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

Anonymous Author(s)

Affiliation

Address

email

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 data, 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. In all, this work provides the first provable analysis for contrastive learning where the guarantees 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 [4, 6, 10, 14, 23, 24, 35, 38, 41, 42, 50–52]. 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 [3, 12, 13, 15, 22] encourage closer representations for augmentations (views) of the same natural data 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 by Arora et al. [3], Lee et al. [28], Tosh et al. [44, 45] provide guarantees under the assumption that two views are somewhat independent conditioned on the label. However, the pair of augmented examples used in practical algorithms usually exhibit a strong correlation, even conditioned on the label. For instance, two augmentations of the same dog image share much more similarity than augmentations of two different random dog images. Thus the existing theory does not 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

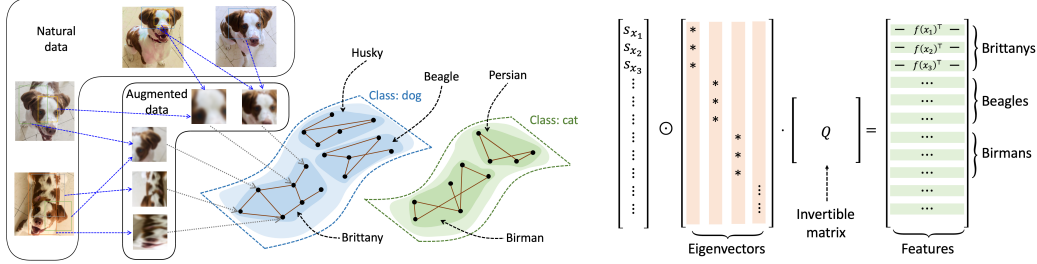


Figure 1: **Left: demonstration of the population augmentation graph.** Two augmented data are connected if they are views of the same natural data. Augmentations of data from different classes in the downstream tasks are assumed to be nearly disconnected, whereas there are more connections within the 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 data 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.

37 realistic data assumption, linear classification using representations learned on a polynomial number
 38 of unlabeled data samples can recover the ground-truth labels of the data with high accuracy.

39 The fundamental data property that we leverage is a notion of continuity of the population data within
 40 the same class. Though a random pair of examples from the same class can be far apart, the pair is
 41 often connected by (many) sequences of examples, where consecutive examples in the sequences are
 42 close neighbors within the same class. This property is more salient when the neighborhood of an
 43 example includes many different types of augmentations. Prior work [49] empirically demonstrates
 44 this type of connectivity property and uses it in the analysis of pseudolabeling algorithms.

45 More formally, we define the *population augmentation graph*, whose vertices are all the augmented
 46 data in the *population* distribution, which can be an exponentially large or infinite set. Two vertices are
 47 connected with an edge if they are augmentations of the same natural example. Our main assumption
 48 is that for some proper $m \in \mathbb{Z}^+$, the sparsest m -partition (Definition 3.4) is large. This intuitively
 49 states that we can’t split the augmentation graph into too many disconnected sub-graphs by only
 50 removing a sparse set of edges. This assumption can be seen as a graph-theoretic version of the
 51 continuity assumption on population data. We also assume that there are very few edges across
 52 different ground-truth classes (Assumption 3.5). Figure 1 (left) illustrates a realistic scenario where
 53 dog and cat are the ground-truth categories, between which edges are very rare. Each breed forms a
 54 sub-graph that has sufficient inner connectivity and thus cannot be further partitioned.

55 Our assumption fundamentally does not require conditional independence and can allow disconnected
 56 sub-graphs within a class. The classes in the downstream task can be also somewhat flexible as
 57 long as they are disconnected in the augmentation graph. For example, when the augmentation
 58 graph consists of m disconnected sub-graphs corresponding to fine-grained classes, our assumptions
 59 allow the downstream task to have any $r \leq m$ coarse-grained classes containing these fine-grained
 60 classes as a sub-partition. Prior work [49] on pseudolabeling algorithms essentially requires an exact
 61 alignment between sub-graphs and downstream classes (i.e., $r = m$). They face this limitation
 62 because their analysis requires fitting discrete pseudolabels on the unlabeled data. We avoid this
 63 difficulty because we consider directly learning continuous representations on the unlabeled data.

64 We apply spectral decomposition—a classical approach for graph partitioning, also known as spectral
 65 clustering [37, 39] in machine learning—to the adjacency matrix defined on the population augmen-
 66 tation graph. We form a matrix where the top- k eigenvectors are the columns and interpret each
 67 row of the matrix as the representation (in \mathbb{R}^k) of an example. Somewhat surprisingly, we show that
 68 this feature extractor can be also recovered (up to some linear transformation) by minimizing the

69 following population objective which is similar to the standard contrastive loss (Section 3.2):

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

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

74 We analyze the linear classification performance of the representations learned by minimizing the
 75 population spectral contrastive loss. Our main result (Theorem 3.7) shows that when the representation
 76 dimension exceeds the maximum number of disconnected sub-graphs, linear classification with
 77 learned representations is guaranteed to have a small error. Our theorem reveals a trend that a larger
 78 representation dimension is needed when there are a larger number of disconnected sub-graphs. Our
 79 analysis relies on novel techniques tailored to linear probe performance, which have not been studied
 80 in the spectral graph theory community to the best of our knowledge.

81 The spectral contrastive loss also works on empirical data. Since our approach optimizes parametric
 82 loss functions, guarantees involving the population loss can be converted to finite sample results
 83 using off-the-shelf generalization bounds. The sample complexity is polynomial in the Rademacher
 84 complexity of the model family and other relevant parameters (Theorem 4.1 and Theorem 4.2).

85 In summary, our main theoretical contributions are: 1) we propose a simple contrastive loss motivated
 86 by spectral decomposition of the population data graph, 2) under simple and realistic assumptions,
 87 we provide downstream classification guarantees for the representation learned by minimizing this
 88 loss on population data, and 3) our analysis is easily applicable to deep networks with polynomial
 89 unlabeled samples via off-the-shelf generalization bounds.

90 In addition, we implement and test the proposed spectral contrastive loss on standard vision benchmark
 91 datasets. We demonstrate that the features learned by our algorithm can match or outperform several
 92 strong baselines [12, 14, 15, 21] when evaluated using a linear probe.

93 2 Additional related works

94 **Empirical works on self-supervised learning.** Self-supervised learning algorithms have been
 95 shown to successfully learn representations that benefit downstream tasks [4, 6, 10, 12, 13, 15, 22–
 96 24, 35, 38, 41, 42, 50–52]. Many recent self-supervised learning algorithms learn features with
 97 siamese networks [8], where two neural networks of shared weights are applied to pairs of augmented
 98 data. Introducing asymmetry to siamese networks either with a momentum encoder like BYOL [21]
 99 or by stopping gradient propagation for one branch of the siamese network like SimSiam [14] has
 100 been shown to effectively avoid collapsing. Contrastive methods [12, 15, 22] minimize the InfoNCE
 101 loss [38], where two views of the same data are attracted while views from different data are repulsed.

102 **Theoretical works on self-supervised learning.** In addition to works [3, 28, 44, 45] discussed in
 103 the introduction, several other works [5, 43, 47, 48] also theoretically study self-supervised learning.
 104 The work Tsai et al. [47] prove that self-supervised learning methods can extract task-relevant
 105 information and discard task-irrelevant information, but lacks guarantees for solving downstream
 106 tasks efficiently with simple (e.g., linear) models. Tian et al. [43] study why non-contrastive self-
 107 supervised learning methods can avoid feature collapse. Cai et al. [9] analyze domain adaptation
 108 algorithms for subpopulation shift with a similar expansion condition as [49] while also allowing
 109 disconnected parts within each class, but require access to ground-truth labels during training. In
 110 contrast, our algorithm doesn’t need labels during pre-training.

111 3 Spectral contrastive learning on population data

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

114 We use $\overline{\mathcal{X}}$ to denote the set of all natural data (raw inputs without augmentation). We assume that
 115 each $\bar{x} \in \overline{\mathcal{X}}$ belongs to one of r classes, and let $y : \overline{\mathcal{X}} \rightarrow [r]$ denote the ground-truth (deterministic)
 116 labeling function. Let $\mathcal{P}_{\overline{\mathcal{X}}}$ be the population distribution over $\overline{\mathcal{X}}$ from which we draw training data and

117 test our final performance. For the ease of exposition, we assume $\bar{\mathcal{X}}$ to be a finite but exponentially
 118 large set (e.g., all real vectors in \mathbb{R}^d with bounded precision).¹

119 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
 120 the distribution of its augmentations. For instance, when \bar{x} represents an image, $\mathcal{A}(\cdot|\bar{x})$ can be the
 121 distribution of common augmentations [12] that includes Gaussian blur, color distortion and random
 122 cropping. We use \mathcal{X} to denote the set of all augmented data, which is the union of supports of all
 123 $\mathcal{A}(\cdot|\bar{x})$ 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
 124 denote $N = |\mathcal{X}|$.

125 We will learn an embedding function $f : \mathcal{X} \rightarrow \mathbb{R}^k$, and then evaluate its quality by the minimum
 126 error achieved with a linear probe. Concretely, a linear classifier has weights $B \in \mathbb{R}^{k \times r}$ and predicts
 127 $g_{f,B}(x) = \arg \max_{i \in [r]} (f(x)^\top B)_i$ for an augmented data x (arg max breaks tie arbitrarily). Then,
 128 given a natural data sample \bar{x} , we ensemble the predictions on augmented data and predict:

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

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

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

130 3.1 Augmentation graph and spectral decomposition

131 Our approach is based on the central concept of **population augmentation graph**, denoted by
 132 $G(\mathcal{X}, w)$, where the vertex set is all augmentation data \mathcal{X} and w denotes the edge weights defined
 133 below. For any two augmented data $x, x' \in \mathcal{X}$, define the weight $w_{xx'}$ as the marginal probability of
 134 generating the pair x and x' from a random natural data $\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}$:

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

135 Therefore, the weights sum to 1 because the total probability mass is 1: $\sum_{x,x' \in \mathcal{X}} w_{xx'} = 1$. The rela-
 136 tive magnitude intuitively captures the closeness between x and x' with respect to the augmentation
 137 transformation. For most of the unrelated x and x' , the value $w_{xx'}$ will be significantly smaller than
 138 the average value. For example, when x and x' are random croppings of a cat and a dog respectively,
 139 $w_{xx'}$ will be essentially zero because no natural data can be augmented into both x and x' . On the
 140 other hand, when x and x' are very close in ℓ_2 -distance or very close in ℓ_2 -distance up to color
 141 distortion, $w_{xx'}$ is nonzero because they may be augmentations of the same image with Gaussian
 142 blur and color distortion. We say that x and x' are connected with an edge if $w_{xx'} > 0$. See Figure 1
 143 (left) for more illustrations.

144 Given the structure of the population augmentation graph, we apply spectral decomposition to the
 145 population graph to construct principled embeddings. The eigenvalue problems are closely related to
 146 graph partitioning as shown in spectral graph theory [17] for both worst-case graphs [11, 25, 29, 33]
 147 and random graphs [1, 30, 34]. In machine learning, spectral clustering [37, 39] is a classical algorithm
 148 that learns embeddings by eigendecomposition on an empirical distance graph and invoking k -means
 149 on the embeddings.

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

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

154 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
 155 matrix with $D_{xx} = w_x$.²

156 Standard spectral graph theory approaches produce vertex embeddings as follows. Let $\gamma_1, \gamma_2, \dots, \gamma_k$
 157 be the k largest eigenvalues of \bar{A} , and v_1, v_2, \dots, v_k be the corresponding unit-norm eigenvectors.

¹This allows us to use sums instead of integrals and avoid non-essential nuances related to calculus.

²We index the matrix A, D by $(x, x') \in \mathcal{X} \times \mathcal{X}$. Generally we index N -dimensional axis by $x \in \mathcal{X}$.

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

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

By the classical theory on low-rank approximation (Eckart–Young–Mirsky theorem [19]), 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 $\hat{F} = F^* \cdot \text{diag}([\sqrt{\gamma_1}, \dots, \sqrt{\gamma_k}])Q$. 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.

The proof can be found in Section C.1.

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 [12] 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 data 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.³

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 a additive constant:*

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

³Though x and x' are simply two independent draws, we call them negative pairs following the literature [3].

200 The proof can be found in Section C.1.

201 We note that spectral contrastive loss is similar to many popular contrastive losses [12, 38, 40, 50]. For
 202 instance, the contrastive loss in SimCLR [12] can be rewritten as (with simple algebraic manipulation)

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

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

207 3.3 Theoretical guarantees for spectral contrastive loss on population data

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

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

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

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

215 We note that when S is a singleton, there is $\phi_G(S) = 1$ due to the definition of w_x . We introduce the
 216 sparsest m -partition to represent the number of edges between m disjoint subsets.

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

$$\rho_m := \min_{S_1, \dots, S_m} \max\{\phi_G(S_1), \dots, \phi_G(S_m)\}$$

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

220 When r is the number of underlying classes, we might expect $\rho_r \approx 0$ since the augmentations from
 221 different classes almost compose a disjoint r -way partition of \mathcal{X} . However, for $m > r$, we can expect
 222 ρ_m to be much larger. For instance, in the extreme case when $m = |\mathcal{X}| = N$, every set S_i is a
 223 singleton, which implies that $\rho_N = 1$.

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

227 **Assumption 3.5** (Labels are recoverable from augmentations). *Let $\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}$ and $y(\bar{x})$ be its label.*
 228 *Let the augmentation $x \sim \mathcal{A}(\cdot|\bar{x})$. We assume that there exists a classifier g that can predict $y(\bar{x})$*
 229 *given x with error at most α . That is, $g(x) = y(\bar{x})$ with probability at least $1 - \alpha$.*

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

232 **Assumption 3.6** (Realizability). *Let \mathcal{F} be a hypothesis class containing functions from \mathcal{X} to \mathbb{R}^k . We*
 233 *assume that at least one of the global minima of $\mathcal{L}(f)$ belongs to \mathcal{F} .*

234 Our main theorem bound from above the linear probe error of the features learned by minimizing the
 235 population spectral contrastive loss.

236 **Theorem 3.7.** *Assume the representation dimension $k \geq 2r$ and Assumption 3.5 holds for $\alpha > 0$.*
 237 *Let \mathcal{F} be a hypothesis class that satisfies Assumption 3.6 and let $f_{\text{pop}}^* \in \mathcal{F}$ be a minimizer of $\mathcal{L}(f)$.*
 238 *Then, we have*

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

Here we use $\tilde{O}(\cdot)$ to hide universal constant factors and logarithmic factor 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 small constant independent of k , whereas $\rho_{\lfloor k/2 \rfloor}$ increases with k and can be much larger than α when k is reasonably large. For instance, when there are t sub-graphs that have sufficient inner connections, we expect ρ_{t+1} to be on the order of a constant because any $t + 1$ partition needs to break one sub-graph into two pieces and incur a large conductance. We characterize the ρ_k 's growth on more concrete distributions in the next subsection.

Previous works on graph partitioning [2, 29, 31] often analyze the so-called 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 [2] 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 3.5) as well as a novel analysis technique that characterizes the linear structure of the representations. In Section C, we provide the proof of Theorem 3.7 as well as its more generalized version where $k/2$ is relaxed to be any constant fraction of k .

3.4 Provable instantiation of Theorem 3.7 to mixture of manifold data

In this section, we exemplify Theorem 3.7 on an example where the natural data distribution is a mixture of manifolds, and the augmentation transformation is adding Gaussian noise.

Example 3.8 (Mixture of manifolds). Suppose $\mathcal{P}_{\bar{\mathcal{X}}}$ is mixture of $r \leq d$ distributions P_1, \dots, P_r , where each P_i is generated by some κ -bi-Lipschitz⁴ generator $Q : \mathbb{R}^{d'} \rightarrow \mathbb{R}^d$ on some latent variable $z \in \mathbb{R}^{d'}$ with $d' \leq d$ which is 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} is $\bar{x} + \xi$ 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 the ground-truth label be the most likely mixture index i that generates x : $y(x) := \arg \max_i P_i(x)$. 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. We apply Theorem 3.7 and get the following theorem:

Theorem 3.9. When $k \geq 2r + 2$, Example 3.8 satisfies Assumption 3.5 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. The proof can be found in Section D.

4 Finite-sample generalization bounds

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

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

⁴A κ bi-Lipschitz function satisfies $\frac{1}{\kappa} \|f(x) - f(y)\|_2 \leq \|x - y\|_2 \leq \kappa \|f(x) - f(y)\|_2$.

It is worth noting that $\hat{\mathcal{L}}_n(f)$ is an unbiased estimator of the population spectral contrastive loss $\mathcal{L}(f)$. (See Claim E.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 $\hat{\mathcal{R}}_n(\mathcal{F}) := \max_{x_1, \dots, x_n \in \mathcal{X}} \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}, i \in [k]} \frac{1}{n} \left(\sum_{j=1}^n \sigma_j f_i(x_j) \right) \right]$,

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 E.1 bounds the population loss of a feature extractor trained with finite data:

Theorem 4.1. *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 , let $\hat{f}_{\text{emp}} \in \mathcal{F}$ be a minimizer of empirical loss $\hat{\mathcal{L}}_n(f)$. Then, with probability at least $1 - \delta$ over the randomness of data, we have*

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

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

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 E.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. *Assume representation dimension $k \geq 4r + 2$, Assumption 3.5 holds for $\alpha > 0$ and Assumption 3.6 holds. Recall γ_i be the i -th largest eigenvalue of the normalized adjacency matrix. Then, for any $\epsilon < \gamma_k^2$ 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.

This theorem shows that the error on the downstream task only grows linearly with the error ϵ during pretraining. We can relax Assumption 3.6 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 E.3.

5 Guarantee for learning linear probe with labeled data

In this section, we provide theoretical guarantees for learning a linear probe with *labeled* data. Theorem 3.7 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 data $x \sim \mathcal{A}(\cdot|\bar{x})$ 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 [36].

The following Theorem 5.1 provides a generalization guarantee for the linear classifier that minimizes capped quadratic loss on a labeled dataset. 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 3.7 where only small 0-1 error is guaranteed. Recall γ_i is the i -th largest eigenvalue of the normalized adjacency matrix. Given a labeled dataset $\{(\bar{x}_i, y(\bar{x}_i))\}_{i=1}^n$ 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 $i \in [n]$.

Theorem 5.1. *In the setting of Theorem 3.7, assume $\gamma_k \geq C_\lambda$ for some $C_\lambda > 0$. Learn a linear probe $\hat{B} \in \arg \min_{\|B\|_F \leq 1/C_\lambda} \sum_{i=1}^n \ell((f_{\text{pop}}^*(x_i), y(\bar{x}_i)), B)$ by minimizing the capped quadratic loss subject to a norm constraint. Then, with probability at least $1 - \delta$ over random data, we have*

$$\Pr_{\bar{x} \sim \mathcal{P}_{\mathcal{X}}} \left(\bar{g}_{f_{\text{pop}}^*, \hat{B}}(\bar{x}) \neq y(\bar{x}) \right) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log k + \frac{r}{C_\lambda} \cdot \sqrt{\frac{k}{n}} + \sqrt{\frac{\log 1/\delta}{n}}.$$

Here the first term is the population error from Theorem 3.7. The last two terms are the generalization gap from standard concentration inequalities for linear classification and are small when the number of labeled data n is polynomial in the feature dimension k . We note that this result reveals a trade-off when choosing the feature dimension k : when n is fixed, a larger k decreases the population contrastive loss while increases the generalization gap for downstream linear classification. The proof of Theorem 5.1 is in Section F.

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 [14] 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 [26] and Tiny-ImageNet [27] in Table 1. Our empirical results show that spectral contrastive learning achieves better performance than two popular baseline algorithms SimCLR [12] and SimSiam [14]. In Table 2 we report results on ImageNet [18] 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 [12]), momentum encoders (BYOL [21] and MoCo [22]) or additional tricks such as stop-gradient (SimSiam [14]).

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 [14]. 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 representations on downstream linear classification. We hope our study can facilitate future theoretical analyses of self-supervised learning and inspire new practical algorithms. For instance, one interesting future direction is to test the topology of the augmentation graph on empirical data distributions and design algorithms using tools from graph theory.

References

- [1] E. Abbe. Community detection and stochastic block models: recent developments, 2017.
- [2] S. Arora, S. Rao, and U. Vazirani. Expander flows, geometric embeddings and graph partitioning. *Journal of the ACM (JACM)*, 56(2):1–37, 2009.
- [3] S. Arora, H. Khandeparkar, M. Khodak, O. Plevrakis, and N. Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *arXiv preprint arXiv:1902.09229*, 2019.
- [4] P. Bachman, R. D. Hjelm, and W. Buchwalter. Learning representations by maximizing mutual information across views. *arXiv preprint arXiv:1906.00910*, 2019.
- [5] Y. Bansal, G. Kaplun, and B. Barak. For self-supervised learning, rationality implies generalization, provably. *arXiv preprint arXiv:2010.08508*, 2020.
- [6] A. Bardes, J. Ponce, and Y. LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- [7] S. 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.
- [8] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah. Signature verification using a "siamese" time delay neural network. *Advances in neural information processing systems*, 6: 737–744, 1993.
- [9] T. Cai, R. Gao, J. D. Lee, and Q. Lei. A theory of label propagation for subpopulation shift. *arXiv preprint arXiv:2102.11203*, 2021.
- [10] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020.
- [11] J. 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.
- [12] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [13] T. Chen, S. Kornblith, K. Swersky, M. Norouzi, and G. Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020.
- [14] X. Chen and K. He. Exploring simple siamese representation learning. *arXiv preprint arXiv:2011.10566*, 2020.
- [15] X. Chen, H. Fan, R. Girshick, and K. He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [16] C.-Y. Chuang, J. Robinson, L. Yen-Chen, A. Torralba, and S. Jegelka. Debiased contrastive learning. *arXiv preprint arXiv:2007.00224*, 2020.
- [17] F. R. Chung and F. C. Graham. *Spectral graph theory*. Number 92. American Mathematical Soc., 1997.
- [18] 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.
- [19] C. Eckart and G. Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- [20] N. Golowich, A. Rakhlin, and O. Shamir. Size-independent sample complexity of neural networks. In *Conference On Learning Theory*, pages 297–299. PMLR, 2018.
- [21] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.

- [22] K. He, H. Fan, Y. Wu, S. Xie, and R. 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.
- [23] O. Henaff. Data-efficient image recognition with contrastive predictive coding. In *International Conference on Machine Learning*, pages 4182–4192. PMLR, 2020.
- [24] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2018.
- [25] R. Kannan, S. Vempala, and A. Vetta. On clusterings: Good, bad and spectral. *Journal of the ACM (JACM)*, 51(3):497–515, 2004.
- [26] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. 2009.
- [27] Y. Le and X. Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7:7, 2015.
- [28] J. D. Lee, Q. Lei, N. Saunshi, and J. Zhuo. Predicting what you already know helps: Provable self-supervised learning. *arXiv preprint arXiv:2008.01064*, 2020.
- [29] J. R. Lee, S. O. Gharan, and L. Trevisan. Multiway spectral partitioning and higher-order cheeger inequalities. *Journal of the ACM (JACM)*, 61(6):1–30, 2014.
- [30] J. Lei, A. Rinaldo, et al. Consistency of spectral clustering in stochastic block models. *Annals of Statistics*, 43(1):215–237, 2015.
- [31] T. Leighton and S. Rao. Multicommodity max-flow min-cut theorems and their use in designing approximation algorithms. *Journal of the ACM (JACM)*, 46(6):787–832, 1999.
- [32] A. Louis and K. 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.
- [33] A. Louis, P. Raghavendra, P. Tetali, and S. 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.
- [34] F. McSherry. Spectral partitioning of random graphs. In *Proceedings 42nd IEEE Symposium on Foundations of Computer Science*, pages 529–537. IEEE, 2001.
- [35] I. Misra and L. v. d. 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.
- [36] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [37] A. Ng, M. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 14:849–856, 2001.
- [38] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [39] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- [40] K. 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.
- [41] Y. Tian, D. Krishnan, and P. Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.

- [42] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola. What makes for good views for contrastive learning. *arXiv preprint arXiv:2005.10243*, 2020.
- [43] Y. Tian, L. Yu, X. Chen, and S. Ganguli. Understanding self-supervised learning with dual deep networks. *arXiv preprint arXiv:2010.00578*, 2020.
- [44] C. Tosh, A. Krishnamurthy, and D. Hsu. Contrastive estimation reveals topic posterior information to linear models. *arXiv:2003.02234*, 2020.
- [45] C. Tosh, A. Krishnamurthy, and D. Hsu. Contrastive learning, multi-view redundancy, and linear models. In *Algorithmic Learning Theory*, pages 1179–1206. PMLR, 2021.
- [46] T. W. Tsai, C. Li, and J. Zhu. Mice: Mixture of contrastive experts for unsupervised image clustering. *arXiv preprint arXiv:2105.01899*, 2021.
- [47] Y.-H. H. Tsai, Y. Wu, R. Salakhutdinov, and L.-P. Morency. Self-supervised learning from a multi-view perspective. *arXiv preprint arXiv:2006.05576*, 2020.
- [48] T. Wang and P. Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR, 2020.
- [49] C. Wei, K. Shen, Y. Chen, and T. Ma. Theoretical analysis of self-training with deep networks on unlabeled data. *arXiv preprint arXiv:2010.03622*, 2020.
- [50] Z. Wu, Y. Xiong, S. X. Yu, and D. 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.
- [51] M. Ye, X. Zhang, P. C. Yuen, and S.-F. 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.
- [52] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section B.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) See Section B.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#)
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) We will include our code in the supplemental material before the supplemental material deadline.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Section 6 and Section A.

- 492 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 493 ments multiple times)? [No] Due to the constraint of time and computing resources,
 494 we couldn't run enough number of trials to compute the error bars. For instance, our
 495 experiment on ImageNet requires running on 8 GPUs for roughly 50 hours every single
 496 trial. We note that many prior works in the computer vision community also typically
 497 don't include error bars.
- 498 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 499 of GPUs, internal cluster, or cloud provider)? [Yes] See Section A
- 500 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 501 (a) If your work uses existing assets, did you cite the creators? [Yes] We use standard
 502 datasets. See Section 6.
- 503 (b) Did you mention the license of the assets? [N/A]
- 504 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 505
- 506 (d) Did you discuss whether and how consent was obtained from people whose data you're
 507 using/curating? [N/A]
- 508 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 509 information or offensive content? [N/A]
- 510 5. If you used crowdsourcing or conducted research with human subjects...
- 511 (a) Did you include the full text of instructions given to participants and screenshots, if
 512 applicable? [N/A]
- 513 (b) Did you describe any potential participant risks, with links to Institutional Review
 514 Board (IRB) approvals, if applicable? [N/A]
- 515 (c) Did you include the estimated hourly wage paid to participants and the total amount
 516 spent on participant compensation? [N/A]

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 [14].

B Limitations and potential negative social impacts

Limitations. This paper provides statistical guarantees for self-supervised learning. One limitation is that we don’t provide guarantees for optimization, and it is unclear whether standard optimization algorithms like SGD can provably reach the global minimum of the spectral contrastive loss. We believe this question is beyond the scope of this work but will be an interesting future direction.

Potential negative social impacts. In the long run, one possible negative impact of our research (and AI research in general) is leading to a large scale of job losses as humans are replaced by machines.

542 However, we believe that day is still far away given the current limitations of AI algorithms. We hope
 543 that our work can lead to more principled AI algorithms, which will overall benefit human society.

544 C Proofs for Section 3

545 We first prove a more generalized version of Theorem 3.7 in section C.2, and then prove Theorem 3.7
 546 in Section C.3.

547 C.1 Proofs of Lemma 3.1 and Lemma 3.2

548 *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
 549 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
 550 augmented data $x \in \bar{\mathcal{X}}$ with representation u_x and linear classifier B . Let $\tilde{B} = Q^{-1}B$, it's easy to
 551 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
 552 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
 553 data $x \in \mathcal{X}$. The equivalence of loss naturally follows. \square

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

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

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

559 C.2 A generalized version of Theorem 3.7

560 For the proof we will follow the convention in literature [29] and define the *normalized Laplacian*
 561 *matrix* as follows:

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

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

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

$$\phi^{\hat{y}} := \sum_{x,x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[\hat{y}(x) \neq \hat{y}(x')]. \quad (8)$$

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

$$\Delta(y, \hat{y}) := \Pr_{x \sim \mathcal{P}_{\bar{\mathcal{X}}}, \tilde{x} \sim \mathcal{A}(\cdot|x)} (\hat{y}(\tilde{x}) \neq y(x)). \quad (9)$$

572 Recall the spectral contrastive loss defined in Section 3.2 is:

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

573 We first state a more general version of Theorem 3.7 as follows.

Theorem C.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 $k \in \mathbb{Z}^+$. Let $k' \geq r$ such that $k + 1 = (1 + \zeta)k'$, where $\zeta \in (0, 1)$ and $k' \in \mathbb{Z}^+$. 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

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

where $\tilde{y}(\bar{x})$ 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 C.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 C.2 except for that the bound depends on λ_{k+1} :

Theorem C.3. Assume the set of augmented data \mathcal{X} is finite. Let $f_{\text{pop}}^* \in \arg \min_{f: \mathcal{X} \rightarrow \mathbb{R}^k} \mathcal{L}(f)$ be a minimizer of the population spectral contrastive loss $\mathcal{L}(f)$ with $k \in \mathbb{Z}^+$. 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

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

where $\tilde{y}(\bar{x})$ 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 C.3 to Section C.4.

To get rid of the dependency on λ_{k+1} , we use following higher-order Cheeger's inequality from [32].

Lemma C.4 (Proposition 1.2 in [32]). 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 C.2 by combining Theorem C.3 and Lemma C.4.

Proof of Theorem C.2. Let $G = (\mathcal{X}, w)$ be the augmentation graph. In Lemma C.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 C.3 finishes the proof. \square

C.3 Proof of Theorem 3.7

We will use the following lemma which gives a connection between $\phi^{\hat{y}}$, $\Delta(y, \hat{y})$ and Assumption 3.5.

605 **Lemma C.5.** Let $G = (\mathcal{X}, w)$ be the augmentation graph, r be the number of underlying classes.
 606 Let S_1, S_2, \dots, S_r be the partition induced by the classifier g in Assumption 3.5. Then, there exists
 607 an extended labeling function \hat{y} such that

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

608 and

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

609 *Proof of Lemma C.5.* We define function $\hat{y} : \mathcal{X} \rightarrow [r]$ as follows: for an augmented data $x \in \mathcal{X}$, we
 610 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 3.5 it is
 611 easy to see $\Delta(y, \hat{y}) \leq \alpha$. On the other hand, we have

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

612 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

613 Now we give the proof of Theorem 3.7 using Lemma C.5 and Theorem C.2.

614 *Proof of Theorem 3.7.* Let S_1, S_2, \dots, S_r be the partition of \mathcal{X} induced by the classifier g given in
 615 Assumption 3.5. Define function $\hat{y} : \mathcal{X} \rightarrow [r]$ as follows: for an augmented data $x \in \mathcal{X}$, we use func-
 616 tion $\hat{y}(x)$ 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 C.2, we
 617 have $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{f_{\text{pop}}, B^*}^*(x) \neq y(\bar{x})) \lesssim \log(k) \cdot \frac{\phi^{\hat{y}}}{\rho_{\lfloor k/2 \rfloor}^2} + \Delta(y, \hat{y})$. By Lemma C.5 we have
 618 $\phi^{\hat{y}} \leq 2\alpha$ and $\Delta(y, \hat{y}) \leq \alpha$, so we have $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{f_{\text{pop}}, B^*}^*(x) \neq y(\bar{x})) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log(k)$.
 619 Notice that by definition of ensembled linear probe predictor, $\bar{g}_{f_{\text{pop}}, B^*}^*(\bar{x}) \neq y(\bar{x})$ happens
 620 only if more than half of the augmentations of \bar{x} predicts differently from $y(\bar{x})$, so we have
 621 $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}} (\bar{g}_{f_{\text{pop}}, B^*}^*(\bar{x}) \neq y(\bar{x})) \leq 2 \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{f_{\text{pop}}, B^*}^*(x) \neq y(\bar{x})) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log(k)$. \square

622 C.4 Proof of Theorem C.3

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

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

629 **Lemma C.6.** Let L be the normalized Laplacian matrix of some graph G . Let $N = |\mathcal{X}|$ be total
 630 number of augmented data, v_i be the i -th smallest unit-norm eigenvector of L with eigenvalue λ_i
 631 (make them orthogonal in case of repeated eigenvalues). Let $R(u) := \frac{u^\top L u}{u^\top u}$ be the Rayleigh quotient
 632 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
 633 $b \in \mathbb{R}^k$ with norm $\|b\|_2 \leq \|u\|_2$ such that

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

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

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

635 We have

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

636 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

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

637 which finishes the proof. \square

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

640 **Claim C.7.** *In the setting of Lemma C.6, let \hat{y} be an extended labeling function. Fix $i \in [r]$. Define
641 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
642 following quantity:*

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

643 Then, we have

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

644 *Proof of Claim C.7.* Let f be any function $\mathcal{X} \rightarrow \mathbb{R}$, define function $u(x) := \sqrt{w_x} \cdot f(x)$. Let
645 $u \in \mathbb{R}^N$ be the vector corresponding to u . Let A be the adjacency matrix with $A_{xx'} = w_{xx'}$ and D
646 be the diagonal matrix with $D_{xx} = w_x$. By definition of Laplacian matrix, we have

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

647 Therefore we have

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

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

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

654 **Lemma C.8.** *Let \hat{F}_{mf} be the matrix form of a feature extractor $\hat{f}_{\text{mf}} : \mathcal{X} \rightarrow \mathbb{R}^k$. Then, \hat{F}_{mf} is a
655 minimizer of $\mathcal{L}_{\text{mf}}(F)$ if and only if*

$$\hat{f}(x) := \frac{1}{\sqrt{w_x}} \cdot \hat{f}_{\text{mf}}(x)$$

656 *is a minimizer of the population spectral contrastive loss $\mathcal{L}(f)$.*

657 *Proof of Lemma C.8.* Notice that

$$\begin{aligned}
\mathcal{L}_{\text{mf}}(F) &= \|(I - L) - FF^\top\|_F^2 \\
&= \sum_{x, x' \in \mathcal{X}} \left(\frac{w_{xx'}}{\sqrt{w_x w_{x'}}} - f(x)^\top f(x') \right)^2 \\
&= \sum_{x, x' \in \mathcal{X}} (f(x)^\top f(x'))^2 - 2 \sum_{x, x' \in \mathcal{X}} \frac{w_{xx'}}{\sqrt{w_x w_{x'}}} f(x)^\top f(x') + \|I - L\|_F^2. \quad (10)
\end{aligned}$$

658 Recall that the definition of spectral contrastive loss is

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

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

$$\begin{aligned}
\mathcal{L}(f) &= -2 \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot f(x)^\top f(x') + \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot (f(x)^\top f(x'))^2 \\
&= \sum_{x, x' \in \mathcal{X}} \left(\left(\frac{f(x)}{\sqrt{w_x}} \right)^\top \left(\frac{f(x')}{\sqrt{w_{x'}}} \right) \right)^2 - 2 \sum_{x, x' \in \mathcal{X}} \frac{w_{xx'}}{\sqrt{w_x w_{x'}}} \left(\frac{f(x)}{\sqrt{w_x}} \right)^\top \left(\frac{f(x')}{\sqrt{w_{x'}}} \right). \quad (11)
\end{aligned}$$

661 Compare Equation (10) and Equation (11), we see their minimizers only differ by a term $\sqrt{w_x}$, which
662 finishes the proof. \square

663 Note that the minimizer of matrix approximation loss is exactly the largest eigenvectors of $I - L$
664 (also the smallest eigenvectors of L) due to Eckart–Young–Mirsky theorem, Lemma C.8 indicates
665 that the minimizer of $\mathcal{L}(f)$ is equivalent to the smallest eigenvectors of L up to data-wise scaling.
666 Now we are ready to prove Theorem C.3 by combining Lemma C.6, Claim C.7 and Lemma C.8.

667 *Proof of Theorem C.3.* Let $F_{\text{sc}} = [v_1, v_2, \dots, v_k]$ be the matrix that contains the smallest k eigen-
668 vectors of L as columns, and $f_{\text{sc}} : \mathcal{X} \rightarrow \mathbb{R}^k$ is the corresponding feature extractor. For each $i \in [r]$,
669 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
670 Lemma C.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

$$\|u_i^{\hat{y}} - F_{\text{sc}} b_i\|_2^2 \leq \frac{R(u_i^{\hat{y}})}{\lambda_{k+1}} \|u_i^{\hat{y}}\|_2^2. \quad (12)$$

671 By Claim C.7, we have

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

672 So we can rewrite Equation (12) as:

$$\begin{aligned}
\|u_i^{\hat{y}} - F_{\text{sc}} b_i\|_2^2 &\leq \frac{\phi_i^{\hat{y}}}{2\lambda_{k+1}} \cdot \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\hat{y}(x) = i] \\
&= \frac{1}{2\lambda_{k+1}} \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[(\hat{y}(x) = i \wedge \hat{y}(x') \neq i) \text{ or } (\hat{y}(x) \neq i \wedge \hat{y}(x') = i)]. \quad (13)
\end{aligned}$$

673 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
674 feature extractor. Define matrix $B \in \mathbb{R}^{N \times k}$ such that $B^\top = [b_1, \dots, b_r]$. Summing Equation (13)
675 over all $i \in [r]$ and by the definition of $\phi^{\hat{y}}$ we have

$$\|U - F_{\text{sc}} B^\top\|_F^2 \leq \frac{1}{2\lambda_{k+1}} \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[\hat{y}(x) \neq \hat{y}(x')] = \frac{\phi^{\hat{y}}}{2\lambda_{k+1}}, \quad (14)$$

676 where

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

677 Now we come back to the feature extractor f_{pop}^* that minimizes the spectral contrastive loss function
 678 $\mathcal{L}(f)$. By Lemma C.8, function $f_{\text{mf}}^*(x) := \sqrt{w_x} \cdot f_{\text{pop}}^*(x)$ is a minimizer of matrix approximation
 679 loss. Let F_{mf}^* be the corresponding matrix. By Eckard-Young-Mirsky theorem, we have

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

680 where Q is an orthonormal matrix and

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

681 Let

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

682 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

$$\begin{aligned} & \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\|\vec{y}(\bar{x}) - B^* f_{\text{pop}}^*(x)\|_2^2 \right] \\ & \leq 2 \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\|\vec{\hat{y}}(x) - B^* f_{\text{pop}}^*(x)\|_2^2 \right] + 2 \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\|\vec{\hat{y}}(x) - \vec{y}(\bar{x})\|_2^2 \right] \\ & = 2 \sum_{x \in \mathcal{X}} w_x \cdot \|\vec{\hat{y}}(x) - B^* f_{\text{pop}}^*(x)\|_2^2 + 4\Delta(y, \hat{y}) \quad (\text{because } w_x \text{ is the probability of } x) \\ & = 2 \sum_{x \in \mathcal{X}} \|u(x) - B^* f_{\text{mf}}^*(x)\|_2^2 + 4\Delta(y, \hat{y}) \quad (\text{because } f_{\text{mf}}^*(x) = \sqrt{w_x} \cdot f_{\text{pop}}^*(x)) \\ & = 2 \|U - F_{\text{mf}}^* B^{*\top}\|_F^2 + 4\Delta(y, \hat{y}) \quad (\text{rewrite in matrix form}) \\ & = 2 \|U - F_{\text{sc}} B^{\top}\|_F^2 + 4\Delta(y, \hat{y}) \quad (\text{by definition of } B^*) \\ & \leq \frac{\phi^{\hat{y}}}{\lambda_{k+1}} + 4\Delta(y, \hat{y}). \quad (\text{by Equation (14)}) \end{aligned}$$

683 To bound the error rate, we first notice that $g_{F_{\text{mf}}^*, B^*}(x) \neq y(\bar{x})$ happens only if $\|\vec{y}(\bar{x}) - B^* f_{\text{sc}}(x)\|_2^2 \geq$
 684 $\frac{1}{2}$, we have for any $x \in \mathcal{X}$,

$$\|\vec{y}(\bar{x}) - B^* \hat{f}_{\text{ma}}(x)\|_2^2 \geq \frac{1}{2} \cdot \mathbb{1}[g_{F_{\text{mf}}^*, B^*}(x) \neq y(\bar{x})]. \quad (15)$$

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

$$\begin{aligned} & \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{f_{\text{pop}}^*, B^*}(x) \neq y(\bar{x}) \right) \\ & \leq 2 \mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left[\|\vec{y}(\bar{x}) - B^* f_{\text{pop}}^*(x)\|_2^2 \right] \quad (\text{by Equation (15)}) \\ & \leq \frac{2\phi^{\hat{y}}}{\lambda_{k+1}} + 8\Delta(y, \hat{y}). \end{aligned}$$

686 Finally we bound the norm of B^* as

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

687

□

D Proofs for Section 3.4

In this section, we give a proof of Theorem 3.9.

The following lemma shows that the augmented graph for Example 3.8 satisfies Assumption 3.5 with some bounded α .

Lemma D.1. *In the setting of Theorem 3.9, the data distribution satisfies Assumption 3.5 with $\alpha \leq \frac{1}{\text{poly}(d')}$.*

Proof of Lemma D.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

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

□

We use the following lemma to give a lower bound for the sparsest m -partition of the augmentation graph in Example 3.8.

Lemma D.2. *In the setting of Theorem 3.9, 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 D.2 can be found in Section D.1. Now we give the proof of Example 3.9.

Proof of Theorem 3.9. The result on α is directly from Lemma D.1. By concentration inequality, there must exists 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 D.2, we have $\rho_{k'} \gtrsim \frac{\sigma}{\kappa\sqrt{d}}$. Set $k' = \lfloor k/2 \rfloor$, we apply Theorem 3.7 and get the bound we need. □

713 **D.1 Proof of Lemma D.2**

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

717 **Claim D.3.** *In the setting of Theorem 3.9, given a set $S \subseteq \mathbb{R}^d$. If $x \in \mathbb{R}^d$ satisfies $\Pr(S|x) :=$
 718 $\Pr_{\tilde{x} \sim \mathcal{A}(\cdot|x)}(\tilde{x} \in S) \geq \frac{2}{3}$. Then, for any x' such that $\|x - x'\|_2 \leq \tau$, we have*

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

719 where

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

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

721 *Proof of Claim D.3.* By the definition of augmentation, we know

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

722 By the definition of c_σ , we have

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

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

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

724 Now we can bound the quantity of our interest:

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

725 □

726 We now give the proof of Lemma D.2.

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

$$\Pr_{x \sim P_i, \tilde{x} \sim \mathcal{A}(\cdot|x)}(\tilde{x} \in S_t) \leq \frac{1}{2}. \quad (16)$$

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

$$\rho_{k'} = \max_{i \in [k']} \phi_G(S_i) \geq \phi_G(S_1) \geq \min_{j \in [r]} \frac{\mathbb{E}_{x \sim P_j} [\Pr(S_1|x)(1 - \Pr(S_1|x))]}{\mathbb{E}_{x \sim P_j} [\Pr(S_1|x)]}, \quad (17)$$

730 where

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

731 WLOG, we assume $j = 1$ minimizes the RHS of Equation (17), so we only need to prove

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

732 We define the following set

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

733 Notice that

$$\begin{aligned} \mathbb{E}_{x \sim P_1} [\Pr(S_1|x)] &= \int_x P_1(x) \Pr(S_1|x) dx \\ &= \int_{x \in R} P_1(x) \Pr(S_1|x) dx + \int_{x \notin R} P_1(x) \Pr(S_1|x) dx. \end{aligned} \quad (18)$$

734 We can consider the following two cases.

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

736 This is the easy case because we have

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

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

738 Define neighbourhood of R as

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

739 We have

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

740 where the second inequality is by Claim D.3. Notice that

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

741 where we use Equation (16). Define set $\tilde{R} := Q^{-1}(R)$ be the set in the ambient space corresponding
742 to R . Define

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

743 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
744 isoperimetric inequality [7], we have

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

745 where

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

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

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

747 So we have

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

748 By Equation (18), either case 1 or case 2 holds. Combining case 1 and case 2, we have

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

749

□

750 E Proofs for Section 4

751 E.1 Proof of Theorem 4.1

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

753 **Definition E.1** (Empirical spectral contrastive loss). *Consider a dataset $\hat{\mathcal{X}} = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\}$*
 754 *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*
 755 *$\hat{P}_{\bar{x}, \bar{x}'}$ be the uniform distribution over data pairs (\bar{x}_i, \bar{x}_j) where $i \neq j$. We define the empirical*
 756 *spectral contrastive loss of a feature extractor f as*

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

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

759 **Claim E.2.** $\hat{\mathcal{L}}_n(f)$ is an unbiased estimator of $\mathcal{L}(f)$, i.e.,

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

760 *Proof.* This is because

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

761

□

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

764 **Definition E.3.** Given dataset $\hat{\mathcal{X}}$, we sample a subset of tuples as follows: first sample a permutation
 765 $\pi : [n] \rightarrow [n]$, then we sample tuples $S = \{(z_i, z_i^+, z_i')\}_{i=1}^{n/2}$ as follows:

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

766 We define the following loss on S :

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

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

768 **Claim E.4.** For given $\hat{\mathcal{X}}$, if we sample S as above, we have:

$$\mathbb{E}_S [\hat{\mathcal{L}}_S(f)] = \hat{\mathcal{L}}_f(\cdot)$$

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

770 The following lemma reveals the relationship between the Rademacher complexity of feature extrac-
 771 tors and the Rademacher complexity of the loss defined on tuples:

772 **Lemma E.5.** Let \mathcal{F} be a hypothesis class of feature extractors from \mathcal{X} to \mathbb{R}^k . Assume $\|f(x)\|_\infty \leq \kappa$
 773 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
 774 of $f(x)$. Let \mathcal{F}_i be the hypothesis containing f_i for all $f \in \mathcal{F}$. For $m \in \mathbb{Z}^+$, let $\hat{\mathcal{R}}_m(\mathcal{F}_i)$ be the
 775 maximal possible empirical Rademacher complexity of \mathcal{F}_i over m data:

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

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

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

Proof.

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

here the second inequality is by Talagrand's lemma. Notice that for any $z_1, z_2 \dots z_m$ and z'_1, z'_2, \dots, z'_m in \mathcal{X} and any $i \in [k]$ we have

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

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

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

781

□

Proof of Theorem 4.1. By Claim E.2 and Claim E.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 E.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 E.5 and notice the fact that $(f(z_j)^\top f(z'_j))^2 - 2f(z_j)^\top f(z'_j)$ 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 $\widehat{\mathcal{X}}$ and S , we have for any $f \in \mathcal{F}$,

$$\mathcal{L}(f) \leq \widehat{\mathcal{L}}_S(f) + (32k^2\kappa^2 + 32k\kappa) \max_{i \in [k]} \widehat{\mathcal{R}}_{n/2}(\mathcal{F}_i) + (4k\kappa^2 + k^2\kappa^4) \cdot \sqrt{\frac{4 \log 2/\delta}{n}}. \quad (19)$$

This means with probability at least $1 - \delta/2$ over random $\widehat{\mathcal{X}}$, we have: with probability at least $1 - \delta/2$ over random tuples S conditioned on $\widehat{\mathcal{X}}$, Equation (19) holds. Since both $\mathcal{L}(f)$ and $\widehat{\mathcal{L}}_n(f)$ 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 $\widehat{\mathcal{X}}$, we have for any $f \in \mathcal{F}$,

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

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 $\widehat{\mathcal{X}}$, we have for any $f \in \mathcal{F}$,

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

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

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

where \hat{f} is minimizer of $\widehat{\mathcal{L}}_n(f)$ 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 $\max_{i \in [k]} \widehat{\mathcal{R}}_{n/2}(\mathcal{F}_i) = \widehat{\mathcal{R}}_{n/2}(\mathcal{F})$ finishes the proof. □

E.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 [20], which gives the following theorem.

Theorem E.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

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

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 $2n$, we have

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

where \hat{f} is the minimizer of $\hat{\mathcal{L}}_{2n}(f)$ 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 E.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 [20], we have

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

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

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

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

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

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

E.3 Proof of Theorem 4.2

In this section we give the proof of Theorem 4.2. We first introduce the following definitions of ϵ -optimal minimizers of matrix approximation loss and population spectral contrastive loss:

Definition E.7. We say a function \hat{f}_{mf} is ϵ -optimal minimizer of matrix approximation loss \mathcal{L}_{mf} if

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

where \widehat{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

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

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

828 **Theorem E.8.** [Generalization of Theorem C.3] Assume the set of augmented data \mathcal{X} is finite. Let
829 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
830 a ϵ -optimal minimizer of the spectral contrastive loss function $\mathcal{L}(f)$ with $k \in \mathbb{Z}^+$. Then, for any
831 labeling function $\hat{y} : \mathcal{X} \rightarrow [r]$ there exists a linear probe $\hat{B} \in \mathbb{R}^{r \times k}$ such that

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

832 where $\phi^{\hat{y}}$ and $\Delta(y, \hat{y})$ are defined in Equations 8 and 9 respectively.

833 The proof of lemma E.8 is deferred to Section E.4.

834 Now we are ready to prove Theorem 4.2 using Theorem E.8.

835 *Proof of Theorem 4.2.* In Theorem E.8 we let $k' = \lfloor \frac{3}{4}k \rfloor$ on the RHS of the bound and get: for any
836 $\hat{y} : \mathcal{X} \rightarrow [r]$ there exists $\hat{B} \in \mathbb{R}^{r \times k}$ such that

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

837 Let S_1, S_2, \dots, S_r be the partition of \mathcal{X} induced by the classifier g in Assumption 3.5. Define
838 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
839 the index of set that x is in, i.e., $x \in S_{\hat{y}(x)}$. Then by Lemma C.5 we have $\phi^{\hat{y}} \leq 2\alpha$ and $\Delta(y, \hat{y}) \leq \alpha$.

840 In Lemma C.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
841 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

842 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

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

843 Notice that by the definition of ensembled linear probe predictor, $\bar{g}_{\hat{f}, \hat{B}}(\bar{x}) \neq y(\bar{x})$ happens
844 only if more than half of the augmentations of \bar{x} predicts differently from $y(\bar{x})$, so we have

845 $\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}} \left(\bar{g}_{\hat{f}, \hat{B}} \neq y(\bar{x}) \right) \leq 2 \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} \left(g_{\hat{f}, \hat{B}}(x) \neq y(\bar{x}) \right)$ which finishes the proof. \square

846 E.4 Proof of Theorem E.8

847 In this section, we give the proof for Theorem E.8.

848 **Lemma E.9** (Generalization of Lemma C.8). Let \hat{F}_{mf} be the matrix form of a feature extractor
849 $\hat{f}_{\text{mf}} : \mathcal{X} \rightarrow \mathbb{R}^k$. Then, \hat{F}_{mf} is a ϵ -optimal minimizer of $\mathcal{L}_{\text{mf}}(F)$ if and only if

$$\hat{f}(x) := \frac{1}{\sqrt{w_x}} \cdot \hat{f}_{\text{mf}}(x)$$

850 is a ϵ -optimal minimizer of spectral contrastive loss $\mathcal{L}(f)$.

851 *Proof of Lemma E.9.* The proof follows the proof of Lemma C.8. \square

852 We will use the following important lemma about ϵ -optimal minimizer of \mathcal{L}_{mf} :

Lemma E.10. Let λ_i be the i -th minimal eigenvalue of the normalized Laplacian matrix L with corresponding unit-norm eigenvector v_i . Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be ϵ -optimal minimizer of \mathcal{L}_{mf} where $\epsilon < (1 - \lambda_k)^2$. Let $F \in \mathbb{R}^{N \times k}$ be the matrix form of f , where $N = |\mathcal{X}|$. 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. For function f with matrix form F , we overload notation $\mathcal{L}_{\text{mf}}(\cdot)$ and use $\mathcal{L}_{\text{mf}}(f)$ to represent $\mathcal{L}_{\text{mf}}(F)$.

We first prove that the column rank of f is k . If the column rank of f is less than k , then there must exists some function $f' : \mathcal{X} \rightarrow \mathbb{R}^{k-1}$ such that $\mathcal{L}_{\text{mf}}(f') = \mathcal{L}_{\text{mf}}(f)$. According to the Eckart–Young–Mirsky Theorem, we have $\min_{f: \mathcal{X} \rightarrow \mathbb{R}^k} \mathcal{L}_{\text{mf}}(f) = \sum_{j=k+1}^N (1 - \lambda_j)^2$ and $\min_{f: \mathcal{X} \rightarrow \mathbb{R}^{k-1}} \mathcal{L}_{\text{mf}}(f) = \sum_{j=k}^N (1 - \lambda_j)^2$. Therefore, $\mathcal{L}_{\text{mf}}(f) = \mathcal{L}_{\text{mf}}(f') \geq (1 - \lambda_k)^2 + \min_{f: \mathcal{X} \rightarrow \mathbb{R}^k} \mathcal{L}_{\text{mf}}(f)$, contradicting with $\epsilon < (1 - \lambda_k)^2$. As a result, the column rank of f has to be k .

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

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

We can rewrite the above lower bound as

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

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:

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

882 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

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

883 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
884 S_j when $j \geq k + 1$. Now notice that

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

885 and also

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

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

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

887 where f_{pop}^* is the minimizer of \mathcal{L}_{mf} , and the last equality is by Eckart–Young–Mirsky Theorem. So we

888 know $\mathcal{L}_{\text{mf}}(f) \geq ((1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2) \left\| \Pi_f^\perp v_i \right\|_2^2 + \mathcal{L}_{\text{mf}}(f_{\text{pop}}^*)$, which implies that $\left\| \Pi_f^\perp v_i \right\|_2^2 \leq$

889 $\frac{\epsilon}{(1 - \lambda_i)^2 - (1 - \lambda_{k+1})^2} \leq \frac{\epsilon}{(\lambda_{k+1} - \lambda_i)^2}.$ \square

890 The following lemma generalizes Lemma C.6.

891 **Lemma E.11** (Generalization of Lemma C.6). *Let L be the normalized Laplacian matrix of graph*
892 *$G = (\mathcal{X}, w)$, where $|\mathcal{X}| = N$. Let $f : \mathcal{X} \rightarrow \mathbb{R}^k$ be an ϵ -optimal minimizer of $\mathcal{L}_{\text{mf}}(f)$ where*
893 *$\epsilon < (1 - \lambda_k)^2$. Let F be the matrix form of f and F_i is the i -th column of F . Let $R(u) := \frac{u^\top L u}{u^\top u}$ be*
894 *the Rayleigh quotient of a vector $u \in \mathbb{R}^N$. Then, for any $k \in \mathbb{Z}^+$ such that $k < N$, there exists a*
895 *vector $b \in \mathbb{R}^k$ such that*

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

896 *Proof.* Let k' be the choice that minimizes the right hand side. We use $p_v(u)$ to denote the projection
897 of u onto the span of $v_1, \dots, v_{k'}$. We use $p_{v,f}(u)$ to denote the projection of $p_v(u)$ onto the span of
898 $f_1, \dots, f_{k'}$. Then we know that

$$\|u - p_{v,f}(u)\|_2^2 \leq 2 \|u - p_v(u)\|_2^2 + 2 \|p_v(u) - p_{v,f}(u)\|_2^2. \quad (20)$$

899 By the proof of Lemma C.6, we know that

$$\|u - p_v(u)\|_2^2 \leq \frac{R(u)}{\lambda_{k'+1}} \|u\|_2^2. \quad (21)$$

900 On the other hand, we have

$$\begin{aligned}
\|p_v(u) - p_{vf}(u)\|_2^2 &= \|\Pi_f^\perp p_v(u)\|_2^2 \\
&= \left\| \sum_{i=1}^{k'} \Pi_f^\perp v_i v_i^\top u \right\|_2^2 \\
&\leq \left(\sum_{i=1}^{k'} \|\Pi_f^\perp v_i\|_2^2 \right) \cdot \left(\sum_{i=1}^{k'} (v_i^\top u)^2 \right) \\
&\leq \frac{k' \epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \|u\|_2^2,
\end{aligned} \tag{22}$$

901 where the first inequality is by Cauchy–Schwarz inequality and the second inequality is by
902 Lemma E.10. Plugging Equation (21) and Equation (22) into Equation (20) finishes the proof. \square

903 Now we prove Theorem E.8 using the above lemmas.

904 *Proof of Theorem E.8.* By Lemma E.9, the ϵ -optimal minimizer of $\mathcal{L}(f)$ is only different from ϵ -
905 optimal minimizer of $\mathcal{L}_{\text{mf}}(f)$ by a positive constant for each x . Since this difference won't influence
906 the prediction accuracy, we only need to prove this theorem assuming \hat{f} is ϵ -optimal minimizer of
907 $\mathcal{L}_{\text{mf}}(f)$.

908 For each $i \in [r]$, we define the function $u_i(x) = \mathbb{1}[\hat{y}(x) = i] \cdot \sqrt{w_x}$. Let $u : \mathcal{X} \rightarrow \mathbb{R}^k$ be the function
909 such that $u(x)$ has u_i at the i -th dimension. By Lemma E.11, there exists a vector $b_i \in \mathbb{R}^k$ such that

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

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

$$\begin{aligned}
\|U - \hat{F} \hat{B}^\top\|_F^2 &\leq \sum_{i=1}^r \min_{1 \leq k' \leq k} \left(\frac{2R(u_i)}{\lambda_{k'+1}} + \frac{2k' \epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) \|u_i\|_2^2 \\
&\leq \min_{1 \leq k' \leq k} \sum_{i=1}^r \left(\frac{2R(u_i)}{\lambda_{k'+1}} \|u_i\|_2^2 + \frac{2k' \epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \|u_i\|_2^2 \right).
\end{aligned} \tag{23}$$

912 Notice that

$$\begin{aligned}
\sum_{i=1}^r R(u_i) \|u_i\|_2^2 &= \sum_{i=1}^r \frac{1}{2} \phi_i^{\hat{y}} \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\hat{y}(x) = i] \\
&= \frac{1}{2} \sum_{i=1}^r \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[(\hat{y}(x) = i \wedge \hat{y}(x') \neq i) \text{ or } (\hat{y}(x) \neq i \wedge \hat{y}(x') = i)] \\
&= \frac{1}{2} \sum_{x, x' \in \mathcal{X}} w_{xx'} \cdot \mathbb{1}[\hat{y}(x) \neq \hat{y}(x')] = \frac{1}{2} \phi^{\hat{y}},
\end{aligned} \tag{24}$$

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

$$\sum_{i=1}^r \|u_i\|_2^2 = \sum_{i=1}^r \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1}[\hat{y}(x) = i] = \sum_{x \in \mathcal{X}} w_x = 1. \tag{25}$$

914 Plugging Equation (24) and Equation (25) into Equation (23) gives us

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

915 Notice that by definition of $u(x)$, we know that prediction $g_{\hat{f}, \hat{B}}(x) \neq \hat{y}(x)$ only happens if

916 $\|u(x) - \hat{B}\hat{f}(x)\|_2^2 \geq \frac{w_x}{2}$. Hence we have

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

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

$$\begin{aligned} \Pr_{x \sim \mathcal{X}}(g_{\hat{f}, \hat{B}}(x) \neq \hat{y}(x)) &= \sum_{x \in \mathcal{X}} w_x \cdot \mathbb{1} [g_{\hat{f}, \hat{B}}(x) \neq \hat{y}(x)] \\ &\leq 2 \cdot \|U - \hat{F}\hat{B}^\top\|_F^2 \leq \min_{1 \leq k' \leq k} \left(\frac{2\phi^{\hat{y}}}{\lambda_{k'+1}} + \frac{4k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right). \end{aligned}$$

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

$$\begin{aligned} \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{\hat{f}, \hat{B}}(x) \neq y(\bar{x})) &\leq \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{\hat{f}, \hat{B}}(x) \neq \hat{y}(x)) + \Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} (y(\bar{x}) \neq \hat{y}(x)) \\ &\leq \min_{1 \leq k' \leq k} \left(\frac{2\phi^{\hat{y}}}{\lambda_{k'+1}} + \frac{4k'\epsilon}{(\lambda_{k+1} - \lambda_{k'})^2} \right) + \Delta(y, \hat{y}). \end{aligned}$$

920

□

921 F Proofs for Section 5

922 In this section we give the proof of Theorem 5.1. We first introduce the following lemma, which
923 states the expected norm of representations:

924 **Lemma F.1.** *Let $f_{\text{pop}}^* : \mathcal{X} \rightarrow \mathbb{R}^k$ be a minimizer of population spectral contrastive loss $\mathcal{L}(f)$. Then,*
925 *we have*

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} [\|f_{\text{pop}}^*(x)\|_2^2] \leq k. \quad (26)$$

926 *Proof of Lemma F.1.* By Lemma C.8 and the definition of w , we have

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} [\|f_{\text{pop}}^*(x)\|_2^2] = \sum_{x \in \mathcal{X}} w_x \|f_{\text{pop}}^*(x)\|_2^2 = \sum_{x \in \mathcal{X}} \|\hat{f}_{\text{ma}}(x)\|_2^2 = \|\hat{F}_{\text{mf}}\|_F^2, \quad (27)$$

927 where \hat{F}_{mf} is a minimizer of the matrix approximation loss defined in Section 3.2. By Eckard-Young-
928 Mirsky theorem, \hat{F}_{mf} looks like

$$\hat{F}_{\text{mf}} = F_{\text{sc}} D_\lambda Q,$$

929 where $F_{\text{sc}} = [v_1, v_2, \dots, v_k]$ contains the k smallest eigenvectors of the laplacian matrix L as
930 columns, Q is an orthonormal matrix and

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

931 So we have

$$\|\hat{F}_{\text{mf}}\|_F^2 = \text{Tr} (F_{\text{sc}} D_\lambda^2 F_{\text{sc}}^\top) \leq \text{Tr} (F_{\text{sc}} F_{\text{sc}}^\top) = k,$$

932 where we use the fact that D_λ has diagonal values less than 1 and v_i is unit-norm. Pluggin this into
933 Equation (27) finishes the proof. □

934 Now we give the proof of Theorem 5.1:

935 *Proof of Theorem 5.1.* Let f_{pop}^* be the minimizer of population spectral contrastive loss $\mathcal{L}(f)$. We
 936 abuse notation and use y_i to denote $y(\bar{x}_i)$, and let $z_i = f_{\text{pop}}^*(x_i)$. We first study the average empirical
 937 Rademacher complexity of the capped quadratic loss on a dataset $\{(z_i, y_i)\}_{i=1}^n$, where (z_i, y_i) is
 938 sampled as in Section 5:

$$\begin{aligned}\widehat{\mathcal{R}}_n(\ell) &:= \mathbb{E}_{\{(z_i, y_i)\}_{i=1}^n} \mathbb{E}_{\sigma} \left[\sup_{\|B\|_F \leq 1/C_\lambda} \frac{1}{n} \left[\sum_{i=1}^n \sigma_i \ell((z_i, y_i), B) \right] \right] \\ &\leq 2r \mathbb{E}_{\{(z_i, y_i)\}_{i=1}^n} \mathbb{E}_{\sigma} \left[\sup_{\|b\|_2 \leq 1/C_\lambda} \frac{1}{n} \left[\sum_{i=1}^n \sigma_i w^\top z_i \right] \right] \\ &\leq \frac{2r}{C_\lambda} \sqrt{\frac{\mathbb{E}[\|z_i\|^2]}{n}} \leq \frac{2r\sqrt{k}}{C_\lambda\sqrt{n}},\end{aligned}$$

939 where the first inequality uses Talagrand's lemma and the fact that ℓ_σ is 2-Lipschitz, the second in-
 940 equality is by standard Rademacher complexity of linear models, the third inequality is by Lemma F.1.

941 By Theorem C.2, there exists a linear probe B^* with norm bound $\|B^*\|_F \leq 1/(1 - \lambda_k) \leq 1/C_\lambda$
 942 such that

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

943 where $(1 + \zeta)k' = k + 1$. Let \widehat{B} be the minimizer of $\sum_{i=1}^n \ell((z_i, y_i), B)$ subject to $\|B\|_F \leq 1/C_\lambda$,
 944 then by standard generalization bound, we have: with probability at least $1 - \delta$, we have

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} [\ell((f_{\text{pop}}^*(x), y(\bar{x})), \widehat{B})] \lesssim \text{poly}(1/\zeta) \log(k+1) \cdot \frac{\phi^{\hat{y}}}{\rho_{k'}^2} + \Delta(y, \hat{y}) + \frac{r\sqrt{k}}{C_\lambda\sqrt{n}} + \sqrt{\frac{\log 1/\delta}{n}}.$$

945 Follow the same steps as in the proof of Theorem 3.7, we can get a genalization bound of

$$\mathbb{E}_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} [\ell((f_{\text{pop}}^*(x), y(\bar{x})), \widehat{B})] \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log k + \frac{r}{C_\lambda} \cdot \sqrt{\frac{k}{n}} + \sqrt{\frac{\log 1/\delta}{n}}.$$

946 Notice that $y(\bar{x}) \neq g_{f_{\text{pop}}^*, \widehat{B}}(x)$ only if $\ell((f_{\text{pop}}^*(x), y(\bar{x})), \widehat{B}) \geq \frac{1}{2}$, we have the error bound

$$\Pr_{\bar{x} \sim \mathcal{P}_{\bar{\mathcal{X}}}, x \sim \mathcal{A}(\cdot|\bar{x})} (g_{f_{\text{pop}}^*, \widehat{B}}(x) \neq y(\bar{x})) \lesssim \frac{\alpha}{\rho_{\lfloor k/2 \rfloor}^2} \cdot \log k + \frac{r}{C_\lambda} \cdot \sqrt{\frac{k}{n}} + \sqrt{\frac{\log 1/\delta}{n}}.$$

947 The result on $\bar{g}_{f_{\text{pop}}^*, \widehat{B}}$ naturally follows by the definition of \bar{g} . □