0, \dots, x_T, \dots where x_t is drawn from $p(\cdot|x_{t-1})$. Our assumption essentially means that this hypothetical Markov chain of sampling neighbors will mix within the same class earlier than it mixes across the entire population (which might not be possible or takes exponential time). More concretely, the assumption that $\rho_{\lfloor k/2 \rfloor}$ is large compared to α in Theorem 3.8 is roughly equivalent to the existence of a (potentially large) T such that x_0 and x_T are still likely to have the same label, but are sufficiently independent conditioned on this label or some hidden variable. Roughly speaking, prior works (Arora et al., 2019, Tosh et al., 2020, 2021) assume probabilistic structure about x_0 and x_1 (instead of x_0 and x_T), e.g., Arora et al. (2019) and Theorem 11 of Tosh et al. (2021) assume that x_0 and x_1 are independent conditioned on the label and/or a hidden variable. Similar Markov chains on augmented data have also been used in previous work (Dao et al., 2019) to study properties of data augmentation.

Several other works (Bansal et al., 2020, Mitrovic et al., 2020, Tian et al., 2020b, Tsai et al., 2020, Wang and Isola, 2020) also theoretically study self-supervised learning. The work Tsai et al. (2020) prove that self-supervised learning methods can extract task-relevant information and discard task-irrelevant information, but lacks guarantees for solving downstream tasks efficiently with simple (e.g., linear) models. Tian et al. (2020b) study why non-contrastive self-supervised learning methods can avoid feature collapse. Zimmermann et al. (2021) prove that for a specific data generating process, contrastive learning can learn representations that recover the latent variable. Cai et al. (2021) analyze domain adaptation algorithms for subpopulation shift with a similar expansion condition as (Wei et al., 2020) while also allowing disconnected parts within each class, but require access to ground-truth labels during training. In contrast, our algorithm doesn’t need labels during pre-training.

Co-training and multi-view learning are related settings which leverage two distinct “views” (i.e., feature subsets) of the data (Balcan et al., 2005, Blum and Mitchell, 1998, Dasgupta et al., 2002). The original co-training algorithms (Blum and Mitchell, 1998, Dasgupta et al., 2002) assume that the two views are independent conditioned on the true label and leverage this independence to obtain accurate pseudolabels for the unlabeled data. Balcan et al. (2005) relax the requirement on independent views of co-training, by using an “expansion” assumption, which is closely related

¹On a technical level, Example 3.10 does not satisfy the requirement regarding the β quantity in Assumption 4.1 of Lee et al. (2020), if (X_1, X_2) in that paper is equal to (x, x^+) here—it requires the label to be correlated with the raw input x , which is not necessarily true in Example 3.10. This can likely be addressed by using a different X_2 .

to our assumption that $\rho_{\lfloor k/2 \rfloor}$ is not too small in Theorem 3.8. Besides recent works (e.g., the work of Toshi et al. (2021)), most co-training or multi-view learning algorithms are quite different from the modern contrastive learning algorithms which use neural network parameterization for vision applications.

Our analysis relies on the normalized adjacency matrix (see Section 3.1), which is closely related to the graph Laplacian regularization that has been studied in the setting of semi-supervised learning (Nadler et al., 2009, Zhu et al., 2003). In their works, the Laplacian matrix is used to define a regularization term that smooths the predictions on unlabeled data. This regularizer is further added to the supervised loss on labeled data during training. In contrast, we use the normalized adjacency matrix to define the unsupervised training objective in this paper.

3 Spectral contrastive learning on population data

In this section, we introduce our theoretical framework, the spectral contrastive loss, and the main analysis of the performance of the representations learned on population data.

We use $\bar{\mathcal{X}}$ to denote the set of all natural data (raw inputs without augmentation). We assume that each $\bar{x} \in \bar{\mathcal{X}}$ belongs to one of r classes, and let $y : \bar{\mathcal{X}} \rightarrow [r]$ denote the ground-truth (deterministic) labeling function. Let $\mathcal{P}_{\bar{\mathcal{X}}}$ be the population distribution over $\bar{\mathcal{X}}$ from which we draw training data and test our final performance. In the main body of the paper, for the ease of exposition, we assume $\bar{\mathcal{X}}$ to be a finite but exponentially large set (e.g., all real vectors in \mathbb{R}^d with bounded precision). This allows us to use sums instead of integrals and avoid non-essential nuances/jargons related to functional analysis. See Section F for the straightforward extensions to the case where $\bar{\mathcal{X}}$ is an infinite compact set (with mild regularity conditions).²

We next formulate data augmentations. Given a natural data sample $\bar{x} \in \bar{\mathcal{X}}$, we use $\mathcal{A}(\cdot|\bar{x})$ to denote the distribution of its augmentations. For instance, when \bar{x} represents an image, $\mathcal{A}(\cdot|\bar{x})$ can be the distribution of common augmentations (Chen et al., 2020a) that includes Gaussian blur, color distortion and random cropping. We use \mathcal{X} to denote the set of all augmented data, which is the union of supports of all $\mathcal{A}(\cdot|\bar{x})$ for $\bar{x} \in \bar{\mathcal{X}}$. As with $\bar{\mathcal{X}}$, we also assume that \mathcal{X} is a finite but exponentially large set, and denote $N = |\mathcal{X}|$. None of the bounds will depend on N — it is only defined and assumed to be finite for the ease of exposition.

We will learn an embedding function $f : \mathcal{X} \rightarrow \mathbb{R}^k$, and then evaluate its quality by the minimum error achieved with a linear probe. Concretely, a linear classifier has weights $B \in \mathbb{R}^{k \times r}$ and predicts $g_{f,B}(x) = \arg \max_{i \in [r]} (f(x)^\top B)_i$ for an augmented datapoint x ($\arg \max$ breaks tie arbitrarily). Then, given a natural data sample \bar{x} , we ensemble the predictions on augmented data and predict:

$$\bar{g}_{f,B}(\bar{x}) := \arg \max_{i \in [r]} \Pr_{x \sim \mathcal{A}(\cdot|\bar{x})} [g_{f,B}(x) = i].$$

We denote the error of the representation and the linear head as:

$$\mathcal{E}(f, B) := \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}} [y(\bar{x}) \neq \bar{g}_{f,B}(\bar{x})].$$

Define the *linear probe* error as the error of the best possible linear classifier on the representations:

$$\mathcal{E}(f) := \min_{B \in \mathbb{R}^{k \times r}} \mathcal{E}(f, B) = \min_{B \in \mathbb{R}^{k \times r}} \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}} [y(\bar{x}) \neq \bar{g}_{f,B}(\bar{x})]. \quad (1)$$

²In Section F, we will deal with an infinite graph, its adjacency operator (instead of adjacency matrix), and the eigenfunctions of the adjacency operator (instead of eigenvectors) essentially in the same way.

3.1 Augmentation graph and spectral decomposition

Our approach is based on the central concept of **population augmentation graph**, denoted by $G(\mathcal{X}, w)$, where the vertex set is all augmentation data \mathcal{X} and w denotes the edge weights defined below. For any two augmented data $x, x' \in \mathcal{X}$, define the weight $w_{xx'}$ as the marginal probability of generating the pair x and x' from a random natural data $\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}$:

$$w_{xx'} := \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}} [\mathcal{A}(x|\bar{x})\mathcal{A}(x'|\bar{x})] \quad (2)$$

Therefore, the weights sum to 1 because the total probability mass is 1: $\sum_{x, x' \in \mathcal{X}} w_{xx'} = 1$. The relative magnitude intuitively captures the closeness between x and x' with respect to the augmentation transformation. For most of the unrelated x and x' , the value $w_{xx'}$ will be significantly smaller than the average value. For example, when x and x' are random croppings of a cat and a dog respectively, $w_{xx'}$ will be essentially zero because no natural data can be augmented into both x and x' . On the other hand, when x and x' are very close in ℓ_2 -distance or very close in ℓ_2 -distance up to color distortion, $w_{xx'}$ is nonzero because they may be augmentations of the same image with Gaussian blur and color distortion. We say that x and x' are connected with an edge if $w_{xx'} > 0$. See Figure 1 (left) for more illustrations.

We emphasize that we *only* work with the population graph rather than the empirical graph (i.e., the corresponding graph constructed with the empirical dataset as the vertex set). The population graph is very sparse but not empty—many similar images exist in the population. In contrast, the empirical graph would be nearly empty, since two images in the empirical dataset almost never share the same augmentation image. Our analysis will apply to minimizing contrastive loss on an empirical dataset (see Section 4), but *not* via analyzing the property of the empirical graph. Instead, we will show that contrastive learning on *empirical* data with parametrized models is similar to decomposing the *population* graph (see technical discussions in Section 5). This is a key difference between our work and classical spectral clustering work—we only require properties of the population graph rather than the empirical graph.

A simplified running example with Gaussian perturbation augmentation. Suppose the natural data is supported on manifolds in Euclidean space, and the data augmentation is adding random noise sampled from $\mathcal{N}(0, \sigma^2 \cdot I_{d \times d})$ where σ is a small quantity (e.g., the norm of the perturbation $\sigma\sqrt{d}$ should be much smaller than the norm of the original datapoint). Then the edge between two augmented datapoints would be have near zero weight unless the two datapoints have small ℓ_2 distance. Hence, the resulting graph is essentially the ϵ -ball proximity graph (Zemel and Carreira-Perpiñán, 2004) or geometric graph (Penrose, 2003) in Euclidean space.

Given the structure of the population augmentation graph, we apply spectral decomposition to the population graph to construct principled embeddings. The eigenvalue problems are closely related to graph partitioning as shown in spectral graph theory (Chung and Graham, 1997) for both worst-case graphs (Cheeger, 1969, Kannan et al., 2004, Lee et al., 2014, Louis et al., 2011) and random graphs (Abbe, 2017, Lei et al., 2015, McSherry, 2001). In machine learning, spectral clustering (Ng et al., 2001, Shi and Malik, 2000) is a classical algorithm that learns embeddings by eigendecomposition on an empirical distance graph and invoking k -means on the embeddings.

We will apply eigendecomposition to the *population* augmentation graph (and then later use linear probe for classification). Let $w_x = \sum_{x' \in \mathcal{X}} w_{xx'}$ be the total weights associated to x , which is often viewed as an analog of the degree of x in weighted graph. A central object in spectral graph theory is the so-called *normalized adjacency matrix*:

$$\bar{A} := D^{-1/2} A D^{-1/2} \quad (3)$$

where $A \in \mathbb{R}^{N \times N}$ is adjacency matrix with entries $A_{xx'} = w_{xx'}$ and $D \in \mathbb{R}^{N \times N}$ is a diagonal matrix with $D_{xx} = w_x$.³

Standard spectral graph theory approaches produce vertex embeddings as follows. Let $\gamma_1, \gamma_2, \dots, \gamma_k$ be the k largest eigenvalues of \bar{A} , and v_1, v_2, \dots, v_k be the corresponding unit-norm eigenvectors. Let $F^* = [v_1, v_2, \dots, v_k] \in \mathbb{R}^{N \times k}$ be the matrix that collects these eigenvectors in columns, and we refer to it as the eigenvector matrix. Let $u_x^* \in \mathbb{R}^k$ be the x -th row of the matrix F^* . It turns out that u_x^* 's can serve as desirable embeddings of x 's because they exhibit clustering structure in Euclidean space that resembles the clustering structure of the graph $G(\mathcal{X}, w)$.

3.2 From spectral decomposition to spectral contrastive learning

The embeddings u_x^* obtained by eigendecomposition are nonparametric—a k -dimensional parameter is needed for every x —and therefore cannot be learned with a realistic amount of data. The embedding matrix F^* cannot be even stored efficiently. Therefore, we will instead parameterize the rows of the eigenvector matrix F^* as a neural net function, and assume embeddings u_x^* can be represented by $f(x)$ for some $f \in \mathcal{F}$, where \mathcal{F} is the hypothesis class containing neural networks. As we'll show in Section 4, this allows us to leverage the extrapolation power of neural networks and learn the representation on a finite dataset.

Next, we design a proper loss function for the feature extractor f , such that minimizing this loss could recover F^* up to some linear transformation. As we will show in Section 4, the resulting population loss function on f also admits an unbiased estimator with finite training samples. Let F be an embedding matrix with u_x on the x -th row, we will first design a loss function of F that can be decomposed into parts about individual rows of F .

We employ the following matrix factorization based formulation for eigenvectors. Consider the objective

$$\min_{F \in \mathbb{R}^{N \times k}} \mathcal{L}_{\text{mf}}(F) := \left\| \bar{A} - FF^\top \right\|_F^2. \quad (4)$$

By the classical theory on low-rank approximation (Eckart–Young–Mirsky theorem (Eckart and Young, 1936)), any minimizer \hat{F} of $\mathcal{L}_{\text{mf}}(F)$ contains scaling of the largest eigenvectors of \bar{A} up to a right transformation—for some orthonormal matrix $R \in \mathbb{R}^{k \times k}$, we have $\hat{F} = F^* \cdot \text{diag}([\sqrt{\gamma_1}, \dots, \sqrt{\gamma_k}])R$. Fortunately, multiplying the embedding matrix by any matrix on the right and any diagonal matrix on the left does not change its linear probe performance, which is formalized by the following lemma.

Lemma 3.1. *Consider an embedding matrix $F \in \mathbb{R}^{N \times k}$ and a linear classifier $B \in \mathbb{R}^{k \times r}$. Let $D \in \mathbb{R}^{N \times N}$ be a diagonal matrix with positive diagonal entries and $Q \in \mathbb{R}^{k \times k}$ be an invertible matrix. Then, for any embedding matrix $\tilde{F} = D \cdot F \cdot Q$, the linear classifier $\tilde{B} = Q^{-1}B$ on \tilde{F} has the same prediction as B on F . As a consequence, we have*

$$\mathcal{E}(F) = \mathcal{E}(\tilde{F}). \quad (5)$$

where $\mathcal{E}(F)$ denotes the linear probe performance when the rows of F are used as embeddings.

Proof of Lemma 3.1. Let $D = \text{diag}(s)$ where $s_x > 0$ for $x \in \mathcal{X}$. Let $u_x, \tilde{u}_x \in \mathbb{R}^k$ be the x -th row of matrices F and \tilde{F} , respectively. Recall that $g_{u, B}(x) = \arg \max_{i \in [r]} (u_x^\top B)_i$ is the prediction on an augmented datapoint $x \in \bar{\mathcal{X}}$ with representation u_x and linear classifier B . Let $\tilde{B} = Q^{-1}B$, it's easy

³We index the matrix A, D by $(x, x') \in \mathcal{X} \times \mathcal{X}$. Generally we index N -dimensional axis by $x \in \mathcal{X}$.

to see that $g_{\tilde{u}, \tilde{B}}(x) = \arg \max_{i \in [r]} (s_x \cdot u_x^\top B)_i$. Notice that $s_x > 0$ doesn't change the prediction since it changes all dimensions of $u_x^\top B$ by the same scale, we have $g_{\tilde{u}, \tilde{B}}(x) = g_{u, B}(x)$ for any augmented datapoint $x \in \mathcal{X}$. The equivalence of loss naturally follows. \square

The main benefit of objective $\mathcal{L}_{\text{mf}}(F)$ is that it's based on the rows of F . Recall that vectors u_x are the rows of F . Each entry of FF^\top is of the form $u_x^\top u_{x'}$, and thus $\mathcal{L}_{\text{mf}}(F)$ can be decomposed into a sum of N^2 terms involving terms $u_x^\top u_{x'}$. Interestingly, if we reparameterize each row u_x by $w_x^{1/2} f(x)$, we obtain a very similar loss function for f that resembles the contrastive learning loss used in practice (Chen et al., 2020a) as shown below in Lemma 3.2. See Figure 1 (right) for an illustration of the relationship between the eigenvector matrix and the representations learned by minimizing this loss.

We formally define the positive and negative pairs to introduce the loss. Let $\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}$ be a random natural datapoint and draw $x \sim \mathcal{A}(\cdot | \bar{x})$ and $x^+ \sim \mathcal{A}(\cdot | \bar{x})$ independently to form a positive pair (x, x^+) . Draw $\bar{x}' \sim \mathcal{P}_{\bar{\mathcal{X}}}$ and $x^- \sim \mathcal{A}(\cdot | \bar{x}')$ independently with \bar{x}, x, x^+ . We call (x, x^-) a negative pair.⁴

Lemma 3.2 (Spectral contrastive loss). *Recall that u_x is the x -th row of F . Let $u_x = w_x^{1/2} f(x)$ for some function f . Then, the loss function $\mathcal{L}_{\text{mf}}(F)$ is equivalent to the following loss function for f , called spectral contrastive loss, up to an additive constant:*

$$\begin{aligned} \mathcal{L}_{\text{mf}}(F) &= \mathcal{L}(f) + \text{const} \\ \text{where } \mathcal{L}(f) &\triangleq -2 \cdot \mathbb{E}_{x, x^+} [f(x)^\top f(x^+)] + \mathbb{E}_{x, x^-} \left[\left(f(x)^\top f(x^-) \right)^2 \right] \end{aligned} \quad (6)$$

Proof of Lemma 3.2. We can expand $\mathcal{L}_{\text{mf}}(F)$ and obtain

$$\begin{aligned} \mathcal{L}_{\text{mf}}(F) &= \sum_{x, x' \in \mathcal{X}} \left(\frac{w_{xx'}}{\sqrt{w_x w_{x'}}} - u_x^\top u_{x'} \right)^2 \\ &= \sum_{x, x' \in \mathcal{X}} \left(\frac{w_{xx'}^2}{w_x w_{x'}} - 2 \cdot w_{xx'} \cdot f(x)^\top f(x') + w_x w_{x'} \cdot \left(f(x)^\top f(x') \right)^2 \right) \end{aligned} \quad (7)$$

Notice that the first term is a constant that only depends on the graph but not the variable f . By the definition of augmentation graph, $w_{xx'}$ is the probability of a random positive pair being (x, x') while w_x is the probability of a random augmented datapoint being x . We can hence rewrite the sum of last two terms in Equation (7) as Equation (6). \square

We note that spectral contrastive loss is similar to many popular contrastive losses (Chen et al., 2020a, Oord et al., 2018, Sohn, 2016, Wu et al., 2018). For instance, the contrastive loss in SimCLR (Chen et al., 2020a) can be rewritten as (with simple algebraic manipulation)

$$-f(x)^\top f(x^+) + \log \left(\exp \left(f(x)^\top f(x^+) \right) + \sum_{i=1}^n \exp \left(f(x)^\top f(x_i) \right) \right).$$

Here x and x^+ are a positive pair and x_1, \dots, x_n are augmentations of other data. Spectral contrastive loss can be seen as removing $f(x)^\top f(x^+)$ from the second term, and replacing the log sum of exponential terms with the average of the squares of $f(x)^\top f(x_i)$. We will show in Section 6 that our loss has a similar empirical performance as SimCLR without requiring a large batch size.

⁴Though x and x^- are simply two independent draws, we call them negative pairs following the literature (Arora et al., 2019).

3.3 Theoretical guarantees for spectral contrastive loss on population data

In this section, we introduce the main assumptions on the data and state our main theoretical guarantee for spectral contrastive learning on population data.

To formalize the idea that G cannot be partitioned into too many disconnected sub-graphs, we introduce the notions of *Dirichlet conductance* and *sparsest m -partition*, which are standard in spectral graph theory. Dirichlet conductance represents the fraction of edges from S to its complement:

Definition 3.3 (Dirichlet conductance). *For a graph $G = (\mathcal{X}, w)$ and a subset $S \subseteq \mathcal{X}$, we define the Dirichlet conductance of S as*

$$\phi_G(S) := \frac{\sum_{x \in S, x' \notin S} w_{xx'}}{\sum_{x \in S} w_x}.$$

We note that when S is a singleton, there is $\phi_G(S) = 1$ due to the definition of w_x . For $i \in \mathbb{Z}^+$, we introduce the sparsest i -partition to represent the number of edges between i disjoint subsets.

Definition 3.4 (Sparsest i -partition). *Let $G = (\mathcal{X}, w)$ be the augmentation graph. For an integer $i \in [2, |\mathcal{X}|]$, we define the sparsest i -partition as*

$$\rho_i := \min_{S_1, \dots, S_i} \max\{\phi_G(S_1), \dots, \phi_G(S_i)\}$$

where S_1, \dots, S_i are non-empty sets that form a partition of \mathcal{X} .

We note that ρ_i increases as i increases.⁵ When r is the number of underlying classes, we might expect $\rho_r \approx 0$ since the augmentations from different classes almost compose a disjoint r -way partition of \mathcal{X} . However, for $i > r$, we can expect ρ_i to be much larger. For instance, in the extreme case when $i = |\mathcal{X}| = N$, every set S_j is a singleton, which implies that $\rho_N = 1$. More generally, as we will show later (Lemma 3.9), ρ_i can be expected to be at least inverse polynomial in data dimension when i is larger than the number of underlying semantic classes in the data.

Assumption 3.5 (at most m clusters). *We assume that $\rho_{m+1} \geq \rho$. A prototypical case would be that there are at most m clusters in the population augmentation graph, and each of them cannot be broken into two subsets both with conductance less than ρ .*

When there are m clusters that have sufficient inner connections (corresponding to, e.g., m semantically coherent subpopulations), we expect ρ_{m+1} to be much larger than ρ_m because any $m + 1$ partition needs to break one sub-graph into two pieces and incur a large conductance. In other words, suppose the graph consists of m clusters, the quantity ρ is characterizing the level of internal connection within each cluster. Furthermore, in many cases we expect ρ_{m+1} to be inverse polynomial in dimension. In the running example of Section 3.1 (where augmentation is adding Gaussian noise), ρ is related to the Cheeger constant or the isoperimetric number of the data manifolds, which in many cases is believed to be at least inverse polynomial in dimension (e.g., see Bobkov et al. (1997) for the Cheeger constant of the Gaussian distribution.) Indeed, in Section 3.4 we will formally lowerbound ρ_{m+1} by the product of the augmentation strength and the Cheeger constant of the subpopulation distributions (Proposition 3.9), and lowerbound the

⁵To see this, consider $3 \leq i \leq |\mathcal{X}|$. Let S_1, \dots, S_i be the partition of \mathcal{X} that minimizes the RHS of Definition 3.4. Define set $S'_{i-1} := S_i \cup S_{i-1}$. It is easy to see that $\phi_G(S'_{i-1}) = \frac{\sum_{x \in S'_{i-1}, x' \notin S'_{i-1}} w_{xx'}}{\sum_{x \in S'_{i-1}} w_x} \leq \frac{\sum_{j=i-1}^i \sum_{x \in S_j, x' \notin S_j} w_{xx'}}{\sum_{j=i-1}^i \sum_{x \in S_j} w_x} \leq \max\{\phi_G(S_{i-1}), \phi_G(S_i)\}$. Notice that $S_1, \dots, S_{i-2}, S'_{i-1}$ are $i-1$ non-empty sets that form a partition of \mathcal{X} , by Definition 3.4 we have $\rho_{i-1} \leq \max\{\phi_G(S_1), \dots, \phi_G(S_{i-2}), \phi_G(S'_{i-1})\} \leq \max\{\phi_G(S_1), \dots, \phi_G(S_i)\} = \rho_i$.

Cheeger constant by inverse polynomial for concrete settings where the data come from a mixture of manifolds (Theorem 3.11).

Assumption 3.5 also implies properties of the graph spectrum. Recall that γ_i is the i -th largest eigenvalue of the normalized adjacency matrix \bar{A} and $\gamma_1 = 1$. According to Cheeger’s inequality (Lemma B.4), Assumption 3.5 implies that $\gamma_{2m} \leq 1 - \Omega(\rho^2/\log m)$, which suggests that there is a gap between γ_1 and γ_{2m} and will be useful in our analysis.

Next, we formalize the assumption that very few edges cross different ground-truth classes. It turns out that it suffices to assume that the labels are recoverable from the augmentations (which is also equivalent to that two examples in different classes can rarely be augmented into the same point).

Assumption 3.6 (Labels are recoverable from augmentations). *Let $\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}$ and $y(\bar{x})$ be its label. Let the augmentation $x \sim \mathcal{A}(\cdot|\bar{x})$. We assume that there exists a classifier g that can predict $y(\bar{x})$ given x with error at most α . That is, $g(x) = y(\bar{x})$ with probability at least $1 - \alpha$.*

A small α in Assumption 3.6 means that different classes are “separated” in the sense that data from different classes have very few (at most $O(\alpha)$) shared augmentations. Alternatively, one can think of this assumption as assuming that the augmentation graph can be partitioned into r clusters each corresponding to augmentations from one class, and there are at most $O(\alpha)$ edges across the clusters. This is typically true for real-world image data like ImageNet, since for any two images from different classes (e.g., images of a Husky and a Birman cat), using the typical data augmentations such as adding noise and random cropping can rarely (with exponentially small probability) lead to the same augmented image.

Typically, both ρ in Assumption 3.5 and α in Assumption 3.7 are small positive values that are much less than 1. However, ρ can be much larger than α . Recall that ρ can be expected to be at least inverse polynomial in dimension. In contrast, α characterizes the separation between classes and are expected to be exponentially small in typical cases. For example, in the running example of Section 3.1 with Gaussian perturbation augmentation, if $\sigma\sqrt{d}$ is smaller than the minimum distance between two subpopulations, we can rarely augment two datapoints from distinct subpopulations into a shared augmentation, and therefore α is expected to be exponentially small. Our analysis below operates in the reasonable regime where ρ^2 is larger than α , which intuitively means that the internal connection within the cluster is bigger than the separation between the clusters.

We also introduce the following assumption which states that some minimizer of the population spectral contrastive loss can be realized by the hypothesis class.

Assumption 3.7 (Expressivity of the hypothesis class). *Let \mathcal{F} be a hypothesis class containing functions from \mathcal{X} to \mathbb{R}^k . We assume that at least one of the global minima of $\mathcal{L}(f)$ belongs to \mathcal{F} .*

Our main theorem bound from above the linear probe error of the feature learned by minimizing the *population* spectral contrastive loss. In Theorem 4.3 we extend this result to the case where both the feature and the linear head are learned from empirical datasets.

Theorem 3.8 (Main theorem for infinite/population pretraining data case). *Assume the representation dimension $k \geq 2r$ and Assumption 3.6 holds for $\alpha > 0$. Let \mathcal{F} be a hypothesis class that satisfies Assumption 3.7 and let $f_{\text{pop}}^* \in \mathcal{F}$ be a minimizer of $\mathcal{L}(f)$. Then, we have*

$$\mathcal{E}(f_{\text{pop}}^*) \leq \tilde{O}\left(\alpha/\rho_{\lfloor k/2 \rfloor}^2\right).$$

In particular, if Assumption 3.5 also holds and $k > 2m$, we have $\mathcal{E}(f_{\text{pop}}^) \leq \tilde{O}(\alpha/\rho^2)$.*

Here we use $\tilde{O}(\cdot)$ to hide universal constant factors and logarithmic factors in k . We note that $\alpha = 0$ when augmentations from different classes are perfectly disconnected in the augmentation graph, in which case the above theorem guarantees the exact recovery of the ground truth. Generally, we expect α to be an extremely (exponentially) small constant independent of k , whereas $\rho_{\lfloor k/2 \rfloor}$ increases with k and can be at least inverse polynomial when k is reasonably large, hence much larger than $\sqrt{\alpha}$. We characterize the ρ_k 's growth on more concrete distributions in the next subsection. When $k > 2m$, as argued below Assumption 3.6, we expect that $\alpha \ll \rho^2 \leq \rho_{m+1}^2$ and thus the error α/ρ^2 is sufficiently small.

Previous works on graph partitioning (Arora et al., 2009, Lee et al., 2014, Leighton and Rao, 1999) often analyze the rounding algorithms that conduct clustering based on the representations of unlabeled data and do not analyze the performance of linear probe (which has access to labeled data). These results provide guarantees on the approximation ratio—the ratio between the conductance of the obtained partition to the best partition—which may depend on graph size (Arora et al., 2009) that can be exponentially large in our setting. The approximation ratio guarantee does not lead to a guarantee on the representations' performance on downstream tasks. Our guarantees are on the linear probe accuracy on the downstream tasks and independent of the graph size. We rely on the formulation of the downstream task's labeling function (Assumption 3.6) as well as a novel analysis technique that characterizes the linear structure of the representations. In Section B, we provide the proof of Theorem 3.8 as well as its more generalized version where $k/2$ is relaxed to be any constant fraction of k . A proof sketch of Theorem 3.8 is given in Section 5.1.

3.4 Provable instantiation of Theorem 3.8 to mixture of manifold data

In this section, we exemplify Theorem 3.8 on examples where the natural data distribution is a mixture of manifolds.

We first show that in the running example given in Section 3.1, Assumption 3.5 holds for some ρ that is closely related to the Cheeger constant of the data manifolds. Recall that the Cheeger constant or isoperimetric number (Buser, 1982) of a distribution μ with density p over \mathbb{R}^d is defined as

$$h_\mu := \inf_{S \subset \mathbb{R}^d} \frac{\int_{\partial S} p(x) dx}{\min\{\int_S p(x) dx, \int_{\mathbb{R}^d \setminus S} p(x) dx\}}. \quad (8)$$

Here the denominator is the smaller one of volumes of S and $\mathbb{R}^d \setminus S$, whereas the numerator is the surface area of S . (See e.g., (Chen, 2021) for the precise definition of the boundary measure ∂S .) The following proposition (proved in Section C.1) shows that ρ scales linearly in the augmentation strength and the Cheeger constant.

Proposition 3.9. *Suppose the natural data distribution $\mathcal{P}_{\bar{\mathcal{X}}}$ is a mixture of m distributions P_1, \dots, P_m supported on disjoint subsets of \mathbb{R}^d , and the data augmentation is Gaussian perturbation sampled from $\mathcal{N}(0, \sigma^2 \cdot I_{d \times d})$. Then,*

$$\lim_{\sigma \rightarrow 0^+} \frac{\rho_{m+1}}{\sigma} \gtrsim \min_{i \in [m]} h_{P_i} \quad (9)$$

That is, ρ_{m+1} is at least linear in the augmentation size σ and the Cheeger constants of subpopulations.

In many cases, the Cheeger constant is at least inverse polynomial in the data dimension (Chen, 2021, Lee and Vempala, 2016). When the manifolds P_i are spherical Gaussian with unit identity covariance, the Cheeger constant is $\Omega(1)$ (Bobkov et al., 1997), and thus the distribution $\mathcal{P}_{\bar{\mathcal{X}}}$ in Proposition 3.9 satisfies Assumption 3.5 with $\rho \gtrsim \sigma$. Furthermore, when the distribution is transformed by

a function with Lipschitzness $\kappa > 0$, the Cheeger constant changes by a factor at most κ . Therefore, Proposition 3.9 also applies to a mixture of manifolds setting defined below.

In the rest of this section, we instantiate Theorem 3.8 on a mixture of manifolds example where the data is generated from a Lipschitz transformation of a mixture of Gaussian distributions, and give an error bound for the downstream classification task.

Example 3.10 (Mixture of manifolds). *Suppose $\mathcal{P}_{\bar{x}}$ is mixture of $r \leq d$ distributions P_1, \dots, P_r , where each P_i is generated by some κ -bi-Lipschitz⁶ generator $Q : \mathbb{R}^{d'} \rightarrow \mathbb{R}^d$ on some latent variable $z \in \mathbb{R}^{d'}$ with $d' \leq d$ which as a mixture of Gaussian distribution:*

$$x \sim P_i \iff x = Q(z), z \sim \mathcal{N}(\mu_i, \frac{1}{d'} \cdot I_{d' \times d'}).$$

Let the data augmentation of a natural data sample \bar{x} be $\bar{x} + \xi$ where $\xi \sim \mathcal{N}(0, \frac{\sigma^2}{d} \cdot I_{d \times d})$ is isotropic Gaussian noise with $0 < \sigma \lesssim \frac{1}{\sqrt{d}}$. We also assume $\min_{i \neq j} \|\mu_i - \mu_j\|_2 \gtrsim \frac{\kappa \cdot \sqrt{\log d}}{\sqrt{d'}}$.

Let $\bar{y}(x)$ be the most likely mixture index i that generates x : $\bar{y}(x) := \arg \max_i P_i(x)$. The simplest downstream task can have label $y(x) = \bar{y}(x)$. More generally, let $r' \leq r$ be the number of labels, and the label $y(x) \in [r']$ in the downstream task be equal to $\pi(\bar{y}(x))$ where π is a function that maps $[r]$ to $[r']$.

We note that the intra-class distance in the latent space is on the scale of $\Omega(1)$, which can be much larger than the distance between class means which is assumed to be $\gtrsim \frac{\kappa \cdot \sqrt{\log d}}{\sqrt{d'}}$. Therefore, distance-based clustering algorithms do not apply. Moreover, in the simple downstream tasks, the label for x could be just the index of the mixture where x comes from. We also allow downstream tasks that merge the r components into r' labels as long as each mixture component gets the same label. We apply Theorem 3.8 and get the following theorem:

Theorem 3.11 (Theorem for the mixture of manifolds example). *When $k \geq 2r + 2$, Example 3.10 satisfies Assumption 3.6 with $\alpha \leq \frac{1}{\text{poly}(d)}$, and has $\rho_{\lfloor k/2 \rfloor} \gtrsim \frac{\sigma}{\kappa \sqrt{d}}$. As a consequence, the error bound is $\mathcal{E}(f_{\text{pop}}^*) \leq \tilde{O}\left(\frac{\kappa^2}{\sigma^2 \cdot \text{poly}(d)}\right)$.*

The theorem above guarantees small error even when σ is polynomially small. In this case, the augmentation noise has a much smaller scale than the data (which is at least on the order of $1/\kappa$). This suggests that contrastive learning can non-trivially leverage the structure of the underlying data and learn good representations with relatively weak augmentation. To the best of our knowledge, it is difficult to apply the theorems in previous works (Arora et al., 2019, Lee et al., 2020, Tosh et al., 2020, 2021, Wei et al., 2020) to this example and get similar guarantees with polynomial dependencies on d, σ, κ . The work of Wei et al. (2020) can apply to the setting where r is known and the downstream label is equal to $\bar{y}(x)$, but cannot handle the case when r is unknown or when two mixture component can have the same label. We refer the reader to the related work section for more discussions and comparisons. The proof can be found in Section C.2.

4 Finite-sample generalization bounds

4.1 Unlabeled sample complexity for pretraining

In Section 3, we provide guarantees for spectral contrastive learning on population data. In this section, we show that these guarantees can be naturally extended to the finite-sample

⁶A κ bi-Lipschitz function satisfies $\frac{1}{\kappa} \|f(x) - f(y)\|_2 \leq \|x - y\|_2 \leq \kappa \|f(x) - f(y)\|_2$.

regime with standard concentration bounds. In particular, given a unlabeled pretraining dataset $\{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_{n_{\text{pre}}}\}$ with $\bar{x}_i \sim \mathcal{P}_{\bar{\mathcal{X}}}$, we learn a feature extractor by minimizing the following *empirical spectral contrastive loss*:

$$\widehat{\mathcal{L}}_{n_{\text{pre}}}(f) := -\frac{2}{n} \sum_{i=1}^{n_{\text{pre}}} \mathbb{E}_{\substack{x \sim \mathcal{A}(\cdot|\bar{x}_i) \\ x^+ \sim \mathcal{A}(\cdot|\bar{x}_i)}} \left[f(x)^\top f(x^+) \right] + \frac{1}{n_{\text{pre}}(n_{\text{pre}} - 1)} \sum_{i \neq j} \mathbb{E}_{\substack{x \sim \mathcal{A}(\cdot|\bar{x}_i) \\ x^- \sim \mathcal{A}(\cdot|\bar{x}_j)}} \left[\left(f(x)^\top f(x^-) \right)^2 \right].$$

It is worth noting that $\widehat{\mathcal{L}}_{n_{\text{pre}}}(f)$ is an unbiased estimator of the population spectral contrastive loss $\mathcal{L}(f)$. (See Claim D.2 for a proof.) Therefore, we can derive generalization bounds via off-the-shelf concentration inequalities. Let \mathcal{F} be a hypothesis class containing feature extractors from \mathcal{X} to \mathbb{R}^k . We extend Rademacher complexity to function classes with high-dimensional outputs and define the Rademacher complexity of \mathcal{F} on n data as $\widehat{\mathcal{R}}_n(\mathcal{F}) := \max_{x_1, \dots, x_n \in \mathcal{X}} \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}, i \in [k]} \frac{1}{n} \left(\sum_{j=1}^n \sigma_j f_i(x_j) \right) \right]$, where σ is a uniform random vector in $\{-1, 1\}^n$ and $f_i(z)$ is the i -th dimension of $f(z)$.

Recall that $f_{\text{pop}}^* \in \mathcal{F}$ is a minimizer of $\mathcal{L}(f)$. The following theorem with proofs in Section D.1 bounds the population loss of a feature extractor trained with finite data:

Theorem 4.1 (Excess contrastive loss). *For some $\kappa > 0$, assume $\|f(x)\|_\infty \leq \kappa$ for all $f \in \mathcal{F}$ and $x \in \mathcal{X}$. Let $f_{\text{pop}}^* \in \mathcal{F}$ be a minimizer of the population loss $\mathcal{L}(f)$. Given a random dataset of size n_{pre} , let $\hat{f}_{\text{emp}} \in \mathcal{F}$ be a minimizer of empirical loss $\widehat{\mathcal{L}}_{n_{\text{pre}}}(f)$. Then, when Assumption 3.7 holds, with probability at least $1 - \delta$ over the randomness of data, we have*

$$\mathcal{L}(\hat{f}_{\text{emp}}) \leq \mathcal{L}(f_{\text{pop}}^*) + c_1 \cdot \widehat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}) + c_2 \cdot \left(\sqrt{\frac{\log 2/\delta}{n_{\text{pre}}}} + \delta \right),$$

where constants $c_1 \lesssim k^2 \kappa^2 + k\kappa$ and $c_2 \lesssim k\kappa^2 + k^2 \kappa^4$.

The Rademacher complexity usually looks like $\widehat{\mathcal{R}}_n(\mathcal{F}) = \sqrt{R/n}$ where R measures the complexity of \mathcal{F} (hence only depends on \mathcal{F}). This suggests that when κ is $O(1)$, the sample complexity for achieving suboptimality ϵ on population loss is $O(k^4 R/\epsilon^2)$. We can apply Theorem 4.1 to any hypothesis class \mathcal{F} of interest (e.g., deep neural networks) and plug in off-the-shelf Rademacher complexity bounds. For instance, in Section D.2 we give a corollary of Theorem 4.1 when \mathcal{F} contains deep neural networks with ReLU activation.

The theorem above shows that we can achieve near-optimal population loss by minimizing empirical loss up to some small excess loss. The following theorem characterizes how the error propagates to the linear probe performance mildly under some spectral gap conditions.

Theorem 4.2 (Minimum downstream error). *In the setting of Theorem 4.1, suppose Assumption 3.5 holds for $\rho > 0$, Assumption 3.6 holds for $\alpha > 0$, Assumption 3.7 holds, and the representation dimension $k \geq \max\{4r + 2, 2m\}$. Then, with $1 - \delta$ probability over the randomness of data, for any $\hat{f}_{\text{emp}} \in \mathcal{F}$ that minimizes the empirical loss $\widehat{\mathcal{L}}_{n_{\text{pre}}}(f)$, we have that*

$$\mathcal{E}(\hat{f}_{\text{emp}}) \lesssim \frac{\alpha}{\rho^2} \cdot \log k + \frac{ck}{\Delta_\gamma^2} \left(\widehat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}) + \sqrt{\frac{\log 2/\delta}{n_{\text{pre}}}} + \delta \right),$$

where $c \lesssim (k\kappa + k\kappa^2 + 1)^2$, and $\Delta_\gamma := \gamma_{\lfloor 3k/4 \rfloor} - \gamma_k$ is the eigenvalue gap between the $\lfloor 3k/4 \rfloor$ -th and the k -th eigenvalue.

This theorem shows that the error on the downstream task only grows linearly with the excess loss during pretraining. Roughly speaking, one can think of Δ_γ as on the order of $1 - \gamma_k$, hence by Cheeger’s inequality it’s larger than ρ^2 . When $\widehat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}) = \sqrt{2R/n_{\text{pre}}}$ and $\kappa \leq O(1)$, we have that the number of unlabeled samples required to achieve ϵ downstream error is $O(m^6 R/\epsilon^2 \rho^4)$. We can relax Assumption 3.7 to approximate realizability in the sense that \mathcal{F} contains some sub-optimal feature extractor under the population spectral loss and pay an additional error term in the linear probe error bound. The proof of Theorem 4.2 can be found in Section D.3.

4.2 Labeled sample complexity for linear probe

In this section, we provide sample complexity analysis for learning a linear probe with *labeled* data. Theorem 3.8 guarantees the existence of a linear probe that achieves a small downstream classification error. However, a priori it is unclear how large the margin of the linear classifier can be, so it is hard to apply margin theory to provide generalization bounds for 0-1 loss. We could in principle control the margin of the linear head, but using capped quadratic loss turns out to suffice and mathematically more convenient. We learn a linear head with the following *capped quadratic loss*: given a tuple $(z, y(\bar{x}))$ where $z \in \mathbb{R}^k$ is a representation of augmented datapoint $x \sim \mathcal{A}(\cdot|\bar{x})$ and $y(\bar{x}) \in [r]$ is the label of \bar{x} , for a linear probe $B \in \mathbb{R}^{k \times r}$ we define loss $\ell((z, y(\bar{x})), B) := \sum_{i=1}^r \min \{ (B^\top z - \vec{y}(\bar{x}))_i^2, 1 \}$, where $\vec{y}(\bar{x})$ is the one-hot embedding of $y(\bar{x})$ as a r -dimensional vector (1 on the $y(\bar{x})$ -th dimension, 0 on other dimensions). This is a standard modification of quadratic loss in statistical learning theory that ensures the boundedness of the loss for the ease of analysis (Mohri et al., 2018).

The following Theorem 4.3 provides a generalization guarantee for the linear classifier that minimizes capped quadratic loss on a labeled downstream dataset of size n_{down} . The key challenge of the proof is showing the existence of a small-norm linear head B that gives small population quadratic loss, which is not obvious from Theorem 4.2 where only small 0-1 error is guaranteed. Given a labeled dataset $\{(\bar{x}_i, y(\bar{x}_i))\}_{i=1}^{n_{\text{down}}}$ where $\bar{x}_i \sim \mathcal{P}_{\bar{\mathcal{X}}}$ and $y(\bar{x}_i)$ is its label, we sample $x_i \sim \mathcal{A}(\cdot|\bar{x}_i)$ for $i \in [n_{\text{down}}]$. Given a norm bound $C_k > 0$, we learn a linear probe \widehat{B} by minimizing the capped quadratic loss subject to a norm constraint:

$$\widehat{B} \in \arg \min_{\|B\|_F \leq C_k} \sum_{i=1}^{n_{\text{down}}} \ell((\hat{f}_{\text{emp}}(x_i), y(\bar{x}_i)), B). \quad (10)$$

Theorem 4.3 (End-to-end error bounds with finite pretraining and downstream samples). *In the setting of Theorem 4.2, choose $C_k > 0$ such that $C_k \geq \frac{2(k+1)}{\gamma_k}$. Then, with probability at least $1 - \delta$ over the randomness of data, for any $\hat{f}_{\text{emp}} \in \mathcal{F}$ that minimizes the empirical pre-training loss $\widehat{\mathcal{L}}_{n_{\text{pre}}}(f)$ and a linear head \widehat{B} learned from Equation (10), we have*

$$\mathcal{E}(\hat{f}_{\text{emp}}, \widehat{B}) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log k + \frac{ck}{\Delta_\gamma^2} \left(\widehat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}) + \sqrt{\frac{\log 2/\delta}{n_{\text{pre}}}} + \delta \right) + \left(rC_k \sqrt{\frac{k}{n_{\text{down}}}} + \sqrt{\frac{\log 1/\delta}{n_{\text{down}}}} \right).$$

Here the first term is an error caused by the property for the population data, which is unavoidable even with infinite pretraining and downstream samples (but it can be small as argued in Section 3.3). The second term is caused by finite pretraining samples, and the third term is caused by finite samples in the linear classification on the downstream task.

Typically, the Rademacher complexity is roughly $\widehat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}) = \sqrt{2R/n_{\text{pre}}}$ where R captures the complexity of the model architecture. Thus, to achieve final linear probe error no more than

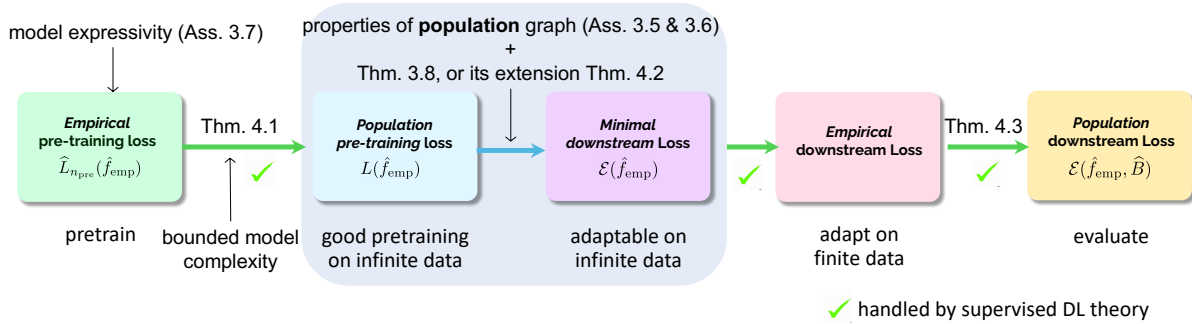


Figure 2: A diagram of our analysis framework. We decompose the problem into a core step that shows a small population pretraining loss implies a small minimal downstream loss (Theorem 3.8, or its extension Theorem 4.2) and a few other somewhat standard steps that link empirical losses to population losses (Theorems 4.1 and Theorem 4.3).

$O(\epsilon)$, we would need to select k such that $\frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log k \leq \epsilon$, and we need $\text{poly}(k, \frac{1}{\Delta_\gamma}, R, \frac{1}{\epsilon})$ pretraining samples and $\text{poly}(k, r, \frac{1}{\gamma_k}, \frac{1}{\epsilon})$ downstream samples.

When $r < k$, the eigengap Δ_γ is on the order of $1 - \gamma_k$ which is larger than ρ^2 by Cheeger inequality. Recall that ρ is at least inverse polynomial in d as argued in Section 3.4, one can expect $\frac{1}{\Delta_\gamma}$ to be at most $\text{poly}(d)$. On the other hand, $\gamma_k \approx 1$ so $\frac{1}{\gamma_k}$ can be thought of as a constant. Thus, the final required number of pretraining samples is $n_{\text{pre}} = \text{poly}(k, d, R, \frac{1}{\epsilon})$ and number of downstream samples is $n_{\text{down}} = \text{poly}(r, k, \frac{1}{\epsilon})$. We note that the downstream sample complexity doesn't depend on the complexity of the hypothesis class R , suggesting that pretraining helps reduce the sample complexity of the supervised downstream task.

The proof of Theorem 4.3 is in Section E.

5 Analysis Framework and Proof Sketch

As discussed before and suggested by the structured of Section 3 and 4, our analysis framework decompose the problem into a key step about the population cases (Section 3) and a few other somewhat standard steps that link empirical losses to population losses (Section 4). As depicted in Figure 2, the core step (Theorem 3.8, or its extension Theorem 4.2) is to show that a small population pretraining loss implies the existence of a linear classifier for the downstream task, that is, a small minimal downstream loss.

We first remark that a feature of our analysis framework is that we link the population pretraining data case to the finite sample case by showing the empirical and population pretraining losses are similar when the feature extractors are a *parameterized* family of models with capacity bounds (the first arrow in Figure 2). Hypothetically, suppose such a connection between population and empirical data case was built through the relationship between the population and empirical graphs, e.g., by proving that the empirical graph has similar spectral properties as the population graph, then the sample complexity will be exponential. Intuitively, this is because the population graph is very sparse, and the empirical graph is with high probability empty if the number of samples is only polynomial in dimension (e.g. consider the case when the augmentation simply adds small perturbation, as in the running example in Section 3.1). The empirical graph essentially follows

the well-studied random geometric graph model (Penrose, 2003), and tends to have no structure in high dimension (Brennan et al., 2020, Bubeck et al., 2016, Liu et al., 2021). The fundamental difference between this hypothetical and our framework is that the empirical graph’s definition does not involve any parameterization, and thus the resemblance between the empirical and population graphs does not leverage the extrapolation (or inductive bias) of the model parameterization as our framework does for the pretraining losses.

We note that the inductive bias of the parameterized model is indeed used in the analysis for finite-sample case. We assume that the model family \mathcal{F} can express the eigenfunctions/eigenvectors of the graph (Assumption 3.7) and also implicitly assume bounds on its Rademacher complexity (in Theorem 4.3).

Once we obtained that the existence of a linear classifier, the remaining steps (the third and fourth arrows in Figure 2) follow from standard supervised learning theory.

In the rest of this section, we will give a proof sketch of the population case, which is the more challenging step.

5.1 Proof Sketch of Theorem 3.8

In this section, we give a proof sketch of Theorem 3.8 in a simplified binary classification setting where there are only two classes in the downstream task.

Recall that N is the size of \mathcal{X} . Recall that w_x is the total weight associated with an augmented datapoint $x \in \mathcal{X}$, which can also be thought of as the probability mass of x as a randomly sampled augmented datapoint. In the scope of this section, for demonstrating the key idea, we also assume that x has uniform distribution, i.e., $w_x = \frac{1}{N}$ for any $x \in \mathcal{X}$.

Let $g : \mathcal{X} \rightarrow \{0, 1\}$ be the Bayes optimal classifier for predicting the label given an augmented datapoint. By Assumption 3.6, g has an error at most α (which is assumed to be small). Thus, we can think of it as the “target” classifier that we aim to recover. We will show that g can be approximated by a linear function on top of the learned features. Recall that v_1, v_2, \dots, v_k are the top- k unit-norm eigenvectors of \bar{A} and the feature u_x^* for x is the x -th row of the eigenvector matrix $F^* = [v_1, v_2, \dots, v_k] \in \mathbb{R}^{N \times k}$. As discussed in Section 3.2, the spectral contrastive loss was designed to compute a variant of the eigenvector matrix F^* up to row scaling and right rotation. More precisely, letting $F_{\text{pop}}^* \in \mathbb{R}^{N \times k}$ be the matrix whose rows contain all the learned embeddings, Section 3.2 shows that $F_{\text{pop}}^* = D \cdot F^* \cdot R$ for a positive diagonal matrix D and an orthonormal matrix R , and Lemma 3.1 shows that these transformations do not affect the feature quality. Therefore, in the rest of the section, it suffices to show that linear models on top of F^* gives the labels of g . Let $\vec{g} \in \{0, 1\}^N$ be the vector that contains the labels of all the data under the optimal g , i.e., $\vec{g}_x = g(x)$. Given a linear head b , note that F^*b gives the prediction (before the threshold function) for all examples. Therefore, it suffices to show the existence of a vector b such that

$$F^*b \approx \vec{g} \tag{11}$$

Let $\mathcal{L} \triangleq I - \bar{A}$ be the normalized Laplacian matrix. Then, v_i ’s are the k smallest unit-norm eigenvectors of \mathcal{L} with eigenvalues $\lambda_i = 1 - \gamma_i$. Elementary derivations can give a well-known, important property of the Laplacian matrix L : the quadratic form $\vec{g}^\top L \vec{g}$ captures the amount of edges across the two groups that are defined by the binary vector \vec{g} (Chung and Graham, 1997, section 1.2):

$$\vec{g}^\top \mathcal{L} \vec{g} = \frac{1}{2} \cdot \sum_{x, x' \in \mathcal{X}} \frac{w_{xx'}}{\sqrt{w_x w_{x'}}} (\vec{g}_x - \vec{g}_{x'})^2 \tag{12}$$

With slight abuse of notation, suppose (x, x^+) is the random variable for a positive pair. Using that w is the density function for the positive pair and the simplification that $w_x = 1/N$, we can rewrite equation (12) as

$$\vec{g}^\top \mathcal{L} \vec{g} = \frac{N}{2} \cdot \mathbb{E}_{x, x^+} [(\vec{g}_x - \vec{g}_{x^+})^2], \quad (13)$$

Note that $\mathbb{E}_{x, x^+} [(\vec{g}_x - \vec{g}_{x^+})^2]$ is the probability that a positive pair have different labels under the Bayes optimal classifier g . Because Assumption 3.6 assumes that the labels can be almost determined by the augmented data, we can show that two augmentations of the same datapoint should rarely produce different labels under the Bayes optimal classifier. We will prove in Lemma B.5 via simple calculation that

$$\vec{g}^\top \mathcal{L} \vec{g} \leq N\alpha \quad (14)$$

(We can sanity-check the special case when $\alpha = 0$, that is, the label is determined by the augmentation. In this case, $g(x) = g(x^+)$ for a positive pair (x, x^+) w.p. 1, which implies $\vec{g}^\top \mathcal{L} \vec{g} = \frac{N}{2} \cdot \mathbb{E}_{x, x^+} [(\vec{g}_x - \vec{g}_{x^+})^2] = 0$.)

Next, we use equation (14) to link \vec{g} to the eigenvectors of L . Let $\lambda_{k+1} \leq \dots \leq \lambda_N$ be the rest of eigenvalues with unit-norm eigenvectors v_{k+1}, \dots, v_N . Let $\Pi \triangleq \sum_{i=1}^k v_i v_i^\top$ and $\Pi_\perp \triangleq \sum_{i=k+1}^N v_i v_i^\top$ be the projection operators onto the subspaces spanned by the first k and the last $N-k$ eigenvectors, respectively. Equation (14) implies that \vec{g} has limited projection to the subspace of Π_\perp :

$$N\alpha \geq \vec{g}^\top \mathcal{L} \vec{g} = (\Pi \vec{g} + \Pi_\perp \vec{g})^\top \mathcal{L} (\Pi \vec{g} + \Pi_\perp \vec{g}) \geq (\Pi_\perp \vec{g})^\top \mathcal{L} (\Pi_\perp \vec{g}) \geq \lambda_{k+1} \|\Pi_\perp \vec{g}\|_2^2, \quad (15)$$

where the first inequality follows from dropping the $\|(\Pi \vec{g})^\top \mathcal{L} \Pi \vec{g}\|_2^2$ and using $\Pi_\perp L \Pi = 0$, and the second inequality is because that Π_\perp only contains eigenvectors with eigenvalue at least λ_{k+1} .

Note that $\Pi \vec{g}$ is in the span of eigenvectors v_1, \dots, v_k , that is, the column-span of F^* . Therefore, there exists $b \in \mathbb{R}^k$ such that $\Pi \vec{g} = F^* b$. As a consequence,

$$\|\vec{g} - F^* b\|_2^2 = \|\Pi_\perp \vec{g}\|_2^2 \leq \frac{N\alpha}{\lambda_{k+1}} \quad (16)$$

By higher-order Cheeger inequality (see Lemma B.4), we have that $\lambda_{k+1} \gtrsim \rho_{\lceil k/2 \rceil}^2$. Then, we obtain the mean-squared error bound:

$$\frac{1}{N} \|\vec{g} - F^* b\|_2^2 \leq \alpha / \rho_{\lceil k/2 \rceil}^2 \quad (17)$$

The steps above demonstrate the gist of the proofs, which are formalized in more generality in Section B.1. We will also need two minor steps to complete the proof of Theorem 3.8. First, we can convert the mean-squared error bound to classification error bound: because $F^* b$ is close to the binary vector \vec{g} in mean-squared error, $\mathbb{1}[F^* b > 1/2]$ is close to \vec{g} in 0-1 error. (See Claim B.9 for the formal argument.) Next, $F^* b$ only gives the prediction of the model given the augmented datapoint. We will show in Section B.2 that averaging the predictions on the augmentations of a data point will not increase the classification error.

6 Experiments

We test spectral contrastive learning on benchmark vision datasets. We minimize the empirical spectral contrastive loss with an encoder network f and sample fresh augmentation in each iteration. The pseudo-code for the algorithm and more implementation details can be found in Section A.

Encoder / feature extractor. The encoder f contains three components: a backbone network, a projection MLP and a projection function. The backbone network is a standard ResNet architecture. The projection MLP is a fully connected network with BN applied to each layer, and ReLU activation applied to each except for the last layer. The projection function takes a vector and projects it to a sphere ball with radius $\sqrt{\mu}$, where $\mu > 0$ is a hyperparameter that we tune in experiments. We find that using a projection MLP and a projection function improves the performance.

Linear evaluation protocol. Given the pre-trained encoder network, we follow the standard linear evaluation protocol (Chen and He, 2020) and train a supervised linear classifier on frozen representations, which are from the ResNet’s global average pooling layer.

Results. We report the accuracy on CIFAR-10/100 (Krizhevsky and Hinton, 2009) and Tiny-ImageNet (Le and Yang, 2015) in Table 1. Our empirical results show that spectral contrastive learning achieves better performance than two popular baseline algorithms SimCLR (Chen et al., 2020a) and SimSiam (Chen and He, 2020). In Table 2 we report results on ImageNet (Deng et al., 2009) dataset, and show that our algorithm achieves similar performance as other state-of-the-art methods. We note that our algorithm is much more principled than previous methods and doesn’t rely on large batch sizes (SimCLR (Chen et al., 2020a)), momentum encoders (BYOL (Grill et al., 2020) and MoCo (He et al., 2020)) or additional tricks such as stop-gradient (SimSiam (Chen and He, 2020)).

Datasets	CIFAR-10			CIFAR-100			Tiny-ImageNet		
Epochs	200	400	800	200	400	800	200	400	800
SimCLR (repro.)	83.73	87.72	90.60	54.74	61.05	63.88	43.30	46.46	48.12
SimSiam (repro.)	87.54	90.31	91.40	61.56	64.96	65.87	34.82	39.46	46.76
Ours	88.66	90.17	92.07	62.45	65.82	66.18	41.30	45.36	49.86

Table 1: Top-1 accuracy under linear evaluation protocol.

	SimCLR	BYOL	MoCo v2	SimSiam	Ours
acc. (%)	66.5	66.5	67.4	68.1	66.97

Table 2: ImageNet linear evaluation accuracy with 100-epoch pre-training. All results but ours are reported from (Chen and He, 2020). We use batch size 384 during pre-training.

7 Conclusion

In this paper, we present a novel theoretical framework of self-supervised learning and provide provable guarantees for the learned representation on downstream linear classification tasks. We hope the framework could facilitate future theoretical analyses of self-supervised pretraining losses and inspire new methods. It does not capture the potential implicit bias of optimizers but does take into account the inductive bias of the models. By abstracting away the effect of optimization, we can focus on the effect of pretraining losses and their interaction with the structure of the population data. Future directions may include designing better pretraining losses and analyzing more fine-grained properties of the learned representations (e.g., as in recent follow-up works (HaoChen et al., 2022, Shen et al., 2022)), by potentially leveraging more advanced techniques from spectral graph theory.

Acknowledgements

We thank Margalit Glasgow, Ananya Kumar, Jason D. Lee, Sang Michael Xie, and Guodong Zhang for helpful discussions. CW acknowledges support from an NSF Graduate Research Fellowship. TM acknowledges support of Google Faculty Award and NSF IIS 2045685. We also acknowledge the support of HAI and the Google Cloud. Toyota Research Institute ("TRI") provided funds to assist the authors with their research but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.

References

- Emmanuel Abbe. Community detection and stochastic block models: recent developments, 2017.
- Sanjeev Arora, Satish Rao, and Umesh Vazirani. Expander flows, geometric embeddings and graph partitioning. *Journal of the ACM (JACM)*, 56(2):1–37, 2009.
- Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *arXiv preprint arXiv:1902.09229*, 2019.
- Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. *arXiv preprint arXiv:1906.00910*, 2019.
- Maria-Florina Balcan, Avrim Blum, and Ke Yang. Co-training and expansion: Towards bridging theory and practice. *Advances in neural information processing systems*, 17:89–96, 2005.
- Yamini Bansal, Gal Kaplun, and Boaz Barak. For self-supervised learning, rationality implies generalization, provably. *arXiv preprint arXiv:2010.08508*, 2020.
- Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100, 1998.
- Sergey G Bobkov et al. An isoperimetric inequality on the discrete cube, and an elementary proof of the isoperimetric inequality in gauss space. *The Annals of Probability*, 25(1):206–214, 1997.
- Matthew Brennan, Guy Bresler, and Dheeraj Nagaraj. Phase transitions for detecting latent geometry in random graphs. *Probability Theory and Related Fields*, 178(3):1215–1289, 2020.
- Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. Signature verification using a " siamese" time delay neural network. *Advances in neural information processing systems*, 6:737–744, 1993.
- Sébastien Bubeck, Jian Ding, Ronen Eldan, and Miklós Z Rácz. Testing for high-dimensional geometry in random graphs. *Random Structures & Algorithms*, 49(3):503–532, 2016.
- Daniel Bump. *Automorphic forms and representations*. Number 55. Cambridge university press, 1998.
- Peter Buser. A note on the isoperimetric constant. In *Annales scientifiques de l'École normale supérieure*, volume 15, pages 213–230, 1982.

- Tianle Cai, Ruiqi Gao, Jason D Lee, and Qi Lei. A theory of label propagation for subpopulation shift. *arXiv preprint arXiv:2102.11203*, 2021.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020.
- Jeff Cheeger. A lower bound for the smallest eigenvalue of the laplacian. In *Proceedings of the Princeton conference in honor of Professor S. Bochner*, pages 195–199, 1969.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020a.
- Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020b.
- Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. *arXiv preprint arXiv:2011.10566*, 2020.
- Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020c.
- Yuansi Chen. An almost constant lower bound of the isoperimetric coefficient in the kls conjecture. *Geometric and Functional Analysis*, 31(1):34–61, 2021.
- Fan RK Chung and Fan Chung Graham. *Spectral graph theory*. Number 92. American Mathematical Soc., 1997.
- Tri Dao, Albert Gu, Alexander Ratner, Virginia Smith, Chris De Sa, and Christopher Ré. A kernel theory of modern data augmentation. In *International Conference on Machine Learning*, pages 1528–1537. PMLR, 2019.
- Sanjoy Dasgupta, Michael L Littman, and David McAllester. Pac generalization bounds for co-training. *Advances in neural information processing systems*, 1:375–382, 2002.
- J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- Luc Devroye, Abbas Mehrabian, and Tommy Reddad. The total variation distance between high-dimensional gaussians. *arXiv preprint arXiv:1810.08693*, 2018.
- Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, pages 297–299. PMLR, 2018.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.
- Heinrich Walter Guggenheimer. *Applicable Geometry: Global and Local Convexity*. RE Krieger Publishing Company, 1977.

- Jeff Z HaoChen, Colin Wei, Ananya Kumar, and Tengyu Ma. Beyond separability: Analyzing the linear transferability of contrastive representations to related subpopulations. *arXiv preprint arXiv:2204.02683*, 2022.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.
- Olivier Henaff. Data-efficient image recognition with contrastive predictive coding. In *International Conference on Machine Learning*, pages 4182–4192. PMLR, 2020.
- R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2018.
- Ravi Kannan, Santosh Vempala, and Adrian Vetta. On clusterings: Good, bad and spectral. *Journal of the ACM (JACM)*, 51(3):497–515, 2004.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7:7, 2015.
- James R Lee, Shayan Oveis Gharan, and Luca Trevisan. Multiway spectral partitioning and higher-order cheeger inequalities. *Journal of the ACM (JACM)*, 61(6):1–30, 2014.
- Jason D Lee, Qi Lei, Nikunj Saunshi, and Jiacheng Zhuo. Predicting what you already know helps: Provable self-supervised learning. *arXiv preprint arXiv:2008.01064*, 2020.
- Yin Tat Lee and Santosh S Vempala. Eldan’s stochastic localization and the kls conjecture: Isoperimetry, concentration and mixing. *arXiv preprint arXiv:1612.01507*, 2016.
- Jing Lei, Alessandro Rinaldo, et al. Consistency of spectral clustering in stochastic block models. *Annals of Statistics*, 43(1):215–237, 2015.
- Tom Leighton and Satish Rao. Multicommodity max-flow min-cut theorems and their use in designing approximation algorithms. *Journal of the ACM (JACM)*, 46(6):787–832, 1999.
- Siqi Liu, Sidhanth Mohanty, Tselil Schramm, and Elizabeth Yang. Testing thresholds for high-dimensional sparse random geometric graphs. *arXiv preprint arXiv:2111.11316*, 2021.
- Anand Louis and Konstantin Makarychev. Approximation algorithm for sparsest k-partitioning. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages 1244–1255. SIAM, 2014.
- Anand Louis, Prasad Raghavendra, Prasad Tetali, and Santosh Vempala. Algorithmic extensions of cheeger’s inequality to higher eigenvalues and partitions. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 315–326. Springer, 2011.
- Frank McSherry. Spectral partitioning of random graphs. In *Proceedings 42nd IEEE Symposium on Foundations of Computer Science*, pages 529–537. IEEE, 2001.
- Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020.

- Jovana Mitrovic, Brian McWilliams, Jacob Walker, Lars Buesing, and Charles Blundell. Representation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*, 2020.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- Boaz Nadler, Nathan Srebro, and Xueyuan Zhou. Semi-supervised learning with the graph laplacian: The limit of infinite unlabelled data. *Advances in neural information processing systems*, 22: 1330–1338, 2009.
- Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 14:849–856, 2001.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Mathew Penrose. *Random geometric graphs*, volume 5. OUP Oxford, 2003.
- Geoffrey Schiebinger, Martin J Wainwright, and Bin Yu. The geometry of kernelized spectral clustering. *The Annals of Statistics*, 43(2):819–846, 2015.
- Kendrick Shen, Robbie Jones, Ananya Kumar, Sang Michael Xie, Jeff Z HaoChen, Tengyu Ma, and Percy Liang. Connect, not collapse: Explaining contrastive learning for unsupervised domain adaptation. *arXiv preprint arXiv:2204.00570*, 2022.
- Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 1857–1865, 2016.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.
- Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning. *arXiv preprint arXiv:2005.10243*, 2020a.
- Yuandong Tian, Lantao Yu, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning with dual deep networks. *arXiv preprint arXiv:2010.00578*, 2020b.
- Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive estimation reveals topic posterior information to linear models. *arXiv:2003.02234*, 2020.
- Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive learning, multi-view redundancy, and linear models. In *Algorithmic Learning Theory*, pages 1179–1206. PMLR, 2021.
- Yao-Hung Hubert Tsai, Yue Wu, Ruslan Salakhutdinov, and Louis-Philippe Morency. Self-supervised learning from a multi-view perspective. *arXiv preprint arXiv:2006.05576*, 2020.
- Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR, 2020.
- Colin Wei, Kendrick Shen, Yining Chen, and Tengyu Ma. Theoretical analysis of self-training with deep networks on unlabeled data. *arXiv preprint arXiv:2010.03622*, 2020.

- Wikipedia contributors. Hilbert–schmidt integral operator — Wikipedia, the free encyclopedia, 2020. URL https://en.wikipedia.org/w/index.php?title=Hilbert%E2%80%9993Schmidt_integral_operator&oldid=986771357. [Online; accessed 21-July-2021].
- Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. Unsupervised data augmentation for consistency training. *arXiv preprint arXiv:1904.12848*, 2019.
- Mang Ye, Xu Zhang, Pong C Yuen, and Shih-Fu Chang. Unsupervised embedding learning via invariant and spreading instance feature. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6210–6219, 2019.
- Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.
- Richard Zemel and Miguel Carreira-Perpiñán. Proximity graphs for clustering and manifold learning. *Advances in neural information processing systems*, 17, 2004.
- Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pages 912–919, 2003.
- Roland S Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel. Contrastive learning inverts the data generating process. *arXiv preprint arXiv:2102.08850*, 2021.

A Experiment details

The pseudo-code for our empirical algorithm is summarized in Algorithm 1.

Algorithm 1 Spectral Contrastive Learning

Require: batch size N , structure of encoder network f

- 1: **for** sampled minibatch $\{\bar{x}_i\}_{i=1}^N$ **do**
 - 2: **for** $i \in \{1, \dots, N\}$ **do**
 - 3: draw two augmentations $x_i = \text{aug}(\bar{x}_i)$ and $x'_i = \text{aug}(\bar{x}_i)$.
 - 4: compute $z_i = f(x_i)$ and $z'_i = f(x'_i)$.
 - 5: compute loss $\mathcal{L} = -\frac{2}{N} \sum_{i=1}^N z_i^\top z'_i + \frac{1}{N(N-1)} \sum_{i \neq j} (z_i^\top z'_j)^2$
 - 6: update f to minimize \mathcal{L}
 - 7: **return** encoder network $f(\cdot)$
-

Our results with different hyperparameters on CIFAR-10/100 and Tiny-ImageNet are listed in Table 3.

Datasets	CIFAR-10			CIFAR-100			Tiny-ImageNet		
Epochs	200	400	800	200	400	800	200	400	800
SimCLR (repro.)	83.73	87.72	90.60	54.74	61.05	63.88	43.30	46.46	48.12
SimSiam (repro.)	87.54	90.31	91.40	61.56	64.96	65.87	34.82	39.46	46.76
Ours ($\mu = 1$)	86.47	89.90	92.07	59.13	63.83	65.52	28.76	33.94	40.82
Ours ($\mu = 3$)	87.72	90.09	91.84	61.05	64.79	66.18	40.06	42.52	49.86
Ours ($\mu = 10$)	88.66	90.17	91.01	62.45	65.82	65.16	41.30	45.36	47.84

Table 3: Top-1 accuracy under linear evaluation protocol.

Additional details about the encoder. For the backbone network, we use the CIFAR variant of ResNet18 for CIFAR-10 and CIFAR-100 experiments and use ResNet50 for Tiny-ImageNet and ImageNet experiments. For the projection MLP, we use a 2-layer MLP with hidden and output dimensions 1000 for CIFAR-10, CIFAR100, and Tiny-ImageNet experiments. We use a 3-layer MLP with hidden and output dimension 8192 for ImageNet experiments. We set $\mu = 10$ in the ImageNet experiment, and set $\mu \in \{1, 3, 10\}$ for the CIFAR-10/100 and Tiny-ImageNet experiments.

Training the encoder. We train the neural network using SGD with momentum 0.9. The learning rate starts at 0.05 and decreases to 0 with a cosine schedule. On CIFAR-10/100 and Tiny-ImageNet we use weight decay 0.0005 and train for 800 epochs with batch size 512. On ImageNet we use weight decay 0.0001 and train for 100 epochs with batch size 384. We use 1 GTX 1080 GPU for CIFAR-10/100 and Tiny-ImageNet experiments, and use 8 GTX 1080 GPUs for ImageNet experiments.

Linear evaluation protocol. We train the linear head using SGD with batch size 256 and weight decay 0 for 100 epochs, learning rate starts at 30.0 and is decayed by 10x at the 60th and 80th epochs.

Image transformation details. We use the same augmentation strategy as described in (Chen and He, 2020).

B Proofs for Section 3

We first prove a more generalized version of Theorem 3.8 in section B.1, and then prove Theorem 3.8 in Section B.2.

B.1 A generalized version of Theorem 3.8

For the proof we will follow the convention in literature (Lee et al., 2014) and define the *normalized Laplacian matrix* as follows:

Definition B.1. Let $G = (\mathcal{X}, w)$ be the augmentation graph defined in Section 3.1. The normalized Laplacian matrix of the graph is defined as $\mathcal{L} = I - D^{-1/2}AD^{-1/2}$, where A is the adjacency matrix with $A_{xx'} = w_{xx'}$ and D is a diagonal matrix with $D_{xx} = w_x$.

It is easy to see that $\mathcal{L} = I - \bar{A}$ where \bar{A} is the normalized adjacency matrix defined in Section 3.1. Therefore, when λ_i is the i -th smallest eigenvalue of \mathcal{L} , $1 - \lambda_i$ is the i -th largest eigenvalue of \bar{A} .

We call a function defined on augmented data $\hat{y} : \mathcal{X} \rightarrow [r]$ an *extended labeling function*. Given an extended labeling function, we define the following quantity that describes the difference between extended labels of two augmented data of the same natural datapoint:

$$\phi^{\hat{y}} := \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[\hat{y}(x) \neq \hat{y}(x')]. \quad (18)$$

We also define the following quantity that describes the difference between extended label of an augmented datapoint and the ground truth label of the corresponding natural datapoint:

$$\Delta(y, \hat{y}) := \Pr_{x \sim \mathcal{P}_{\bar{\mathcal{X}}}, \tilde{x} \sim \mathcal{A}(\cdot|x)} (\hat{y}(\tilde{x}) \neq y(x)). \quad (19)$$

Recall the spectral contrastive loss defined in Section 3.2 is:

$$\mathcal{L}(f) := \mathbb{E}_{\substack{x_1 \sim \mathcal{P}_{\bar{\mathcal{X}}}, x_2 \sim \mathcal{P}_{\bar{\mathcal{X}}}, \\ x \sim \mathcal{A}(\cdot|x_1), x^+ \sim \mathcal{A}(\cdot|x_1), x' \sim \mathcal{A}(\cdot|x_2)}} \left[-2 \cdot f(x)^\top f(x^+) + \left(f(x)^\top f(x') \right)^2 \right].$$

We first state a more general version of Theorem 3.8 as follows.

Theorem B.2. Assume the set of augmented data \mathcal{X} is finite. Let $f_{\text{pop}}^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathbb{R}^k} \mathcal{L}(f)$ be a minimizer of the population spectral contrastive loss $\mathcal{L}(f)$ with $k \in \mathcal{Z}^+$. Let $k' \geq r$ such that $k + 1 = (1 + \zeta)k'$, where $\zeta \in (0, 1)$ and $k' \in \mathcal{Z}^+$. Then, there exists a linear probe $B^* \in \mathbb{R}^{r \times k}$ and a universal constant c such that the linear probe predictor satisfies

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\left\| \vec{y}(\bar{x}) - B^* f_{\text{pop}}^*(x) \right\|_2^2 \right] \leq c \cdot \left(\text{poly}(1/\zeta) \cdot \log(k + 1) \cdot \frac{\phi^{\hat{y}}}{\rho_{k'}^2} + \Delta(y, \hat{y}) \right),$$

where $\vec{y}(\bar{x})$ is the one-hot embedding of $y(\bar{x})$ and $\rho_{k'}$ is the sparsest m -partition defined in Definition 3.4. Furthermore, the error of the linear probe predictor can be bounded by

$$\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{f_{\text{pop}}^*, B^*}(x) \neq y(\bar{x}) \right) \leq 2c \cdot \left(\text{poly}(1/\zeta) \cdot \log(k + 1) \cdot \frac{\phi^{\hat{y}}}{\rho_{k'}^2} + \Delta(y, \hat{y}) \right).$$

Also, if we let λ_i be the i -th smallest eigenvalue of the normalized Laplacian matrix of the graph of the augmented data, we can find a matrix B^* satisfying the above equations with norm bound $\|B^*\|_F \leq 1/(1 - \lambda_k)$.

We provide the proof for Theorem B.2 below.

Let $\lambda_1, \lambda_2, \dots, \lambda_k, \lambda_{k+1}$ be the $k+1$ smallest eigenvalues of the Laplacian matrix L . The following theorem gives a theoretical guarantee similar to Theorem B.2 except for that the bound depends on λ_{k+1} :

Theorem B.3. *Assume the set of augmented data \mathcal{X} is finite. Let $f_{\text{pop}}^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathbb{R}^k}$ be a minimizer of the population spectral contrastive loss $\mathcal{L}(f)$ with $k \in \mathcal{Z}^+$. Then, for any labeling function $\hat{y} : \mathcal{X} \rightarrow [r]$ there exists a linear probe $B^* \in \mathbb{R}^{r \times k}$ with norm $\|B^*\|_F \leq 1/(1 - \lambda_k)$ such that*

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\left\| \vec{y}(\bar{x}) - B^* f_{\text{pop}}^*(x) \right\|_2^2 \right] \leq \frac{\phi^{\hat{y}}}{\lambda_{k+1}} + 4\Delta(y, \hat{y}),$$

where $\vec{y}(\bar{x})$ is the one-hot embedding of $y(\bar{x})$. Furthermore, the error can be bounded by

$$\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{f_{\text{pop}}^*, B^*}(x) \neq y(\bar{x}) \right) \leq \frac{2\phi^{\hat{y}}}{\lambda_{k+1}} + 8\Delta(y, \hat{y}).$$

We defer the proof of Theorem B.3 to Section B.3.

To get rid of the dependency on λ_{k+1} , we use following higher-order Cheeger's inequality from (Louis and Makarychev, 2014).

Lemma B.4 (Proposition 1.2 in (Louis and Makarychev, 2014)). *Let $G = (V, w)$ be a weight graph with $|V| = N$. Then, for any $t \in [N]$ and $\zeta > 0$ such that $(1 + \zeta)t \in [N]$, there exists a partition S_1, S_2, \dots, S_t of V with*

$$\phi_G(S_i) \lesssim \text{poly}(1/\zeta) \sqrt{\lambda_{(1+\zeta)t} \log t},$$

where $\phi_G(\cdot)$ is the Dirichlet conductance defined in Definition 3.3.

Now we prove Theorem B.2 by combining Theorem B.3 and Lemma B.4.

Proof of Theorem B.2. Let $G = (\mathcal{X}, w)$ be the augmentation graph. In Lemma B.4 let $(1 + \zeta)t = k + 1$ and $t = k'$ we have: there exists partition $S_1, \dots, S_{k'} \subset \mathcal{X}$ such that $\phi_G(S_i) \lesssim \text{poly}(1/\zeta) \sqrt{\lambda_{k+1} \log(k+1)}$ for $\forall i \in [k']$. By Definition 3.4, we have $\rho_{k'} \leq \max_{i \in [k']} \phi_G(S_i) \lesssim \text{poly}(1/\zeta) \sqrt{\lambda_{k+1} \log(k+1)}$, which leads to $\frac{1}{\lambda_{k+1}} \lesssim \text{poly}(1/\zeta) \cdot \log(k+1) \cdot \frac{1}{\rho_{k'}^2}$. Plugging this bound to Theorem B.3 finishes the proof. \square

B.2 Proof of Theorem 3.8

We will use the following lemma which gives a connection between $\phi^{\hat{y}}, \Delta(y, \hat{y})$ and Assumption 3.6.

Lemma B.5. *Let $G = (\mathcal{X}, w)$ be the augmentation graph, r be the number of underlying classes. Let S_1, S_2, \dots, S_r be the partition induced by the classifier g in Assumption 3.6. Then, there exists an extended labeling function \hat{y} such that*

$$\Delta(y, \hat{y}) \leq \alpha$$

and

$$\phi^{\hat{y}} = \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[\hat{y}(x) \neq \hat{y}(x')] \leq 2\alpha.$$

Proof of Lemma B.5. We define function $\hat{y} : \mathcal{X} \rightarrow [r]$ as follows: for an augmented data $x \in \mathcal{X}$, we use function $\hat{y}(x)$ to represent the index of set that x is in, i.e., $x \in S_{\hat{y}(x)}$. By Assumption 3.6 it is easy to see $\Delta(y, \hat{y}) \leq \alpha$. On the other hand, we have

$$\begin{aligned}
\phi^{\hat{y}} &= \sum_{x, x' \in \mathcal{X}} w_{xx'} \mathbb{1} [\hat{y}(x) \neq \hat{y}(x')] \\
&= \sum_{x, x' \in \mathcal{X}} \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}} [\mathcal{A}(x|\bar{x})\mathcal{A}(x'|\bar{x}) \cdot \mathbb{1} [\hat{y}(x) \neq \hat{y}(x')]] \\
&\leq \sum_{x, x' \in \mathcal{X}} \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}} [\mathcal{A}(x|\bar{x})\mathcal{A}(x'|\bar{x}) \cdot (\mathbb{1} [\hat{y}(x) \neq y(\bar{x})] + \mathbb{1} [\hat{y}(x') \neq y(\bar{x})])] \\
&= 2 \cdot \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}} [\mathcal{A}(x|\bar{x}) \cdot \mathbb{1} [\hat{y}(x) \neq y(\bar{x})]] \\
&= 2 \cdot \Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} (x \notin S_{y(\bar{x})}) = 2\alpha.
\end{aligned}$$

Here the inequality is because when $\hat{y}(x) \neq \hat{y}(x')$, there must be $\hat{y}(x) \neq y(\bar{x})$ or $\hat{y}(x') \neq y(\bar{x})$. \square

Now we give the proof of Theorem 3.8 using Lemma B.5 and Theorem B.2.

Proof of Theorem 3.8. Let S_1, S_2, \dots, S_r be the partition of \mathcal{X} induced by the classifier g given in Assumption 3.6. Define function $\hat{y} : \mathcal{X} \rightarrow [r]$ as follows: for an augmented datapoint $x \in \mathcal{X}$, we use function $\hat{y}(x)$ to represent the index of set that x is in, i.e., $x \in S_{\hat{y}(x)}$. Let $k' = \lfloor \frac{k}{2} \rfloor$ in Theorem B.2, we have $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{f_{\text{pop}}, B^*}^*(x) \neq y(\bar{x})) \lesssim \log(k) \cdot \frac{\phi^{\hat{y}}}{\rho_{\lfloor k/2 \rfloor}^2} + \Delta(y, \hat{y})$. By Lemma B.5 we have $\phi^{\hat{y}} \leq 2\alpha$ and $\Delta(y, \hat{y}) \leq \alpha$, so we have $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{f_{\text{pop}}, B^*}^*(x) \neq y(\bar{x})) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log(k)$. Notice that by definition of ensembled linear probe predictor, $\bar{g}_{f_{\text{pop}}, B^*}(\bar{x}) \neq y(\bar{x})$ happens only if more than half of the augmentations of \bar{x} predicts differently from $y(\bar{x})$, so we have $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}} (\bar{g}_{f_{\text{pop}}, B^*}(\bar{x}) \neq y(\bar{x})) \leq 2 \Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{f_{\text{pop}}, B^*}^*(x) \neq y(\bar{x})) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log(k)$. \square

B.3 Proof of Theorem B.3

The proof of Theorem B.3 contains two steps. First, we show that when the feature extractor is composed of the minimal eigenvectors of the normalized Laplacian matrix L , we can achieve good linear probe accuracy. Then we show that minimizing $\mathcal{L}(f)$ gives us a feature extractor equally good as the eigenvectors.

For the first step, we use the following lemma which shows that the smallest eigenvectors of \mathcal{L} can approximate any function on \mathcal{X} up to an error proportional to the Rayleigh quotient of the function.

Lemma B.6. *Let \mathcal{L} be the normalized Laplacian matrix of some graph G . Let $N = |\mathcal{X}|$ be total number of augmented data, v_i be the i -th smallest unit-norm eigenvector of \mathcal{L} with eigenvalue λ_i (make them orthogonal in case of repeated eigenvalues). Let $R(u) := \frac{u^\top \mathcal{L} u}{u^\top u}$ be the Rayleigh quotient of a vector $u \in \mathbb{R}^N$. Then, for any $k \in \mathbb{Z}^+$ such that $k < N$ and $\lambda_{k+1} > 0$, there exists a vector $b \in \mathbb{R}^k$ with norm $\|b\|_2 \leq \|u\|_2$ such that*

$$\left\| u - \sum_{i=1}^k b_i v_i \right\|_2^2 \leq \frac{R(u)}{\lambda_{k+1}} \|u\|_2^2.$$

Proof of Lemma B.6. We can decompose the vector u in the eigenvector basis as:

$$u = \sum_{i=1}^N \zeta_i v_i.$$

We have

$$R(u) = \frac{\sum_{i=1}^N \lambda_i \zeta_i^2}{\|u\|_2^2}.$$

Let $b \in \mathbb{R}^k$ be the vector such that $b_i = \zeta_i$. Obviously we have $\|b\|_2^2 \leq \|u\|_2^2$. Noticing that

$$\left\| u - \sum_{i=1}^k b_i v_i \right\|_2^2 = \sum_{i=k+1}^N \zeta_i^2 \leq \frac{R(u)}{\lambda_{k+1}} \|u\|_2^2,$$

which finishes the proof. \square

We also need the following claim about the Rayleigh quotient $R(u)$ when u is a vector defined by an extended labeling function \hat{y} .

Claim B.7. *In the setting of Lemma B.6, let \hat{y} be an extended labeling function. Fix $i \in [r]$. Define function $u_i^{\hat{y}}(x) := \sqrt{w_x} \cdot \mathbb{1}[\hat{y}(x) = i]$ and $u_i^{\hat{y}}$ is the corresponding vector in \mathbb{R}^N . Also define the following quantity:*

$$\phi_i^{\hat{y}} := \frac{\sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[(\hat{y}(x) = i \wedge \hat{y}(x') \neq i) \text{ or } (\hat{y}(x) \neq i \wedge \hat{y}(x') = i)]}{\sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\hat{y}(x) = i]}.$$

Then, we have

$$R(u_i^{\hat{y}}) = \frac{1}{2} \phi_i^{\hat{y}}.$$

Proof of Claim B.7. Let f be any function $\mathcal{X} \rightarrow \mathbb{R}$, define function $u(x) := \sqrt{w_x} \cdot f(x)$. Let $u \in \mathbb{R}^N$ be the vector corresponding to u . Let A be the adjacency matrix with $A_{xx'} = w_{xx'}$ and D be the diagonal matrix with $D_{xx} = w_x$. By definition of Laplacian matrix, we have

$$\begin{aligned} u^\top \mathcal{L}u &= \|u\|_2^2 - u^\top D^{-1/2} A D^{-1/2} u \\ &= \sum_{x \in \mathcal{X}} w_x f(x)^2 - \sum_{x, x' \in \mathcal{X}} w_{xx'} f(x) f(x') \\ &= \frac{1}{2} \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot (f(x) - f(x'))^2. \end{aligned}$$

Therefore we have

$$\begin{aligned} R(u) &= \frac{u^\top \mathcal{L}u}{u^\top u} \\ &= \frac{1}{2} \cdot \frac{\sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot (f(x) - f(x'))^2}{\sum_{x \in \mathcal{X}} w_x \cdot f(x)^2}. \end{aligned}$$

Setting $f(x) = \mathbb{1}[\hat{y}(x) = i]$ finishes the proof. \square

To see the connection between the feature extractor minimizing the population spectral contrastive loss $\mathcal{L}(f)$ and the feature extractor corresponding to eigenvectors of the Laplacian matrix, we use the following lemma which states that the minimizer of the matrix approximation loss defined in Section 3.2 is equivalent to the minimizer of population spectral contrastive loss up to a data-wise scaling.

Lemma B.8. Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be a feature extractor, matrix $F \in \mathbb{R}^{N \times k}$ be such that its x -th row is $\sqrt{w_x} \cdot f(x)$. Then, F is a minimizer of $\mathcal{L}_{\text{mf}}(F)$ if and only if f is a minimizer of the population spectral contrastive loss $\mathcal{L}(f)$.

Proof of Lemma B.8. Notice that

$$\begin{aligned} \mathcal{L}_{\text{mf}}(F) &= \left\| (I - \mathcal{L}) - FF^\top \right\|_F^2 \\ &= \sum_{x, x' \in \mathcal{X}} \left(\frac{w_{xx'}}{\sqrt{w_x w_{x'}}} - \sqrt{w_x w_{x'}} f(x)^\top f(x') \right)^2 \\ &= \sum_{x, x' \in \mathcal{X}} w_x w_{x'} \left(f(x)^\top f(x') \right)^2 - 2 \sum_{x, x' \in \mathcal{X}} w_{xx'} f(x)^\top f(x') + \|I - \mathcal{L}\|_F^2. \end{aligned} \quad (20)$$

Recall that the definition of spectral contrastive loss is

$$\mathcal{L}(f) := -2 \cdot \mathbb{E}_{x, x^+} \left[f(x)^\top f(x^+) \right] + \mathbb{E}_{x, x^-} \left[\left(f(x)^\top f(x^-) \right)^2 \right],$$

where (x, x^+) is a random positive pair, (x, x^-) is a random negative pair. We can rewrite the spectral contrastive loss as

$$\mathcal{L}(f) = -2 \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot f(x)^\top f(x') + \sum_{x, x' \in \mathcal{X}} w_x w_{x'} \cdot \left(f(x)^\top f(x') \right)^2. \quad (21)$$

Compare Equation (20) and Equation (21), we see they only differ by a constant, which finishes the proof. \square

Note that the minimizer of matrix approximation loss is exactly the largest eigenvectors of $I - L$ (also the smallest eigenvectors of L) due to Eckart–Young–Mirsky theorem, Lemma B.8 indicates that the minimizer of $\mathcal{L}(f)$ is equivalent to the smallest eigenvectors of \mathcal{L} up to data-wise scaling.

The following claim shows the relationship between quadratic loss and prediction error.

Claim B.9. Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be a feature extractor, $B \in \mathbb{R}^{k \times k}$ be a linear head. Let $g_{f, B}$ be the predictor defined in Section 3. Then, for any $x \in \mathcal{X}$ and label $y \in [k]$, we have

$$\|\vec{y} - Bf(x)\|_2^2 \geq \frac{1}{2} \cdot \mathbb{1}[y \neq g_{f, B}(x)],$$

where \vec{y} is the one-hot embedding of y .

Proof. When $y \neq g_{f, B}(x)$, by the definition of $g_{f, B}$ we know that there exists another $y' \neq y$ such that $(Bf(x))_{y'} \geq (Bf(x))_y$. In this case,

$$\|\vec{y} - Bf(x)\|_2^2 \geq (1 - (Bf(x))_y)^2 + (Bf(x))_{y'}^2 \quad (22)$$

$$\geq \frac{1}{2} (1 - (Bf(x))_y + (Bf(x))_{y'})^2 \quad (23)$$

$$\geq \frac{1}{2}, \quad (24)$$

where the first inequality is by omitting all the dimensions in the ℓ_2 norm other than the y -th and y' -th dimensions, the second inequality is by Jensen's inequality, and the third inequality is because $(Bf(x))_{y'} \geq (Bf(x))_y$. This proves the inequality in the claim when $y \neq g_{f, B}(x)$. Finally, we finish the proof by noticing that the inequality in this claim obviously holds when $y = g_{f, B}(x)$. \square

Now we are ready to prove Theorem B.3 by combining Lemma B.6, Claim B.7, Lemma B.8 and Claim B.9.

Proof of Theorem B.3. Let $F_{\text{sc}} = [v_1, v_2, \dots, v_k]$ be the matrix that contains the smallest k eigenvectors of \mathcal{L} as columns. For each $i \in [r]$, we define function $u_i^{\hat{y}}(x) := \sqrt{w_x} \cdot \mathbb{1}[\hat{y}(x) = i]$ and $u_i^{\hat{y}}$ be the corresponding vector in \mathbb{R}^N . By Lemma B.6, there exists a vector $b_i \in \mathbb{R}^k$ with norm bound $\|b_i\|_2 \leq \|u_i^{\hat{y}}\|_2$ such that

$$\|u_i^{\hat{y}} - F_{\text{sc}} b_i\|_2^2 \leq \frac{R(u_i^{\hat{y}})}{\lambda_{k+1}} \|u_i^{\hat{y}}\|_2^2. \quad (25)$$

By Claim B.7, we have

$$R(u_i^{\hat{y}}) = \frac{1}{2} \phi_i^{\hat{y}} = \frac{1}{2} \cdot \frac{\sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[(\hat{y}(x) = i \wedge \hat{y}(x') \neq i) \text{ or } (\hat{y}(x) \neq i \wedge \hat{y}(x') = i)]}{\sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\hat{y}(x) = i]}.$$

So we can rewrite Equation (25) as:

$$\begin{aligned} \|u_i^{\hat{y}} - F_{\text{sc}} b_i\|_2^2 &\leq \frac{\phi_i^{\hat{y}}}{2\lambda_{k+1}} \cdot \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\hat{y}(x) = i] \\ &= \frac{1}{2\lambda_{k+1}} \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[(\hat{y}(x) = i \wedge \hat{y}(x') \neq i) \text{ or } (\hat{y}(x) \neq i \wedge \hat{y}(x') = i)]. \end{aligned} \quad (26)$$

Let matrix $U = [u_1^{\hat{y}}, \dots, u_r^{\hat{y}}]$ contains all $u_i^{\hat{y}}$ as columns, and let $u : \mathcal{X} \rightarrow \mathbb{R}^r$ be the corresponding feature extractor. Define matrix $B \in \mathbb{R}^{N \times k}$ such that $B^\top = [b_1, \dots, b_r]$. Summing Equation (26) over all $i \in [r]$ and by the definition of $\phi^{\hat{y}}$ we have

$$\|U - F_{\text{sc}} B^\top\|_F^2 \leq \frac{1}{2\lambda_{k+1}} \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[\hat{y}(x) \neq \hat{y}(x')] = \frac{\phi^{\hat{y}}}{2\lambda_{k+1}}, \quad (27)$$

where

$$\|B\|_F^2 = \sum_{i=1}^r \|b_i\|_2^2 \leq \sum_{i=1}^r \|u_i^{\hat{y}}\|_2^2 = \sum_{x \in \mathcal{X}} w_x = 1.$$

Now we come back to the feature extractor f_{pop}^* that minimizes the spectral contrastive loss function $\mathcal{L}(f)$. By Lemma B.8, matrix F^* that contains $\sqrt{w_x} \cdot f_{\text{pop}}^*(x)$ as its x -th row is a minimizer of $\mathcal{L}_{\text{mf}}(F)$. By Eckard-Young-Mirsky theorem, we have

$$F^* = F_{\text{sc}} D_\lambda Q,$$

where Q is an orthonormal matrix and

$$D_\lambda = \begin{bmatrix} \sqrt{1 - \lambda_1} & & & \\ & \sqrt{1 - \lambda_2} & & \\ & & \dots & \\ & & & \sqrt{1 - \lambda_k} \end{bmatrix}.$$

Let

$$B^* = B D_\lambda^{-1} Q^{-1},$$

and let $\vec{y}(\bar{x})$ be the one-hot embedding of $y(\bar{x})$, $\vec{\hat{y}}(x)$ be the one-hot embedding of $\hat{y}(x)$, we have

$$\begin{aligned}
& \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\left\| \vec{y}(\bar{x}) - B^* f_{\text{pop}}^*(x) \right\|_2^2 \right] \\
& \leq 2 \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\left\| \vec{\hat{y}}(x) - B^* f_{\text{pop}}^*(x) \right\|_2^2 \right] + 2 \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\left\| \vec{\hat{y}}(x) - \vec{y}(\bar{x}) \right\|_2^2 \right] \\
& = 2 \sum_{x \in \mathcal{X}} w_x \cdot \left\| \vec{\hat{y}}(x) - B^* f_{\text{pop}}^*(x) \right\|_2^2 + 4\Delta(y, \hat{y}) \quad (\text{because } w_x \text{ is the probability of } x) \\
& = 2 \left\| U - F^* B^{*\top} \right\|_F^2 + 4\Delta(y, \hat{y}) \quad (\text{rewrite in matrix form}) \\
& = 2 \left\| U - F_{\text{sc}} B^\top \right\|_F^2 + 4\Delta(y, \hat{y}) \quad (\text{by definition of } B^*) \\
& \leq \frac{\phi^{\hat{y}}}{\lambda_{k+1}} + 4\Delta(y, \hat{y}). \quad (\text{by Equation (27)})
\end{aligned}$$

To bound the error rate, we first notice that Claim B.9 tells us that for any $x \in \mathcal{X}$,

$$\left\| \vec{y}(\bar{x}) - B^* f_{\text{pop}}^*(x) \right\|_2^2 \geq \frac{1}{2} \cdot \mathbb{1} \left[g_{f_{\text{pop}}^*, B^*}(x) \neq y(\bar{x}) \right]. \quad (28)$$

Now we bound the error rate on \mathcal{X} as follows:

$$\begin{aligned}
& \Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{f_{\text{pop}}^*, B^*}(x) \neq y(\bar{x}) \right) \\
& \leq 2 \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\left\| \vec{y}(\bar{x}) - B^* f_{\text{pop}}^*(x) \right\|_2^2 \right] \quad (\text{by Equation (28)}) \\
& \leq \frac{2\phi^{\hat{y}}}{\lambda_{k+1}} + 8\Delta(y, \hat{y}).
\end{aligned}$$

Finally we bound the norm of B^* as

$$\|B^*\|_F^2 = \text{Tr} \left(B^* B^{*\top} \right) = \text{Tr} \left(B D_\lambda^{-2} B^\top \right) \leq \frac{1}{1 - \lambda_k} \|B\|_F^2 = \frac{1}{1 - \lambda_k}.$$

□

C Proofs for Section 3.4

C.1 Proof of Proposition 3.9

Proof of Proposition 3.9. Let B_σ be the uniform distribution over a ball with radius σ . Let S_1, S_2, \dots, S_{m+1} be a partition of the Euclidean space. There must be some $i \in [m+1]$ such that $\Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i] \leq \frac{1}{2}$ for all $j \in [m]$. Thus, we know that

$$\rho_{m+1} \geq \phi_G(S_i) \geq \min_{j \in [m]} \frac{\Pr_{x \sim P_j, \xi \sim B_\sigma, \xi' \sim B_\sigma} [x + \xi \in S_i \wedge x + \xi' \notin S_i]}{\Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i]}. \quad (29)$$

For $x \in \mathbb{R}^d$, we use $P(S_i|x)$ as a shorthand for $\Pr_{\xi \sim B_\sigma} (x + \xi \in S_i)$. Let

$$R := \left\{ x \mid \Pr(S_i|x) \geq \frac{2}{3} \right\}. \quad (30)$$

On one hand, suppose $\int_{x \notin R} P_j(x)P(S_i|x)dx \geq \frac{1}{2} \Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i]$, we can lower bound the numerator in the RHS of Equation (29) as

$$\Pr_{x \sim P_j, \xi \sim B_\sigma, \xi' \sim B_\sigma} [x + \xi \in S_i \wedge x + \xi' \notin S_i] \geq \int_{x \notin R} P_j(x)P(S_i|x)(1 - P(S_i|x))dx \quad (31)$$

$$\geq \frac{1}{6} \Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i], \quad (32)$$

hence the RHS of Equation (29) is at least $1/6$.

On the other hand, suppose $\int_{x \notin R} P_j(x)P(S_i|x)dx < \frac{1}{2} \Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i]$, we have

$$\int_{x \in R} P_j(x)P(S_i|x)dx = \Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i] - \int_{x \notin R} P_j(x)P(S_i|x)dx \geq \frac{1}{2} \Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i], \quad (33)$$

hence the denominator of the RHS of Equation (29) can be upper bounded by

$$\Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i] \leq 2 \int_{x \in R} P(S_i|x)P_j(x)dx \leq 2 \int_{x \in R} P_j(x)dx. \quad (34)$$

Define

$$N(R) := \left\{ x \mid \|x - a\|_2 \leq \frac{\sigma}{6} \text{ for some } a \in R \right\}. \quad (35)$$

For two Gaussian distributions with variance $\sigma^2 \cdot \mathcal{I}_{d \times d}$ and centers at most $\frac{\sigma}{6}$ far from each other, their TV-distance is at most $\frac{1}{6}$ (see the first equation on Page 5 of (Devroye et al., 2018)), hence for any $x \in N(R)$, we have $P(S_i|x) \geq \frac{2}{3} - \frac{1}{6} = \frac{1}{2}$. We can now lower bound the numerator in the RHS of Equation (29) as:

$$\begin{aligned} \Pr_{x \sim P_j, \xi \sim B_\sigma, \xi' \sim B_\sigma} [x + \xi \in S_i \wedge x + \xi' \notin S_i] &\geq \int_{x \in N(R) \setminus R} P_j(x)P(S_i|x)(1 - P(S_i|x))dx \\ &\geq \frac{1}{6} \int_{x \in N(R) \setminus R} P_j(x)dx. \end{aligned} \quad (36)$$

Notice that $\Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i] \leq \frac{1}{2}$ and by the definition of R , we know $\int_{x \in R} P_j(x)dx \leq \frac{3}{4}$, thus

$$\int_{x \in R} P_j(x)dx \leq \frac{3}{4} \leq 3 \int_{x \notin R} P_j(x)dx. \quad (37)$$

Combine Equation (34), Equation (36) and Equation (37) gives:

$$\frac{\Pr_{x \sim P_j, \xi \sim B_\sigma, \xi' \sim B_\sigma} [x + \xi \in S_i \wedge x + \xi' \notin S_i]}{\Pr_{x \sim P_j, \xi \sim B_\sigma} [x + \xi \in S_i]} \geq \frac{1}{36} \frac{\int_{x \in N(R) \setminus R} P_j(x)P(S_i|x)dx}{\min\{\int_{x \in R} P_j(x)dx, \int_{x \in R} P_j(x)dx\}}. \quad (38)$$

Notice that (using the definition of surface area (Guggenheimer, 1977, chapter 4))

$$\lim_{\sigma \rightarrow 0^+} \frac{1}{\sigma} \cdot \frac{\int_{x \in N(R) \setminus R} P_j(x)P(S_i|x)dx}{\min\{\int_{x \in R} P_j(x)dx, \int_{x \in R} P_j(x)dx\}} \geq \frac{1}{6} h_{P_j}, \quad (39)$$

we have that as $\sigma \rightarrow 0^+$,

$$\frac{\rho_{m+1}}{\sigma} \geq \frac{1}{216} \min_{j \in [m]} h_{P_j}, \quad (40)$$

which finishes the proof. \square

C.2 Proof of Theorem 3.11

In this section, we give a proof of Theorem 3.11.

The following lemma shows that the augmented graph for Example 3.10 satisfies Assumption 3.6 with some bounded α .

Lemma C.1. *In the setting of Theorem 3.11, the data distribution satisfies Assumption 3.6 with $\alpha \leq \frac{1}{\text{poly}(d')}$.*

Proof of Lemma C.1. For any $z \sim \mathcal{N}(\mu_i, \frac{1}{d'} \cdot I_{d' \times d'})$ and any $j \neq i$, by the tail bound of gaussian distribution we have

$$\Pr_{z \sim \mathcal{N}(\mu_i, \frac{1}{d'} \cdot I_{d' \times d'})} \left((z - \mu_i)^\top \left(\frac{\mu_j - \mu_i}{\|\mu_j - \mu_i\|_2} \right) \lesssim \frac{\sqrt{\log d}}{\sqrt{d'}} \right) \geq 1 - \frac{1}{\text{poly}(d)}.$$

Also, for $\xi \sim \mathcal{N}(0, \frac{1}{d} \cdot I_{d \times d})$, when $\sigma \leq \frac{1}{\sqrt{d}}$ we have

$$\Pr_{\xi \sim \mathcal{N}(0, \frac{\sigma^2}{d} \cdot I_{d' \times d'})} \left(\|\xi\|_2 \lesssim \frac{\sqrt{\log d}}{\sqrt{d}} \right) \geq 1 - \frac{1}{\text{poly}(d)}.$$

Notice that $\|Q^{-1}(Q(z) + \xi) - z\|_2 \leq \kappa \|\xi\|$, we can set $\|\mu_i - \mu_j\| \gtrsim \kappa \frac{\sqrt{\log d}}{\sqrt{d}}$. Therefore, when $\|\mu_i - \mu_j\| \gtrsim \kappa \frac{\sqrt{\log d}}{\sqrt{d}}$ we can combine the above two cases and have

$$\Pr_{z \sim \mathcal{N}(\mu_i, \frac{1}{d'} \cdot I_{d' \times d'}), \xi \sim \mathcal{N}(0, \frac{\sigma^2}{d} \cdot I_{d' \times d'})} (P_i(z) > P_j(Q^{-1}(Q(z) + \xi))) \geq 1 - \frac{1}{\text{poly}(d)}.$$

Since $r \leq d$, we have

$$\Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot | \bar{x})} (y(x) \neq y(\bar{x})) \geq 1 - \frac{1}{\text{poly}(d)}.$$

□

We use the following lemma to give a lower bound for the sparsest m -partition of the augmentation graph in Example 3.10.

Lemma C.2. *In the setting of Theorem 3.11, for any $k' > r$ and $\tau > 0$, we have*

$$\rho_{k'} \geq \frac{c_{\tau/\kappa}}{18} \cdot \exp \left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d} \right),$$

where

$$c_{\sigma} := \sigma \cdot \Phi_d^{-1} \left(\frac{2}{3} \right)$$

with $\Phi_d(z) := \Pr_{\xi \sim \mathcal{N}(0, \frac{1}{d} I_{d \times d})} (\|\xi\|_2 \leq z)$, and

$$c_{\tau/\kappa} := \min_{p \in [0, \frac{3}{4}]} \frac{\Phi(\Phi^{-1}(p) + \tau\sqrt{d}/\kappa)}{p} - 1$$

with $\Phi(z) := \int_{-\infty}^z \frac{e^{-u^2/2}}{\sqrt{2\pi}} du$.

The proof of Lemma C.2 can be found in Section C.3. Now we give the proof of Example 3.11.

Proof of Theorem 3.11. The result on α is directly from Lemma C.1. By concentration inequality, there must exist some universal constant $C > 0$ such that for any $d \geq C$, we have $1 - \Phi_d(\sqrt{\frac{3}{2}}) \leq \frac{1}{3}$. When this happens, we have $\Phi_d^{-1}(\frac{2}{3}) \leq \sqrt{\frac{3}{2}}$. Since for $d \leq C$ we can just treat d as constant, we have $\Phi_d^{-1}(\frac{2}{3}) \lesssim 1$. Set $\tau = \sigma/d$ in Lemma C.2, we have $\rho_{k'} \gtrsim \frac{\sigma}{\kappa\sqrt{d}}$. Set $k' = \lfloor k/2 \rfloor$, we apply Theorem 3.8 and get the bound we need. \square

C.3 Proof of Lemma C.2

In this section we give a proof for Lemma C.2. We first introduce the following claim which states that for a given subset of augmented data, any two data close in L_2 norm cannot have a very different chance of being augmented into this set.

Claim C.3. *In the setting of Theorem 3.11, given a set $S \subseteq \mathbb{R}^d$. If $x \in \mathbb{R}^d$ satisfies $\Pr(S|x) := \Pr_{\tilde{x} \sim \mathcal{A}(\cdot|x)}(\tilde{x} \in S) \geq \frac{2}{3}$. Then, for any x' such that $\|x - x'\|_2 \leq \tau$, we have*

$$\Pr(S|x') \geq \frac{1}{3} \cdot \exp\left(-\frac{2c_\sigma\tau + \tau^2}{2\sigma^2}\right),$$

where

$$c_\sigma := \sigma \cdot \Phi_d^{-1}\left(\frac{2}{3}\right),$$

with $\Phi_d(z) := \Pr_{\xi \sim \mathcal{N}(0, \frac{1}{d} \cdot I_{d \times d})}(\|\xi\|_2 \leq z)$.

Proof of Claim C.3. By the definition of augmentation, we know

$$\Pr(S|x) = \mathbb{E}_{\xi \sim \mathcal{N}(0, \frac{\sigma^2}{d} \cdot I_{d \times d})}[\mathbb{1}[x + \xi \in S]].$$

By the definition of c_σ , we have

$$\Pr_{\xi \sim \mathcal{N}(0, \frac{\sigma^2}{d} \cdot I_{d \times d})}(\|\xi\|_2 \leq c_\sigma) = \frac{2}{3}.$$

Since $\Pr(S|x) \geq \frac{2}{3}$ by assumption, we have

$$\mathbb{E}_{\xi \sim \mathcal{N}(0, \frac{\sigma^2}{d} \cdot I_{d \times d})}[P(S|x + \xi) \cdot \mathbb{1}[\|\xi\|_2 \leq c_\sigma]] \geq \frac{1}{3}.$$

Now we can bound the quantity of our interest:

$$\begin{aligned} \Pr(S|x') &= \frac{1}{(2\pi\sigma^2/d)^{d/2}} \int_{\xi} e^{-\frac{\|\xi\|_2^2}{2\sigma^2/d}} P(S|x' + \xi) d\xi \\ &= \frac{1}{(2\pi\sigma^2/d)^{d/2}} \int_{\xi} e^{-\frac{\|\xi+x-x'\|_2^2}{2\sigma^2/d}} P(S|x + \xi) d\xi \\ &\geq \frac{1}{(2\pi\sigma^2/d)^{d/2}} \int_{\xi} e^{-\frac{\|\xi+x-x'\|_2^2}{2\sigma^2/d}} P(S|x + \xi) \cdot \mathbb{1}[\|\xi\|_2 \leq c_\sigma] d\xi \\ &\geq \frac{1}{(2\pi\sigma^2/d)^{d/2}} \int_{\xi} e^{-\frac{2c_\sigma\tau + \tau^2 + \|\xi\|_2^2}{2\sigma^2/d}} P(S|x + \xi) \cdot \mathbb{1}[\|\xi\|_2 \leq c_\sigma] d\xi \\ &= e^{-\frac{2c_\sigma\tau + \tau^2}{2\sigma^2/d}} \cdot \mathbb{E}_{\xi \sim \mathcal{N}(0, \frac{\sigma^2}{d} \cdot I_{d \times d})}[P(S|x + \xi) \cdot \mathbb{1}[\|\xi\|_2 \leq c_\sigma]] \\ &\geq \frac{1}{3} \cdot \exp\left(-\frac{2c_\sigma\tau + \tau^2}{2\sigma^2/d}\right). \end{aligned}$$

\square

We now give the proof of Lemma C.2.

Proof of Lemma C.2. Let $S_1, \dots, S_{k'}$ be the disjoint sets that gives $\rho_{k'}$ in Definition 3.4. First we notice that when $k' > r$, there must exist $t \in [k']$ such that for all $i \in [r]$, we have

$$\Pr_{x \sim P_i, \tilde{x} \sim \mathcal{A}(\cdot|x)} (\tilde{x} \in S_t) \leq \frac{1}{2}. \quad (41)$$

WLOG, we assume $t = 1$. So we know that

$$\rho_{k'} = \max_{i \in [k']} \phi_G(S_i) \geq \phi_G(S_1) \geq \min_{j \in [r]} \frac{\mathbb{E}_{x \sim P_j} [\Pr(S_1|x)(1 - \Pr(S_1|x))]}{\mathbb{E}_{x \sim P_j} [\Pr(S_1|x)]}, \quad (42)$$

where

$$\Pr(S|x) := \Pr_{\tilde{x} \sim \mathcal{A}(\cdot|x)} (\tilde{x} \in S).$$

WLOG, we assume $j = 1$ minimizes the RHS of Equation (42), so we only need to prove

$$\frac{\mathbb{E}_{x \sim P_1} [\Pr(S_1|x)(1 - \Pr(S_1|x))]}{\mathbb{E}_{x \sim P_1} [\Pr(S_1|x)]} \geq \frac{c_{\tau/\kappa}}{18} \cdot \exp\left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d}\right).$$

We define the following set

$$R := \left\{ x \mid \Pr(S_1|x) \geq \frac{2}{3} \right\}.$$

Notice that

$$\begin{aligned} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)] &= \int_x P_1(x) \Pr(S_1|x) dx \\ &= \int_{x \in R} P_1(x) \Pr(S_1|x) dx + \int_{x \notin R} P_1(x) \Pr(S_1|x) dx. \end{aligned} \quad (43)$$

We can consider the following two cases.

Case 1: $\int_{x \notin R} P_1(x) \Pr(S_1|x) dx \geq \frac{1}{2} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)]$.

This is the easy case because we have

$$\begin{aligned} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)(1 - \Pr(S_1|x))] &\geq \int_{x \notin R} P_1(x) \Pr(S_1|x)(1 - \Pr(S_1|x)) dx \\ &\geq \frac{1}{3} \int_{x \notin R} P_1(x) \Pr(S_1|x) dx \\ &\geq \frac{1}{6} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)]. \end{aligned}$$

Case 2: $\int_{x \in R} P_1(x) \Pr(S_1|x) dx \geq \frac{1}{2} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)]$.

Define neighbourhood of R as

$$N(R) := \left\{ x \mid \|x - a\|_2 \leq \tau \text{ for some } a \in R \right\}.$$

We have

$$\begin{aligned} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)(1 - \Pr(S_1|x))] &\geq \int_{x \in N(R) \setminus R} P_1(x) \Pr(S_1|x)(1 - \Pr(S_1|x)) dx \\ &\geq \frac{1}{9} \cdot \exp\left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d}\right) \cdot \int_{x \in N(R) \setminus R} P_1(x) dx, \end{aligned}$$

where the second inequality is by Claim C.3. Notice that

$$\int_{x \in R} P_1(x) dx \leq \frac{3}{2} \int_{x \in R} P_1(x) \Pr(S_1|x) dx \leq \frac{3}{2} \int_x P_1(x) \Pr(S_1|x) dx \leq \frac{3}{4},$$

where we use Equation (41). Define set $\tilde{R} := Q^{-1}(R)$ be the set in the ambient space corresponding to R . Define

$$\tilde{N}(\tilde{R}) := \left\{ x' \in \mathbb{R}^d \mid \|x' - a\|_2 \leq \frac{\tau}{\kappa} \text{ for some } a \in \tilde{R} \right\}$$

Due to Q being κ -bi-lipschitz, it is easy to see $\tilde{N}(\tilde{R}) \subseteq Q^{-1}(N(R))$. According to the Gaussian isoperimetric inequality (Bobkov et al., 1997), we have

$$\int_{x \in N(R) \setminus R} P_1(x) dx \geq c_{\tau/\kappa} \int_{x \in R} P_1(x) dx,$$

where

$$c_{\tau/\kappa} := \min_{0 \leq p \leq 3/4} \frac{\Phi(\Phi^{-1}(p) + \tau\sqrt{d}/\kappa) - p}{p} - 1,$$

with $\Phi(\cdot)$ is the Gaussian CDF function defined as

$$\Phi(z) := \int_{-\infty}^z \frac{e^{-u^2/2}}{\sqrt{2\pi}} du.$$

So we have

$$\begin{aligned} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)(1 - \Pr(S_1|x))] &\geq \frac{c_{\tau/\kappa}}{9} \cdot \exp\left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d}\right) \cdot \int_{x \in R} P_1(x) dx \\ &\geq \frac{c_{\tau/\kappa}}{9} \cdot \exp\left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d}\right) \cdot \int_{x \in R} P_1(x) \Pr(S_1|x) dx \\ &\geq \frac{c_{\tau/\kappa}}{18} \cdot \exp\left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d}\right) \cdot \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)]. \end{aligned}$$

By Equation (43), either case 1 or case 2 holds. Combining case 1 and case 2, we have

$$\begin{aligned} \frac{\mathbb{E}_{x \sim P_1} [\Pr(S_1|x)(1 - \Pr(S_1|x))]}{\mathbb{E}_{x \sim P_1} [\Pr(S_1|x)]} &\geq \min \left\{ \frac{1}{6}, \frac{c_{\tau/\kappa}}{18} \cdot \exp\left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d}\right) \right\} \\ &= \frac{c_{\tau/\kappa}}{18} \cdot \exp\left(-\frac{2c_{\sigma}\tau + \tau^2}{2\sigma^2/d}\right). \end{aligned}$$

□

D Proofs for Section 4

D.1 Proof of Theorem 4.1

We restate the empirical spectral contrastive loss defined in Section 4 as follows:

Definition D.1 (Empirical spectral contrastive loss). *Consider a dataset $\hat{\mathcal{X}} = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\}$ containing n data points i.i.d. sampled from $\mathcal{P}_{\bar{\mathcal{X}}}$. Let $\hat{\mathcal{P}}_{\mathcal{X}}$ be the uniform distribution over $\hat{\mathcal{X}}$. Let $\hat{P}_{\bar{x}, \bar{x}'}$ be the uniform distribution over data pairs (\bar{x}_i, \bar{x}_j) where $i \neq j$. We define the empirical spectral contrastive loss of a feature extractor f as*

$$\hat{\mathcal{L}}_n(f) := -2\mathbb{E}_{\substack{\bar{x} \sim \hat{\mathcal{P}}_{\mathcal{X}}, \\ x \sim \mathcal{A}(\cdot|\bar{x}), x' \sim \mathcal{A}(\cdot|\bar{x})}} \left[f(x)^\top f(x') \right] + \mathbb{E}_{\substack{(\bar{x}, \bar{x}') \sim \hat{P}_{\bar{x}, \bar{x}'}, \\ x \sim \mathcal{A}(\cdot|\bar{x}), x' \sim \mathcal{A}(\cdot|\bar{x}')}} \left[\left(f(x)^\top f(x') \right)^2 \right].$$

The following claim shows that $\widehat{\mathcal{L}}_n(f)$ is an unbiased estimator of population spectral contrastive loss.

Claim D.2. $\widehat{\mathcal{L}}_n(f)$ is an unbiased estimator of $\mathcal{L}(f)$, i.e.,

$$\mathbb{E}_{\widehat{\mathcal{X}}} \left[\widehat{\mathcal{L}}_n(f) \right] = \mathcal{L}(f).$$

Proof. This is because

$$\begin{aligned} \mathbb{E}_{\widehat{\mathcal{X}}} \left[\widehat{\mathcal{L}}_n(f) \right] &= -2 \cdot \mathbb{E}_{\widehat{\mathcal{X}}} \left[\mathbb{E}_{\substack{\bar{x} \sim \widehat{\mathcal{P}}_{\widehat{\mathcal{X}}}, \\ x \sim \mathcal{A}(\cdot|\bar{x}), x' \sim \mathcal{A}(\cdot|\bar{x})}} \left[f(x)^\top f(x') \right] \right] + \mathbb{E}_{\widehat{\mathcal{X}}} \left[\mathbb{E}_{\substack{(\bar{x}, \bar{x}') \sim \widehat{P}_{\bar{x}, \bar{x}'}, \\ x \sim \mathcal{A}(\cdot|\bar{x}), x' \sim \mathcal{A}(\cdot|\bar{x}')}} \left[\left(f(x)^\top f(x') \right)^2 \right] \right] \\ &= -2 \mathbb{E}_{\substack{\bar{x} \sim \mathcal{P}_{\widehat{\mathcal{X}}}, \\ x \sim \mathcal{A}(\cdot|\bar{x}), x' \sim \mathcal{A}(\cdot|\bar{x})}} \left[f(x)^\top f(x') \right] + \mathbb{E}_{\substack{\bar{x} \sim \mathcal{P}_{\widehat{\mathcal{X}}}, \bar{x}' \sim \mathcal{P}_{\widehat{\mathcal{X}}}, \\ x \sim \mathcal{A}(\cdot|\bar{x}), x' \sim \mathcal{A}(\cdot|\bar{x}')}} \left[\left(f(x)^\top f(x') \right)^2 \right] = \mathcal{L}(f). \end{aligned}$$

□

To make use of the Radmacher complexity theory, we need to write the empirical loss as the sum of i.i.d. terms, which is achieved by the following sub-sampling scheme:

Definition D.3. Given dataset $\widehat{\mathcal{X}}$, we sample a subset of tuples as follows: first sample a permutation $\pi : [n] \rightarrow [n]$, then we sample tuples $S = \{(z_i, z_i^+, z_i')\}_{i=1}^{n/2}$ as follows:

$$\begin{aligned} z_i &\sim \mathcal{A}(\cdot|\bar{x}_{\pi(2i-1)}), \\ z_i^+ &\sim \mathcal{A}(\cdot|\bar{x}_{\pi(2i-1)}), \\ z_i' &\sim \mathcal{A}(\cdot|\bar{x}_{\pi(2i)}). \end{aligned}$$

We define the following loss on S :

$$\widehat{\mathcal{L}}_S(f) := \frac{1}{n/2} \sum_{i=1}^{n/2} \left[\left(f(z_i)^\top f(z_i') \right)^2 - 2f(z_i)^\top f(z_i^+) \right].$$

It is easy to see that $\widehat{\mathcal{L}}_S(f)$ is an unbiased estimator of $\widehat{\mathcal{L}}_n(f)$:

Claim D.4. For given $\widehat{\mathcal{X}}$, if we sample S as above, we have:

$$\mathbb{E}_S \left[\widehat{\mathcal{L}}_S(f) \right] = \widehat{\mathcal{L}}_n(f).$$

Proof. This is obvious by the definition of $\widehat{\mathcal{L}}_S(f)$ and $\widehat{\mathcal{L}}_n(f)$. □

The following lemma reveals the relationship between the Rademacher complexity of feature extractors and the Rademacher complexity of the loss defined on tuples:

Lemma D.5. Let \mathcal{F} be a hypothesis class of feature extractors from \mathcal{X} to \mathbb{R}^k . Assume $\|f(x)\|_\infty \leq \kappa$ for all $x \in \mathcal{X}$. For $i \in [k]$, define $f_i : \mathcal{X} \rightarrow \mathbb{R}$ be the function such that $f_i(x)$ is the i -th dimension of $f(x)$. Let \mathcal{F}_i be the hypothesis containing f_i for all $f \in \mathcal{F}$. For $n \in \mathbb{Z}^+$, let $\widehat{\mathcal{R}}_n(\mathcal{F}_i)$ be the maximal possible empirical Rademacher complexity of \mathcal{F}_i over n data:

$$\widehat{\mathcal{R}}_n(\mathcal{F}_i) := \max_{\{x_1, x_2, \dots, x_n\}} \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f_i(x_j) \right) \right],$$

where x_1, x_2, \dots, x_m are in \mathcal{X} , and σ is a uniform random vector in $\{-1, 1\}^n$. Then, the empirical Rademacher complexity on any n tuples $\{(z_i, z_i^+, z_i^-)\}_{i=1}^n$ can be bounded by

$$\mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j \left(\left(f(z_j)^\top f(z_j') \right)^2 - 2f(z_j)^\top f(z_j^+) \right) \right) \right] \leq (16k^2\kappa^2 + 16k\kappa) \cdot \max_{i \in [k]} \widehat{\mathcal{R}}_n(\mathcal{F}_i).$$

Proof.

$$\begin{aligned} & \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j \left(\left(f(z_j)^\top f(z_j') \right)^2 - 2f(z_j)^\top f(z_j^+) \right) \right) \right] \\ & \leq \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j \left(f(z_j)^\top f(z_j') \right)^2 \right) \right] + 2\mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f(z_j)^\top f(z_j^+) \right) \right] \\ & \leq 2k\kappa \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f(z_j)^\top f(z_j') \right) \right] + 2\mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f(z_j)^\top f(z_j^+) \right) \right] \\ & \leq (2k^2\kappa + 2k) \max_{\substack{z_1, z_2, \dots, z_n \\ z_1', z_2', \dots, z_n'}} \max_{i \in [k]} \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f_i(z_j) f_i(z_j') \right) \right], \end{aligned}$$

here the second inequality is by Talagrand's lemma. Notice that for any z_1, z_2, \dots, z_n and z_1', z_2', \dots, z_n' in \mathcal{X} and any $i \in [k]$ we have

$$\begin{aligned} & \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f_i(z_j) f_i(z_j') \right) \right] \\ & \leq \frac{1}{2} \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j \left(f_i(z_j) + f_i(z_j') \right)^2 \right) \right] + \frac{1}{2} \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j \left(f_i(z_j) - f_i(z_j') \right)^2 \right) \right] \\ & \leq 4\kappa \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f_i(z_j) \right) \right] + 4\kappa \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f_i(z_j') \right) \right], \end{aligned}$$

where the first inequality is by Talagrand's lemma. Combine these two equations and we get:

$$\begin{aligned} & \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j \left(\left(f(z_j)^\top f(z_j') \right)^2 - 2f(z_j)^\top f(z_j^+) \right) \right) \right] \\ & \leq (16k^2\kappa^2 + 16k\kappa) \max_{z_1, z_2, \dots, z_n} \max_{i \in [k]} \mathbb{E}_\sigma \left[\sup_{f_i \in \mathcal{F}_i} \left(\frac{1}{n} \sum_{j=1}^n \sigma_j f_i(z_j) \right) \right]. \end{aligned}$$

□

Proof of Theorem 4.1. By Claim D.2 and Claim D.4, we know that $\mathbb{E}_S[\widehat{\mathcal{L}}_S(f)] = \mathcal{L}(f)$, where S is sampled by first sampling $\widehat{\mathcal{X}}$ then sample S according to Definition D.3. Notice that when $\widehat{\mathcal{X}}$ contains n i.i.d. samples natural data, the set of random tuples S contains n i.i.d tuples. Therefore, we can apply generalization bound with Rademacher complexity to get a uniform convergence bound. In

particular, by Lemma D.5 and notice the fact that $\left(f(z_j)^\top f(z'_j)\right)^2 - 2f(z_j)^\top f(z_j^+)$ always take values in range $[-2k\kappa^2, 2k\kappa^2 + k^2\kappa^4]$, we apply standard generalization analysis based on Rademacher complexity and get: with probability at least $1 - \delta^2/4$ over the randomness of $\hat{\mathcal{X}}$ and S , we have for any $f \in \mathcal{F}$,

$$\mathcal{L}(f) \leq \hat{\mathcal{L}}_S(f) + (32k^2\kappa^2 + 32k\kappa) \max_{i \in [k]} \hat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}_i) + (4k\kappa^2 + k^2\kappa^4) \cdot \sqrt{\frac{4 \log 2/\delta}{n_{\text{pre}}}}. \quad (44)$$

This means with probability at least $1 - \delta/2$ over random $\hat{\mathcal{X}}$, we have: with probability at least $1 - \delta/2$ over random tuples S conditioned on $\hat{\mathcal{X}}$, Equation (44) holds. Since both $\mathcal{L}(f)$ and $\hat{\mathcal{L}}_{n_{\text{pre}}}(f)$ take value in range $[-2k\kappa^2, 2k\kappa^2 + k^2\kappa^4]$, we have: with probability at least $1 - \delta/2$ over random $\hat{\mathcal{X}}$, we have for any $f \in \mathcal{F}$,

$$\mathcal{L}(f) \leq \hat{\mathcal{L}}_{n_{\text{pre}}}(f) + (32k^2\kappa^2 + 32k\kappa) \cdot \max_{i \in [k]} \hat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}_i) + (4k\kappa^2 + k^2\kappa^4) \cdot \left(\sqrt{\frac{4 \log 2/\delta}{n_{\text{pre}}}} + \frac{\delta}{2} \right).$$

Since negating the functions in a function class doesn't change its Rademacher complexity, we also have the other direction: with probability at least $1 - \delta/2$ over random $\hat{\mathcal{X}}$, we have for any $f \in \mathcal{F}$,

$$\mathcal{L}(f) \geq \hat{\mathcal{L}}_{n_{\text{pre}}}(f) - (32k^2\kappa^2 + 32k\kappa) \cdot \max_{i \in [k]} \hat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}_i) + (4k\kappa^2 + k^2\kappa^4) \cdot \left(\sqrt{\frac{4 \log 2/\delta}{n_{\text{pre}}}} + \frac{\delta}{2} \right).$$

Combine them together we get the excess risk bound: with probability at least $1 - \delta$, we have

$$\mathcal{L}(\hat{f}) \leq \mathcal{L}(f_{\mathcal{F}}^*) + (64k^2\kappa^2 + 64k\kappa) \cdot \max_{i \in [k]} \hat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}_i) + (8k\kappa^2 + 2k^2\kappa^4) \cdot \left(\sqrt{\frac{4 \log 2/\delta}{n_{\text{pre}}}} + \frac{\delta}{2} \right),$$

where \hat{f} is minimizer of $\hat{\mathcal{L}}_{n_{\text{pre}}}(f)$ in \mathcal{F} and $f_{\mathcal{F}}^*$ is minimizer of $\mathcal{L}(f)$ in \mathcal{F} . Set $c_1 = 64k^2\kappa^2 + 64k\kappa$ and $c_2 = 16k\kappa^2 + 4k^2\kappa^4$ and notice that $\max_{i \in [k]} \hat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F}_i) = \hat{\mathcal{R}}_{n_{\text{pre}}/2}(\mathcal{F})$ finishes the proof. \square

D.2 Generalization bound for spectral contrastive learning with deep neural networks

In this section, we exemplify Theorem 4.1 with the norm-controlled Rademacher complexity bound introduced in (Golowich et al., 2018), which gives the following theorem.

Theorem D.6. *Assume \mathcal{X} is a subset of Euclidean space \mathbb{R}^d and $\|x\|_2 \leq C_x$ for any $x \in \mathcal{X}$. Let \mathcal{F} be a hypothesis class of norm-controlled l -layer deep neural networks defined as*

$$\{x \rightarrow P_\kappa(W_l \sigma(W_{l-1} \sigma(\dots \sigma(W_1 x)))) : \|W_i\|_F \leq C_{w,i}\}$$

where $\sigma(\cdot)$ is element-wise ReLU activation, $P_\kappa(\cdot)$ is element-wise projection to interval $[-\kappa, \kappa]$ for some $\kappa > 0$, $C_{w,i}$ is the norm bound of the i -th layer, W_l has k rows and W_1 has d columns. Then, with probability at least $1 - \delta$ over randomness of a dataset with size $2n_{\text{pre}}$, we have

$$\mathcal{L}(\hat{f}) \leq \mathcal{L}_{\mathcal{F}}^* + c_1 \cdot \frac{C_x C_w \sqrt{l}}{\sqrt{n_{\text{pre}}}} + c_2 \cdot \left(\sqrt{\frac{\log 1/\delta}{n_{\text{pre}}}} + \delta \right),$$

where \hat{f} is the minimizer of $\hat{\mathcal{L}}_{2n_{\text{pre}}}(f)$ in \mathcal{F} , $\mathcal{L}_{\mathcal{F}}^*$ is the minimal $\mathcal{L}(f)$ achievable by any function $f \in \mathcal{F}$, $C_w := \prod_{i=1}^l C_{w,i}$, constants $c_1 \lesssim k^2\kappa^2 + k\kappa$ and $c_2 \lesssim k\kappa^2 + k^2\kappa^4$.

Proof of Theorem D.6. Consider the following hypothesis class of real-valued neural networks:

$$\mathcal{F}_{\text{real}} \triangleq \left\{ x \rightarrow \widehat{W}_l \sigma(W_{l-1} \sigma(\cdots \sigma(W_1 x))) : \|W_i\|_F \leq C_{w,i} \right\}$$

where $\sigma(\cdot)$ is element-wise ReLU activation and $C_{w,i}$ is the norm bound of the i -th layer defined in the theorem, W_l has k rows and \widehat{W}_l is a vector. By Theorem 1 of (Golowich et al., 2018), we have

$$\widehat{\mathcal{R}}_{n_{\text{pre}}}(\mathcal{F}_{\text{real}}) \leq \frac{C_x(\sqrt{2 \log(2)l} + 1)C_w}{\sqrt{n_{\text{pre}}}}.$$

Let the projection version of this hypothesis class be:

$$\mathcal{F}_{\text{real+proj}} \triangleq \left\{ x \rightarrow P_\kappa(\widehat{W}_l \sigma(W_{l-1} \sigma(\cdots \sigma(W_1 x)))) : \|W_i\|_F \leq C_{w,i} \right\},$$

where $P_\kappa(\cdot)$ projects a real number into interval $[-C_w, C_w]$. Notice that $P_\kappa(\cdot)$ is 1-Lipschitz, by Telegrand's lemma we have

$$\widehat{\mathcal{R}}_{n_{\text{pre}}}(\mathcal{F}_{\text{real+proj}}) \leq \frac{C_x(\sqrt{2 \log(2)l} + 1)C_w}{\sqrt{n_{\text{pre}}}}.$$

For each $i \in [k]$, define function $f_i : \mathcal{X} \rightarrow \mathbb{R}$ such that $f_i(x)$ is the i -th dimension of $f(x)$, define \mathcal{F}_i be the hypothesis class including all f_i for $f \in \mathcal{F}$. Then when \mathcal{F} is the composition of deep neural networks and projection function as defined in the theorem, it is obvious to see that $\mathcal{F}_i = \mathcal{F}_{\text{real+proj}}$ for all $i \in [k]$. Therefore, by Theorem 4.1 we have

$$\mathcal{L}(\hat{f}) \leq \mathcal{L}_{\mathcal{F}}^* + c_1 \cdot \frac{C_x(\sqrt{2 \log(2)l} + 1)C_w}{\sqrt{n_{\text{pre}}}} + c_2 \cdot \left(\sqrt{\frac{\log 2/\delta}{n_{\text{pre}}}} + \delta \right),$$

and absorbing the constants into c_1 finishes the proof. \square

D.3 Proof of Theorem 4.2

In this section we give the proof of Theorem 4.2. We will first prove the following theorem that characterize the error propagation from pre-training to the downstream task.

Theorem D.7 (Error propagation from pre-training to the downstream task). *Assume representation dimension $k \geq 4r + 2$, Assumption 3.6 holds for $\alpha > 0$ and Assumption 3.7 holds. Recall γ_i be the i -th largest eigenvalue of the normalized adjacency matrix. Then, for any $\epsilon > 0$ and $\hat{f}_{\text{emp}} \in \mathcal{F}$ such that $\mathcal{L}(\hat{f}_{\text{emp}}) < \mathcal{L}(f_{\text{pop}}^*) + \epsilon$, we have:*

$$\mathcal{E}(\hat{f}_{\text{emp}}) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log k + \frac{k\epsilon}{\Delta_\gamma^2},$$

where $\Delta_\gamma := \gamma_{\lfloor 3k/4 \rfloor} - \gamma_k$ is the eigenvalue gap between the $\lfloor 3k/4 \rfloor$ -th and the k -th eigenvalue. Furthermore, there exists a linear head $\widehat{B} \in \mathbb{R}^{k \times r}$ that achieves this error and has norm bound

$$\|\widehat{B}\|_F \leq \frac{2(k+1)}{1-\lambda_k}. \quad (45)$$

We first introduce the following definitions of ϵ -optimal minimizers of matrix approximation loss and population spectral contrastive loss:

Definition D.8. We say a function \hat{f}_{mf} is ϵ -optimal minimizer of matrix approximation loss \mathcal{L}_{mf} if

$$\mathcal{L}_{\text{mf}}(\hat{F}_{\text{mf}}) \leq \min_F \mathcal{L}_{\text{mf}}(F) + \epsilon,$$

where \hat{F}_{mf} is \hat{f}_{mf} written in the matrix form. We say a function \hat{f} is ϵ -optimal minimizer of spectral contrastive loss \mathcal{L} if

$$\mathcal{L}(\hat{f}) \leq \min_f \mathcal{L}(f) + \epsilon.$$

We introduce the following generalized version of Theorem B.3, which captures the main effects of error in the representation.

Theorem D.9. [Generalization of Theorem B.3] Assume the set of augmented data \mathcal{X} is finite. Let λ_i be the i -th smallest eigenvalue of the normalized Laplacian matrix. Let $\hat{f} \in \arg \min_{f: \mathcal{X} \rightarrow \mathbb{R}^k}$ be a ϵ -optimal minimizer of the spectral contrastive loss function $\mathcal{L}(f)$ with $k \in \mathbb{Z}^+$. Then, for any labeling function $\hat{y} : \mathcal{X} \rightarrow [r]$ there exists a linear probe $\hat{B} \in \mathbb{R}^{r \times k}$ with norm bound $\|\hat{B}\|_F \leq \frac{2(k+1)}{1-\lambda_k}$ such that

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\left\| \vec{y}(\bar{x}) - \hat{B} \hat{f}_{\text{emp}}(x) \right\|_2^2 \right] \lesssim \min_{1 \leq k' \leq k} \left(\frac{\phi^{\hat{y}}}{\lambda_{k'+1}} + \frac{k' \epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) + \Delta(y, \hat{y}),$$

and

$$\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{\hat{f}, \hat{B}}(x) \neq y(\bar{x}) \right) \lesssim \min_{1 \leq k' \leq k} \left(\frac{\phi^{\hat{y}}}{\lambda_{k'+1}} + \frac{k' \epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) + \Delta(y, \hat{y}),$$

where $\phi^{\hat{y}}$ and $\Delta(y, \hat{y})$ are defined in Equations 18 and 19 respectively.

The proof of Theorem D.9 is deferred to Section D.4.

Now we are ready to prove Theorem 4.2 using Theorem D.9.

Proof of Theorem D.7. In Theorem D.9 we let $k' = \lfloor \frac{3}{4}k \rfloor$ on the RHS of the bound and get: for any $\hat{y} : \mathcal{X} \rightarrow [r]$ there exists $\hat{B} \in \mathbb{R}^{r \times k}$ such that

$$\Pr_{x \sim \mathcal{P}_{\bar{\mathcal{X}}}, \bar{x} \sim \mathcal{A}(\cdot|x)} \left(g_{\hat{f}, \hat{B}}(\bar{x}) \neq y(x) \right) \lesssim \frac{\phi^{\hat{y}}}{\lambda_{\lfloor \frac{3}{4}k \rfloor + 1}} + \frac{k\epsilon}{(\lambda_{k+1} - \lambda_{\lfloor \frac{3}{4}k \rfloor})^2} + \Delta(y, \hat{y}).$$

Let S_1, S_2, \dots, S_r be the partition of \mathcal{X} induced by the classifier g in Assumption 3.6. Define function $\hat{y} : \mathcal{X} \rightarrow [r]$ as follows: for an augmented datapoint $x \in \mathcal{X}$, we use function $\hat{y}(x)$ to represent the index of set that x is in, i.e., $x \in S_{\hat{y}(x)}$. Then by Lemma B.5 we have $\phi^{\hat{y}} \leq 2\alpha$ and $\Delta(y, \hat{y}) \leq \alpha$. In Lemma B.4 let $(1 + \zeta)t = \lfloor \frac{3}{4}k \rfloor + 1$ and $t = \lfloor \frac{k}{2} \rfloor$, then there is $\zeta \geq 0.5$, so we have: there exists a partition $S_1, \dots, S_{\lfloor \frac{k}{2} \rfloor} \subset \mathcal{X}$ such that $\phi_G(S_i) \lesssim \sqrt{\lambda_{\lfloor \frac{3}{4}k \rfloor + 1} \log(k)}$ for $\forall i \in [\lfloor \frac{k}{2} \rfloor]$. By Definition 3.4, we have $\rho_{\lfloor \frac{k}{2} \rfloor} \lesssim \sqrt{\lambda_{\lfloor \frac{3}{4}k \rfloor + 1} \log(k)}$, which leads to $\frac{1}{\lambda_{\lfloor \frac{3}{4}k \rfloor + 1}} \lesssim \frac{\log(k)}{\rho_{\lfloor \frac{k}{2} \rfloor}^2}$. So we have

$$\begin{aligned} \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{\hat{f}, \hat{B}}(x) \neq y(\bar{x}) \right) &\lesssim \frac{\alpha}{\rho_{\lfloor \frac{k}{2} \rfloor}^2} \cdot \log(k) + \frac{k\epsilon}{(\lambda_{k+1} - \lambda_{\lfloor \frac{3}{4}k \rfloor})^2} \\ &\lesssim \frac{\alpha}{\rho_{\lfloor \frac{k}{2} \rfloor}^2} \cdot \log(k) + \frac{k\epsilon}{(\lambda_k - \lambda_{\lfloor \frac{3}{4}k \rfloor})^2}. \end{aligned}$$

Notice that by the definition of ensemble linear probe predictor, $\bar{g}_{\hat{f}, \hat{B}}(\bar{x}) \neq y(\bar{x})$ happens only if more than half of the augmentations of \bar{x} predicts differently from $y(\bar{x})$, so we have $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}} \left(\bar{g}_{\hat{f}, \hat{B}} \neq y(\bar{x}) \right) \leq 2 \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{\hat{f}, \hat{B}}(x) \neq y(\bar{x}) \right)$ which finishes the proof. \square

Proof of Theorem 4.2. Theorem 4.2 is a direct corollary of Theorem 4.1 and Theorem D.7. \square

D.4 Proof of Theorem D.9

In this section, we give the proof for Theorem D.9.

Lemma D.10 (Generalization of Lemma B.8). *Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be a feature extractor, matrix $F \in \mathbb{R}^{N \times k}$ be such that its x -th row is $\sqrt{w_x} \cdot f(x)$. Then, F is an ϵ -optimal minimizer of $\mathcal{L}_{\text{mf}}(F)$ if and only if f is an ϵ -optimal minimizer of the population spectral contrastive loss $\mathcal{L}(f)$.*

Proof of Lemma D.10. The proof follows the proof of Lemma B.8. \square

We will use the following two lemmas about ϵ -optimal minimizer of \mathcal{L}_{mf} :

Lemma D.11. *Let λ_i be the i -th minimal eigenvalue of the normalized Laplacian matrix \mathcal{L} with corresponding unit-norm eigenvector v_i . Let $F \in \mathbb{R}^{N \times k}$ be an ϵ -optimal minimizer of \mathcal{L}_{mf} . Let $\Pi_f v_i$ be the projection of v_i onto the column span of F . Then, there exists vector $b \in \mathbb{R}^k$ with norm bound $\|b\| \leq \|F\|_F / (1 - \lambda_i)$ such that*

$$\|\Pi_f v_i - Fb\|_2^2 \leq \frac{\epsilon}{(1 - \lambda_i)^2}. \quad (46)$$

Furthermore, the norm of F is bounded by

$$\|F\|_F^2 \leq 2(k + \epsilon). \quad (47)$$

Proof of Lemma D.11. Since columns of $\bar{A} - \Pi_f \bar{A}$ and columns of $\Pi_f \bar{A} - FF^\top$ are in orthogonal subspaces, we have

$$\|\bar{A} - FF^\top\|_F^2 = \|\bar{A} - \Pi_f \bar{A}\|_F^2 + \|\Pi_f \bar{A} - FF^\top\|_F^2. \quad (48)$$

On one hand, since $\Pi_f \bar{A}$ is a rank- k matrix, we know that $\|\bar{A} - \Pi_f \bar{A}\|_F^2 \geq \min_F \mathcal{L}_{\text{mf}}(F)$. On the other hand, by the definition of ϵ -optimal minimizer, we have $\|\bar{A} - FF^\top\|_F^2 \leq \min_F \mathcal{L}_{\text{mf}}(F) + \epsilon$. Thus, we have

$$\|\Pi_f \bar{A} - FF^\top\|_F^2 \leq \epsilon. \quad (49)$$

Since $\bar{A} = \sum_{i=1}^N (1 - \lambda_i) v_i v_i^\top$, we have $v_i = \frac{1}{1 - \lambda_i} \bar{A} v_i$. Thus,

$$\Pi_f v_i = \frac{1}{1 - \lambda_i} \Pi_f (\bar{A} v_i) = \frac{1}{1 - \lambda_i} FF^\top v_i + \frac{1}{1 - \lambda_i} (\Pi_f \bar{A} - FF^\top) v_i. \quad (50)$$

Let $b = \frac{1}{1 - \lambda_i} FF^\top v_i$, we have

$$\|\Pi_f v_i - Fb\|_2^2 = \frac{1}{(1 - \lambda_i)^2} \left\| (\Pi_f \bar{A} - FF^\top) v_i \right\|_2^2 \quad (51)$$

$$\leq \frac{1}{(1 - \lambda_i)^2} \left\| \Pi_f \bar{A} - FF^\top \right\|_F^2 \quad (52)$$

$$\leq \frac{\epsilon}{(1 - \lambda_i)^2}. \quad (53)$$

To bound the norm of F , we first notice that

$$\|\Pi_f \bar{A}\|_F^2 = \text{Tr}(\bar{A}^2 \Pi_f) \leq \text{Tr}(\Pi_f) = k, \quad (54)$$

where the inequality uses that fact that \bar{A}^2 has operator norm at most 1. Combine this result with $\|\Pi_f \bar{A} - FF^\top\|_F^2 \leq \epsilon$ we have

$$\|FF^\top\|_F \leq \sqrt{k} + \sqrt{\epsilon}. \quad (55)$$

Since FF^\top has rank at most k , we can write its SVD decomposition as $FF^\top = U\Sigma U^\top$ where $U \in \mathbb{R}^{N \times k}$ and $\Sigma \in \mathbb{R}^{k \times k}$. As a result, we have

$$\|F\|_F^2 = \text{Tr}(FF^\top) = \text{Tr}(\Sigma) \leq \sqrt{k} \sqrt{\text{Tr}(\Sigma^2)} = \sqrt{k} \|FF^\top\|_F \leq k + \sqrt{k\epsilon} \leq 2(k + \epsilon). \quad (56)$$

□

Lemma D.12. *Let λ_i be the i -th minimal eigenvalue of the normalized Laplacian matrix \mathcal{L} with corresponding unit-norm eigenvector v_i . Let $F \in \mathbb{R}^{N \times k}$ be an ϵ -optimal minimizer of \mathcal{L}_{mf} . Let $\Pi_f^\perp v_i$ be the projection of v_i onto the subspace orthogonal to the column span of F . Then, for $i \leq k$ we have*

$$\|\Pi_f^\perp v_i\|_2^2 \leq \frac{\epsilon}{(\lambda_{k+1} - \lambda_i)^2}.$$

Proof. Recall normalized adjacency matrix $\bar{A} = I - L$. We use \bar{A}_i to denote the i -th column of \bar{A} . We use \hat{A} to denote matrix FF^\top and \hat{A}_i to denote the i -th column of \hat{A} . Let z_1, \dots, z_k be unit-norm orthogonal vectors in the column span of F . Since the column span of \hat{A} is the same as the column span of F , we know columns of \hat{A} are in $\text{span}\{z_1, \dots, z_k\}$. Let z_{k+1}, \dots, z_N be unit-norm orthogonal vectors such that together with z_1, \dots, z_k they form an orthonormal basis of \mathbb{R}^N . We use Π_f and Π_f^\perp to denote matrices $\sum_{j=1}^k z_j z_j^\top$ and $\sum_{j=k+1}^N z_j z_j^\top$ respectively, then for any vector $v \in \mathbb{R}^N$, vectors $\Pi_f v$ and $\Pi_f^\perp v$ are the projections of v onto the column span of F and its orthogonal space respectively.

We first give a lower bound of $\mathcal{L}_{\text{mf}}(F)$ as follows:

$$\begin{aligned} \mathcal{L}_{\text{mf}}(F) &= \|\bar{A} - \hat{A}\|_F^2 = \sum_{j=1}^N \|\bar{A}_j - \hat{A}_j\|_2^2 \geq \sum_{j=1}^N \|\bar{A}_j - \Pi_f \bar{A}_j\|_2^2 \\ &= \sum_{j=1}^N \left\| \bar{A}_j - \left(\sum_{t=1}^k z_t z_t^\top \right) \bar{A}_j \right\|_2^2 = \sum_{j=1}^N \left\| \left(\sum_{t=k+1}^N z_t z_t^\top \right) \bar{A}_j \right\|_2^2 \\ &= \left\| \left(\sum_{t=k+1}^N z_t z_t^\top \right) \bar{A} \right\|_F^2 = \|\Pi_f^\perp \bar{A}\|_F^2. \end{aligned}$$

where the first equality is by definition of $\mathcal{L}_{\text{mf}}(F)$, the second equality is by writing the Frobenius norm square as the sum of column norm square, the inequality is because \hat{A}_j must be in the span of z_1, \dots, z_k while $\Pi_f \bar{A}_j$ is the vector in this span that is closest to \bar{A}_j , the third equality is writing the projection function in the matrix form, the fourth equality is because z_1, \dots, z_d are an orthonormal basis, the fifth equality is rewriting to Frobenius norm, and the last equality is by definition of Π_f^\perp .

Notice that

$$\|\Pi_f^\perp \bar{A}\|_F^2 = \text{Tr}(\bar{A}^\top \Pi_f^\perp \Pi_f^\perp \bar{A}) = \text{Tr}(\bar{A}^\top \Pi_f^\perp \bar{A}) = \text{Tr}(\bar{A} \bar{A}^\top \Pi_f^\perp).$$

We can rewrite the above lower bound as

$$\mathcal{L}_{\text{mf}}(F) \geq \text{Tr}(\bar{A} \bar{A}^\top \Pi_f^\perp) = \text{Tr} \left(\sum_{j=1}^N (1 - \lambda_j)^2 v_j v_j^\top \sum_{t=k+1}^N z_t z_t^\top \right) = \sum_{j=1}^N \sum_{t=k+1}^N (1 - \lambda_j)^2 \langle v_j, z_t \rangle^2.$$

We define variable $S_j \triangleq \sum_{t=1}^j \sum_{l=k+1}^d \langle v_t, z_l \rangle^2$ for any $j \in [N]$. Also denote $\lambda_{d+1} = 1$. We have the following equality:

$$\sum_{j=1}^N \sum_{t=k+1}^N (1 - \lambda_j)^2 \langle v_j, z_t \rangle^2 = \sum_{j=1}^N ((1 - \lambda_j)^2 - (1 - \lambda_{j+1})^2) S_j.$$

Notice that $S_j \geq 0$ and also when $i \leq j \leq k$, we have $S_j \geq \left\| \Pi_f^\perp v_i \right\|_2^2$, we have

$$\sum_{j=1}^N \sum_{t=k+1}^N (1 - \lambda_j)^2 \langle v_j, z_t \rangle^2 \geq ((1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2) \left\| \Pi_f^\perp v_i \right\|_2^2 + \sum_{j=k+1}^N ((1 - \lambda_j)^2 - (1 - \lambda_{j+1})^2) S_j,$$

where we replace every S_j with 0 when $j < k$, replace S_j with $\left\| \Pi_f^\perp v_i \right\|_2^2$ when $i \leq j \leq k$, and keep S_j when $j \geq k + 1$. Now notice that

$$S_N = \sum_{t=1}^N \sum_{l=k+1}^N \langle v_t, z_l \rangle^2 = \sum_{l=k+1}^N \sum_{t=1}^N \langle v_t, z_l \rangle^2 = \sum_{l=k+1}^N \|z_l\|_2^2 = N - k,$$

and also

$$S_{j+1} - S_j = \sum_{l=k+1}^N \langle v_{j+1}, z_l \rangle^2 \leq \sum_{l=1}^N \langle v_{j+1}, z_l \rangle^2 = 1,$$

there must be $S_j \geq j - k$ when $j \geq k + 1$. So we have

$$\begin{aligned} & \sum_{j=1}^N \sum_{t=k+1}^N (1 - \lambda_j)^2 \langle v_j, z_t \rangle^2 \\ & \geq ((1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2) \left\| \Pi_f^\perp v_i \right\|_2^2 + \sum_{j=k+1}^N ((1 - \lambda_j)^2 - (1 - \lambda_{j+1})^2) (j - k) \\ & = ((1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2) \left\| \Pi_f^\perp v_i \right\|_2^2 + \sum_{j=k+1}^N (1 - \lambda_j)^2 \\ & = ((1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2) \left\| \Pi_f^\perp v_i \right\|_2^2 + \min_{F \in \mathbb{R}^{N \times k}} \mathcal{L}_{\text{mf}}(F), \end{aligned}$$

where the last equality is by Eckart–Young–Mirsky Theorem. So we know

$$\mathcal{L}_{\text{mf}}(F) \geq ((1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2) \left\| \Pi_f^\perp v_i \right\|_2^2 + \min_{F \in \mathbb{R}^{N \times k}} \mathcal{L}_{\text{mf}}(F), \quad (57)$$

which implies that

$$\left\| \Pi_f^\perp v_i \right\|_2^2 \leq \frac{\epsilon}{(1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2} \leq \frac{\epsilon}{(\lambda_{k+1} - \lambda_i)^2}. \quad (58)$$

□

The following lemma generalizes Lemma B.6.

Lemma D.13 (Generalization of Lemma B.6). *Let \mathcal{L} be the normalized Laplacian matrix of graph $G = (\mathcal{X}, w)$, where $|\mathcal{X}| = N$. Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be an ϵ -optimal minimizer of $\mathcal{L}_{\text{mf}}(f)$. Let F be the matrix form of*

f and F_i is the i -th column of F . Let $R(u) := \frac{u^\top \mathcal{L}u}{u^\top u}$ be the Rayleigh quotient of a vector $u \in \mathbb{R}^N$. Then, for any $k \in \mathcal{Z}^+$ such that $k < N$, there exists a vector $b \in \mathbb{R}^k$ such that

$$\|u - Fb\|_2^2 \leq \min_{1 \leq k' \leq k} \left(\frac{3R(u)}{\lambda_{k'+1}} + \frac{6k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) \|u\|_2^2.$$

Furthermore, the norm of b is upper bounded by

$$\|b\|_2 \leq \frac{2(k+1)}{1 - \lambda_k} \|u\|_2. \quad (59)$$

Proof. Let k' be the choice that minimizes the right hand side. We use $p_v(u)$ to denote the projection of u onto the span of $v_1, \dots, v_{k'}$. We denote the coefficients as $p_v(u) = \sum_{i=1}^{k'} \rho_i v_i$. For every $i \in [k']$, let b_i be the vector in Lemma D.11. Define vector $b = \sum_{i=1}^{k'} \rho_i b_i$.

We use $p_{v,f}(u)$ to denote the projection of $p_v(u)$ onto the span of f_1, \dots, f_k . Then we know that

$$\|u - Fb\|_2^2 \leq 3 \|u - p_v(u)\|_2^2 + 3 \|p_v(u) - p_{v,f}(u)\|_2^2 + 3 \|p_{v,f}(u) - Fb\|_2^2. \quad (60)$$

By the proof of Lemma B.6, we know that

$$\|u - p_v(u)\|_2^2 \leq \frac{R(u)}{\lambda_{k'+1}} \|u\|_2^2. \quad (61)$$

For the second term, we have

$$\begin{aligned} \|p_v(u) - p_{v,f}(u)\|_2^2 &= \left\| \Pi_f^\perp p_v(u) \right\|_2^2 \\ &= \left\| \sum_{i=1}^{k'} \Pi_f^\perp v_i v_i^\top u \right\|_2^2 \\ &\leq \left(\sum_{i=1}^{k'} \left\| \Pi_f^\perp v_i \right\|_2^2 \right) \cdot \left(\sum_{i=1}^{k'} (v_i^\top u)^2 \right) \\ &\leq \frac{k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \|u\|_2^2, \end{aligned} \quad (62)$$

where the first inequality is by Cauchy-Schwarz inequality and the second inequality is by Lemma D.12.

For the third term, we have

$$\begin{aligned} \|p_{v,f}(u) - Fb\|_2^2 &= \left\| \sum_{i=1}^{k'} \rho_i (\Pi_f v_i - Fb_i) \right\|_2^2 \\ &\leq k' \sum_{i=1}^{k'} \rho_i^2 \|\Pi_f v_i - Fb_i\|_2^2 \\ &\leq \frac{k'\epsilon}{(1 - \lambda_{k'})^2} \|u\|_2^2, \end{aligned} \quad (63)$$

where the first inequality is by Cauchy-Schwarz inequality, and the second inequality is by Lemma D.11. Plugging Equation (61), Equation (62), and Equation (63) into Equation (60) finishes the proof.

To bound the norm of b , we use Lemma D.11 and have

$$\|b\|_2^2 = \left\| \sum_{i=1}^{k'} \rho_i b_i \right\|_2^2 \leq k' \sum_{i=1}^{k'} \rho_i^2 \|b_i\|_2^2 \leq \frac{k' \|u\|_2^2}{(1 - \lambda_{k'})^2} \|F\|_F^2 \leq \frac{2k'(k+1)}{(1 - \lambda_{k'})^2} \|u\|_2^2. \quad (64)$$

□

Now we prove Theorem D.9 using the above lemmas.

Proof of Theorem D.9. Let $\widehat{F} \in \mathbb{R}^{N \times k}$ be such that its x -th row is $\sqrt{w_x} \cdot \widehat{f}(x)$. By Lemma D.10, \widehat{F} is an ϵ -optimal minimizer of $\mathcal{L}_{\text{mf}}(F)$.

For each $i \in [r]$, we define the function $u_i(x) = \mathbb{1}[\widehat{y}(x) = i] \cdot \sqrt{w_x}$. Let $u : \mathcal{X} \rightarrow \mathbb{R}^k$ be the function such that $u(x)$ has u_i at the i -th dimension. By Lemma D.13, there exists a vector $b_i \in \mathbb{R}^k$ such that

$$\|u_i - \widehat{F} b_i\|_2^2 \leq \min_{1 \leq k' \leq k} \left(\frac{3R(u_i)}{\lambda_{k'+1}} + \frac{6k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) \|u_i\|_2^2$$

Let matrices $U = [u_1, \dots, u_r]$ and $\widehat{B}^\top = [b_1, \dots, b_r]$. We sum the above equation over all $i \in [r]$ and get

$$\begin{aligned} \|U - \widehat{F} \widehat{B}^\top\|_F^2 &\leq \sum_{i=1}^r \min_{1 \leq k' \leq k} \left(\frac{3R(u_i)}{\lambda_{k'+1}} + \frac{6k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) \|u_i\|_2^2 \\ &\leq \min_{1 \leq k' \leq k} \sum_{i=1}^r \left(\frac{3R(u_i)}{\lambda_{k'+1}} \|u_i\|_2^2 + \frac{6k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \|u_i\|_2^2 \right). \end{aligned} \quad (65)$$

Notice that

$$\begin{aligned} \sum_{i=1}^r R(u_i) \|u_i\|_2^2 &= \sum_{i=1}^r \frac{1}{2} \phi_i^{\widehat{y}} \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\widehat{y}(x) = i] \\ &= \frac{1}{2} \sum_{i=1}^r \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[(\widehat{y}(x) = i \wedge \widehat{y}(x') \neq i) \text{ or } (\widehat{y}(x) \neq i \wedge \widehat{y}(x') = i)] \\ &= \frac{1}{2} \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[\widehat{y}(x) \neq \widehat{y}(x')] = \frac{1}{2} \phi^{\widehat{y}}, \end{aligned} \quad (66)$$

where the first equality is by Claim B.7. On the other hand, we have

$$\sum_{i=1}^r \|u_i\|_2^2 = \sum_{i=1}^r \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\widehat{y}(x) = i] = \sum_{x \in \mathcal{X}} w_x = 1. \quad (67)$$

Plugging Equation (66) and Equation (67) into Equation (65) gives us

$$\|U - \widehat{F} \widehat{B}^\top\|_F^2 \leq \min_{1 \leq k' \leq k} \left(\frac{3\phi^{\widehat{y}}}{2\lambda_{k'+1}} + \frac{3k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right).$$

Notice that by definition of $u(x)$, we know that prediction $g_{\widehat{f}, \widehat{B}}(x) \neq \widehat{y}(x)$ only happens if $\|u(x) - \widehat{B} \widehat{f}(x)\|_2^2 \geq \frac{w_x}{2}$. Hence we have

$$\sum_{x \in \mathcal{X}} \frac{1}{2} w_x \cdot \mathbb{1}[g_{\widehat{f}, \widehat{B}}(x) \neq \widehat{y}(x)] \leq \sum_{x \in \mathcal{X}} \|u(x) - \widehat{B} \widehat{f}(x)\|_2^2 = \|U - \widehat{F} \widehat{B}^\top\|_F^2.$$

Now we are ready to bound the error rate on \mathcal{X} :

$$\begin{aligned} \Pr_{x \sim \mathcal{X}}(g_{\hat{f}, \hat{B}}(x) \neq \hat{y}(x)) &= \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1} \left[g_{\hat{f}, \hat{B}}(x) \neq \hat{y}(x) \right] \\ &\leq 2 \cdot \left\| U - \hat{f} \hat{B}^\top \right\|_F^2 \leq \min_{1 \leq k' \leq k} \left(\frac{3\phi^{\hat{y}}}{\lambda_{k'+1}} + \frac{6k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right). \end{aligned}$$

Here for the equality we are using the fact that $\Pr(x) = w_x$. We finish the proof by noticing that by the definition of $\Delta(y, \hat{y})$:

$$\begin{aligned} \Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{\hat{f}, \hat{B}}(x) \neq y(\bar{x}) \right) &\leq \Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{\hat{f}, \hat{B}}(x) \neq \hat{y}(x) \right) + \Pr_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(y(\bar{x}) \neq \hat{y}(x) \right) \\ &\leq \min_{1 \leq k' \leq k} \left(\frac{3\phi^{\hat{y}}}{\lambda_{k'+1}} + \frac{6k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) + \Delta(y, \hat{y}). \end{aligned}$$

The norm of \hat{B} can be bounded using Lemma D.13 as:

$$\left\| \hat{B} \right\|_F \leq \frac{2(k+1)}{1 - \lambda_k} \sqrt{\sum_{i=1}^r \|u_i\|_2^2} = \frac{2(k+1)}{1 - \lambda_k}. \quad (68)$$

□

E Proofs for Section 4.2

In this section we give the proof of Theorem 4.3.

Proof of Theorem 4.3. Let \hat{f}_{emp} be the minimizer of the empirical spectral contrastive loss. Let $\epsilon = \mathcal{L}(\hat{f}_{\text{emp}}) - \mathcal{L}(f_{\text{pop}}^*)$. We abuse notation and use y_i to denote $y(\bar{x}_i)$, and let $z_i = \hat{f}_{\text{emp}}(x_i)$. We first study the average empirical Rademacher complexity of the capped quadratic loss on a dataset $\{(z_i, y_i)\}_{i=1}^{n_{\text{down}}}$, where (z_i, y_i) is sampled as in Section 4.2:

$$\begin{aligned} \hat{\mathcal{R}}_{n_{\text{down}}}(\ell) &:= \mathbb{E}_{\{(z_i, y_i)\}_{i=1}^{n_{\text{down}}}} \mathbb{E}_{\sigma} \left[\sup_{\|B\|_F \leq C_k} \frac{1}{n_{\text{down}}} \left[\sum_{i=1}^{n_{\text{down}}} \sigma_i \ell((z_i, y_i), B) \right] \right] \\ &\leq 2r \mathbb{E}_{\{(z_i, y_i)\}_{i=1}^{n_{\text{down}}}} \mathbb{E}_{\sigma} \left[\sup_{\|b\|_2 \leq C_k} \frac{1}{n_{\text{down}}} \left[\sum_{i=1}^{n_{\text{down}}} \sigma_i w^\top z_i \right] \right] \\ &\leq 2r C_k \sqrt{\frac{\mathbb{E}[\|z_i\|_2^2]}{n_{\text{down}}}} \leq 2r C_k \sqrt{\frac{2(k+\epsilon)}{n_{\text{down}}}}, \end{aligned}$$

where the first inequality uses Talagrand's lemma and the fact that ℓ_σ is 2-Lipschitz, the second inequality is by standard Rademacher complexity of linear models, the third inequality is by the feature norm bound in Lemma D.11.

By Theorem D.9 and follow the proof of Theorem D.7, we know that there exists a linear probe \hat{B}^* with norm bound $\left\| \hat{B}^* \right\|_F \leq C_k$ such that

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{x}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\ell \left((\hat{f}_{\text{emp}}(x), y(\bar{x})), \hat{B}^* \right) \right] \lesssim \frac{\alpha}{\rho_{\lfloor \frac{k}{2} \rfloor}^2} \cdot \log(k) + \frac{k\epsilon}{(\lambda_k - \lambda_{\lfloor \frac{3}{4}k \rfloor})^2}.$$

Let \hat{B} be the minimizer of $\sum_{i=1}^{n_{\text{down}}} \ell((z_i, y_i), B)$ subject to $\|B\|_F \leq C_k$, then by standard generalization bound, we have: with probability at least $1 - \delta$, we have

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{X}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\ell \left((\hat{f}_{\text{emp}}(x), y(\bar{x})), \hat{B} \right) \right] \lesssim \frac{\alpha}{\rho_{\lfloor \frac{k}{2} \rfloor}^2} \cdot \log(k) + \frac{k\epsilon}{(\lambda_k - \lambda_{\lfloor \frac{3}{4}k \rfloor})^2} + \frac{rC_k\sqrt{k} + \epsilon}{\sqrt{n_{\text{down}}}} + \sqrt{\frac{\log 1/\delta}{n_{\text{down}}}}.$$

Notice that $y(\bar{x}) \neq g_{\hat{f}_{\text{emp}}, \hat{B}}(x)$ only if $\ell \left((\hat{f}_{\text{emp}}(x), y(\bar{x})), \hat{B} \right) \geq \frac{1}{2}$, we have that when $\epsilon < 1$ the error bound

$$\Pr_{\bar{x} \sim \mathcal{P}_{\bar{X}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{\hat{f}_{\text{emp}}, \hat{B}}(x) \neq y(\bar{x}) \right) \lesssim \frac{\alpha}{\rho_{\lfloor \frac{k}{2} \rfloor}^2} \cdot \log(k) + \frac{k\epsilon}{(\lambda_k - \lambda_{\lfloor \frac{3}{4}k \rfloor})^2} + \frac{rC_k\sqrt{k}}{\sqrt{n_{\text{down}}}} + \sqrt{\frac{\log 1/\delta}{n_{\text{down}}}}.$$

The result on $\bar{g}_{\hat{f}_{\text{emp}}, \hat{B}}$ naturally follows by the definition of \bar{g} . When $\epsilon > 1$ clearly the bound is also true since LHS is always smaller than 1, so we know that the above bound is true for any ϵ . Plug in the bound for ϵ from Theorem 4.1 finishes the proof. \square

F Formal statements for population with infinite supports

In the main body of the paper, we make the simplifying assumption that the set of augmented data \mathcal{X} is finite (but could be exponential in dimension). Although this is a reasonable assumption given that modern computers store data with finite bits so the possible number of all data has to be finite, one might wonder whether our theory can be generalized to the case where \mathcal{X} is infinite (e.g., the entire Euclidean space \mathbb{R}^d for some integer $d > 0$). In this section, we show that our theory can be straightforwardly extended to the case when \mathcal{X} has infinite supports with some additional regularity conditions. In fact, almost all proofs remain the same as long as we replace sum by integral, finite graph by an infinite graph, adjacency matrix by adjacency operator, and eigenvectors by the eigenfunctions.

For simplicity, we consider the case when $\mathcal{X} = \mathbb{R}^d$ is the set of all augmented data.⁷ The weight matrix $w_{xx'}$ now becomes a weight function $w : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. As usual, let $w(x, x')$ be the marginal probability of generating the pair x and x' from a random natural datapoint $\bar{x} \sim \mathcal{P}_{\bar{X}}$. Or in other words, w is the p.d.f. of the joint distribution of a random positive pair. For any $u \in \mathcal{X}$, define the marginal weight function $w(u) \triangleq \int w(u, z) dz$. A sufficient (but not necessary) condition for our theory to hold is as follows:

Assumption F.1 (Regularity conditions). *The distribution w satisfies the following conditions:*

- (i) For any $u \in \mathcal{X}$, the marginal distribution is well-defined and bounded $w(u) = \int w(u, z) dz < \infty$.
- (ii) There exists $B > 0$ such that for every $u, v \in \mathcal{X}$, the conditional probability with respect to one variable is upper bounded by the marginal probability of the other variable $\frac{w(u, v)}{w(u)} \leq B \cdot w(v)$.

We note that our bound does not depend on value of B —we only the existence of B for a qualitative purpose. When the regularity conditions above hold, we will show that there exists an eigenfunction of the infinite adjacency graph is an analog to the eigenvectors of Laplacian that we introduced in Section B.

Let $L_2(\mathbb{R}^d)$ be the set of all L_2 integratable functions $L_2(\mathbb{R}^d) \triangleq \{f : \mathbb{R}^d \rightarrow \mathbb{R} \mid \int f(z)^2 dz < \infty\}$. For functions $f, g \in L_2(\mathbb{R}^d)$, define their inner product as $\langle f, g \rangle \triangleq \int f(z)g(z) dz$. Note that $\ell_2(\mathbb{R}^d)$ is a Hilbert space.

⁷When \mathcal{X} is a subset of \mathbb{R}^d equipped with a base measure μ , then we will need to replace every dx by $d\mu$ in the formulation below.

To generalize the Laplacian matrix and eigenvectors to the infinite-size \mathcal{X} setting, we consider the notions of Laplacian operators and eigenfunctions. Let $H : L_2(\mathbb{R}^d) \rightarrow L_2(\mathbb{R}^d)$ be a linear operator, a function $f \in L_2(\mathbb{R}^d)$ is an eigenfunction of H if $H(f)(u) = \lambda f(u)$ for any $u \in \mathcal{X}$, where $\lambda \in \mathbb{R}$ is the corresponding eigenvalue. We define the Laplacian operator as $L : L_2(\mathbb{R}^d) \rightarrow L_2(\mathbb{R}^d)$ such that for every $u \in \mathcal{X}$ and function $f \in L_2(\mathbb{R}^d)$, we have

$$L(f)(u) = f(u) - \int \frac{w(u,v)}{\sqrt{w(u)w(v)}} f(v) dv. \quad (69)$$

The following theorem shows the existence of eigenfunctions of the Laplacian operator.

Theorem F.2 (Existence of Eigenfunctions). *When Assumption F.1 is satisfied, there exists an orthonormal basis $\{f_i\}_{i=1}^\infty$ of $L_2(\mathbb{R}^d)$ such that $L(f_i) = \lambda_i f_i$. Furthermore, the eigenvalues satisfy $\lambda_i \in [0, 1]$ and $\lambda_i \leq \lambda_{i+1}$ for any $i \geq 0$.*

Proof of Theorem F.2. Define kernel function $k(u, v) \triangleq \frac{w(u,v)}{\sqrt{w(u)w(v)}}$, we have

$$\int k(u, v)^2 dudv = \int \frac{w(u, v)^2}{w(u)w(v)} dudv \leq B \int w(u, v) dudv = B < \infty. \quad (70)$$

Let I be the identity operator, then $L - I$ is a Hilbert–Schmidt integral operator ([Wikipedia contributors, 2020](#)), so the spectral theorem ([Bump, 1998](#)) applies to $L - I$ hence also applies to L . By the spectral theorem, there exists an orthonormal basis $\{f_i\}_{i=1}^\infty$ of $L_2(\mathbb{R}^d)$ such that $L(f_i) = \lambda_i f_i$.

Notice that

$$\lambda_i = \langle f_i, L(f_i) \rangle = \langle f_i, f_i \rangle - \int \frac{w(u, v)}{\sqrt{w(u)w(v)}} f_i(u) f_i(v) dudv. \quad (71)$$

On the one hand, since $w(u, v) \geq 0$ and $\langle f_i, f_i \rangle = 1$, we have $\lambda_i \leq 1$. On the other hand, notice that by Cauchy-Schwartz inequality,

$$\int \frac{w(u, v)}{\sqrt{w(u)w(v)}} f_i(u) f_i(v) dudv \leq \sqrt{\int f_i(u)^2 \frac{w(u, v)}{w(u)} dudv \cdot \int f_i(v)^2 \frac{w(u, v)}{w(v)} dudv} = \langle f_i, f_i \rangle, \quad (72)$$

so $\lambda_i \geq 0$, which finishes the proof. \square

Given the existence of eigenfunctions guaranteed by Theorem F.2, our results Theorem 3.8, Theorem 4.2 and Theorem 4.3 can all be easily generalized to the infinite-size \mathcal{X} case following exactly the same proof. For example, in the context of Lemma 3.2, u_x will be replaced by $u(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ which belongs to $L_2(\mathbb{R}^d)$, and $f(x) = w(x)^{1/2}u(x)$ as a result belongs to $L_2(w)$. Let $\mathcal{L}_{\text{mf}}(f) = \langle u, Lu \rangle_{L_2(\mathbb{R}^d)} = \langle f, Lf \rangle_{L_2(w)}$. The rest of the derivations follows by replacing the sum in equation (7) by integral (w.r.t to Lebesgue measure). More details on the normalized Laplacian operator and spectral clustering can be found in ([Schiebinger et al., 2015](#)).

We omit the proof for simplicity.