

# Learning Diverse and Discriminative Representations via the Principle of Maximal Coding Rate Reduction

Yaodong Yu<sup>†</sup> Kwan Ho Ryan Chan<sup>†</sup> Chong You<sup>†</sup> Chaobing Song<sup>‡</sup> Yi Ma<sup>†</sup>

<sup>†</sup>Department of EECS, University of California, Berkeley

<sup>‡</sup>Tsinghua-Berkeley Shenzhen Institute, Tsinghua University

## Abstract

To learn intrinsic low-dimensional structures from high-dimensional data that most discriminate between classes, we propose the principle of *Maximal Coding Rate Reduction* (MCR<sup>2</sup>), an information-theoretic measure that maximizes the coding rate difference between the whole dataset and the sum of each individual class. We clarify its relationships with most existing frameworks such as cross-entropy, information bottleneck, information gain, contractive and contrastive learning, and provide theoretical guarantees for learning diverse and discriminative features. The coding rate can be accurately computed from finite samples of degenerate subspace-like distributions and can learn intrinsic representations in supervised, self-supervised, and unsupervised settings in a unified manner. Empirically, the representations learned using this principle alone are significantly more robust to label corruptions in classification than those using cross-entropy, and can lead to state-of-the-art results in clustering mixed data from self-learned invariant features.

## 1 Context and Motivation

Given a random vector  $\mathbf{x} \in \mathbb{R}^D$  which is drawn from a mixture of, say  $k$ , distributions  $\mathcal{D} = \{\mathcal{D}_j\}_{j=1}^k$ , one of the most fundamental problems in machine learning is how to effectively and efficiently *learn the distribution* from a finite set of i.i.d samples, say  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$ . To this end, we *seek a good representation* through a continuous mapping,  $f(\mathbf{x}, \theta) : \mathbb{R}^D \rightarrow \mathbb{R}^d$ , that captures intrinsic structures of  $\mathbf{x}$  and best facilitates subsequent tasks such as classification or clustering.

**Supervised learning of discriminative representations.** To ease the task of learning  $\mathcal{D}$ , in the popular supervised setting, a true class label, represented as a one-hot vector  $\mathbf{y}_i \in \mathbb{R}^k$ , is given for each sample  $\mathbf{x}_i$ . Extensive studies have shown that for many practical datasets (images, audios, and natural languages, etc.), the mapping from the data  $\mathbf{x}$  to its class label  $\mathbf{y}$  can be effectively modeled by training a deep network [GBC16], here denoted as  $f(\mathbf{x}, \theta) : \mathbf{x} \mapsto \mathbf{y}$  with network parameters  $\theta \in \Theta$ . This is typically done by minimizing the *cross-entropy loss* over a training set  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$ , through backpropagation over the network parameters  $\theta$ :

$$\min_{\theta \in \Theta} \text{CE}(\theta, \mathbf{x}, \mathbf{y}) \doteq -\mathbb{E}[\langle \mathbf{y}, \log[f(\mathbf{x}, \theta)] \rangle] \approx -\frac{1}{m} \sum_{i=1}^m \langle \mathbf{y}_i, \log[f(\mathbf{x}_i, \theta)] \rangle. \quad (1)$$

Despite its effectiveness and enormous popularity, there are two serious limitations with this approach: 1) It aims only to predict the labels  $\mathbf{y}$  even if they might be mislabeled. Empirical studies show that deep networks, used as a “black box,” can even fit random labels [ZBH<sup>+</sup>17]. 2) With such an end-to-end data fitting, it is not clear to what extent the intermediate features learned by the network capture the intrinsic structures of the data that make meaningful classification possible in the first place.<sup>1</sup> The precise geometric and statistical properties of the learned features are also often

\*The first two authors contributed equally to this work.

<sup>1</sup>despite plenty of empirical efforts in trying to illustrate or interpreting the so-learned features [ZF14].

obscured, which leads to the lack of interpretability and subsequent performance guarantees (e.g., generalizability, transferability, and robustness, etc.) in deep learning. Therefore, the goal of this paper is to address such limitations of current learning frameworks by reformulating the objective towards learning *explicitly meaningful* representations for the data  $\mathbf{x}$ .

**Minimal discriminative features via information bottleneck.** One popular approach to interpret the role of deep networks is to view outputs of intermediate layers of the network as selecting certain latent features  $\mathbf{z} = f(\mathbf{x}, \theta) \in \mathbb{R}^d$  of the data that are discriminative among multiple classes. Learned representations  $\mathbf{z}$  then facilitate the subsequent classification task for predicting the class label  $\mathbf{y}$  by optimizing a classifier  $g(\mathbf{z})$ :

$$\mathbf{x} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{z}(\theta) \xrightarrow{g(\mathbf{z})} \mathbf{y}.$$

The *information bottleneck* (IB) formulation [TZ15] further hypothesizes that the role of the network is to learn  $\mathbf{z}$  as the minimal sufficient statistics for predicting  $\mathbf{y}$ . Formally, it seeks to maximize the mutual information  $I(\mathbf{z}, \mathbf{y})$ <sup>2</sup> between  $\mathbf{z}$  and  $\mathbf{y}$  while minimizing  $I(\mathbf{x}, \mathbf{z})$  between  $\mathbf{x}$  and  $\mathbf{z}$ :

$$\max_{\theta \in \Theta} \text{IB}(\mathbf{x}, \mathbf{y}, \mathbf{z}(\theta)) \doteq I(\mathbf{z}(\theta), \mathbf{y}) - \beta I(\mathbf{x}, \mathbf{z}(\theta)), \quad \beta > 0. \quad (2)$$

This framework has been successful in describing certain behaviors of deep networks.<sup>3</sup> But by being task-dependent (depending on the label  $\mathbf{y}$ ) and seeking a *minimal* set of most informative features for the task at hand (for predicting the label  $\mathbf{y}$  only), the network sacrifices generalizability, robustness, or transferability.<sup>4</sup> To address this, our framework uses label  $\mathbf{y}$  only as side information to assist learning discriminative features, hence making learned features more robust to mislabeled data.

**Contractive learning of generative representations.** Complementary to the above supervised discriminative approach, *auto-encoding* [BH89, Kra91] is another popular *unsupervised* (label-free) framework used to learn good latent representations. The idea is to learn a compact latent representation  $\mathbf{z} \in \mathbb{R}^d$  that adequately regenerates the original data  $\mathbf{x}$  to certain extent, say through optimizing some decoder or generator  $g(\mathbf{z}, \eta)$ <sup>5</sup>:

$$\mathbf{x} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{z}(\theta) \xrightarrow{g(\mathbf{z}, \eta)} \hat{\mathbf{x}}(\theta, \eta). \quad (3)$$

Typically, such representations are learned in an end-to-end fashion by imposing certain heuristics on geometric or statistical “compactness” of  $\mathbf{z}$ , such as its dimension, energy, or volume. For example, the *contractive* autoencoder [RVM<sup>+</sup>11] penalizes local volume expansion of learned features approximated by the Jacobian  $\|\frac{\partial \mathbf{z}}{\partial \theta}\|$ . Another key design factor of this approach is the choice of a proper, but often elusive, metric that can measure the desired *similarity* between  $\mathbf{x}$  and the decoded  $\hat{\mathbf{x}}$ , either between sample pairs  $\mathbf{x}_i$  and  $\hat{\mathbf{x}}_i$ <sup>6</sup> or between the two distributions  $\mathcal{D}_{\mathbf{x}}$  and  $\mathcal{D}_{\hat{\mathbf{x}}}$ .<sup>7</sup>

Representations learned through this framework can be arguably rich enough to regenerate the data to a certain extent. But depending on the choice of the regularizing heuristics on  $\mathbf{z}$  and similarity metrics on  $\mathbf{x}$  (or  $\mathcal{D}_{\mathbf{x}}$ ), the objective is typically task-dependent and often grossly approximated [RVM<sup>+</sup>11, GPAM<sup>+</sup>14]. When the data contain complicated *multi-modal* structures, naive heuristics or inaccurate metrics may fail to capture all internal subclass structures<sup>8</sup> or to explicitly discriminate among them for classification or clustering purposes. To address this, we propose a principled measure (on  $\mathbf{z}$ ) to learn representations that promotes multi-class discriminative property from data of mixed structures, which works in both supervised and unsupervised settings.

**This work: Learning diverse and discriminative representations.** Whether the given data  $\mathbf{X}$  of a mixed distribution  $\mathcal{D}$  can be effectively classified depends on how separable (or discriminative)

<sup>2</sup>Mutual information is defined to be  $I(\mathbf{z}, \mathbf{y}) \doteq H(\mathbf{z}) - H(\mathbf{z} | \mathbf{y})$  where  $H(\mathbf{z})$  is the entropy of  $\mathbf{z}$  [CT06].

<sup>3</sup>given one can overcome some caveats associated with this framework [KTVK18] and practical difficulties such as how to accurately evaluate mutual information with finitely samples of degenerate distributions.

<sup>4</sup>in case the labels can be corrupted or the learned features be tackled.

<sup>5</sup>hence the auto-encoding [BH89, Kra91] can be viewed as a nonlinear extension to the classical PCA [Jol02].

<sup>6</sup>for tasks such as denoising, in which the metric can be chosen to the  $\ell^p$ -norm between samples of  $\mathbf{x}$  and  $\hat{\mathbf{x}}$ :  $\min_{\theta, \eta} \mathbb{E}[\|\mathbf{x} - \hat{\mathbf{x}}\|_p]$ , where typically  $p = 1$  or  $2$ , for tasks such as image denoising.

<sup>7</sup>the distance between distributions of  $\mathbf{x}$  and  $\hat{\mathbf{x}}$ , say the KL divergence  $\text{KL}(\mathcal{D}_{\mathbf{x}} || \mathcal{D}_{\hat{\mathbf{x}}})$ , is very difficult to evaluate when the data distributions are discrete and degenerate. In practice, it can only be approximated with the help of an additional discriminative network, known as GAN [GPAM<sup>+</sup>14, ACB17].

<sup>8</sup>One consequence of this is the phenomenon of *mode collapsing* in learning generative models for data that have mixed multi-modal structures; see [LPZM20] and references therein.

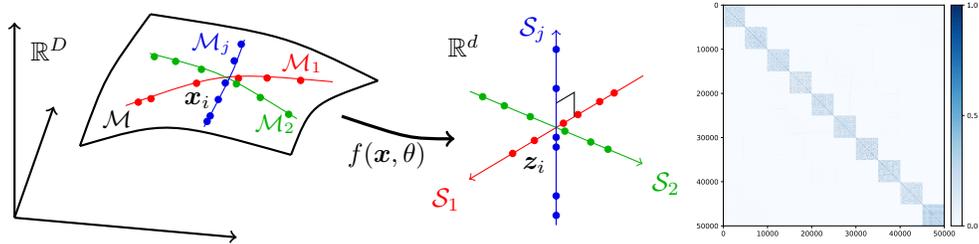


Figure 1: **Left and Middle:** The distribution  $\mathcal{D}$  of high-dim data  $\mathbf{x} \in \mathbb{R}^D$  is supported on a manifold  $\mathcal{M}$  and its classes on low-dim submanifolds  $\mathcal{M}_j$ , we learn a map  $f(\mathbf{x}, \theta)$  such that  $\mathbf{z}_i = f(\mathbf{x}_i, \theta)$  are on a union of maximally uncorrelated subspaces  $\{\mathcal{S}_j\}$ . **Right:** Cosine similarity between learned features by our method for the CIFAR10 training dataset. Each class has 5,000 samples and their features span a subspace of over 10 dimensions (see Figure 3(c)).

the component distributions  $\mathcal{D}_j$  are (or can be made). One popular working assumption is that the distribution of each class has relatively *low-dimensional* intrinsic structures.<sup>9</sup> Hence we may assume the distribution  $\mathcal{D}_j$  of each class has a support on a low-dimensional submanifold, say  $\mathcal{M}_j$  with dimension  $d_j \ll D$ , and the distribution  $\mathcal{D}$  of  $\mathbf{x}$  is supported on the mixture of those submanifolds,  $\mathcal{M} = \cup_{j=1}^k \mathcal{M}_j$ , in the high-dimensional ambient space  $\mathbb{R}^D$ , as illustrated in Figure 1 left.

With the manifold assumption in mind, we want to learn a mapping  $\mathbf{z} = f(\mathbf{x}, \theta)$  that maps each of the submanifolds  $\mathcal{M}_j \subset \mathbb{R}^D$  to a *linear* subspace  $\mathcal{S}_j \subset \mathbb{R}^d$  (see Figure 1 middle). To do so, we require our learned representation to have the following properties:

1. *Between-Class Discriminative:* Features of samples from different classes/clusters should be highly *uncorrelated* and belong to different low-dimensional linear subspaces.
2. *Within-Class Compressible:* Features of samples from the same class/cluster should be relatively *correlated* in a sense that they belong to a low-dimensional linear subspace.
3. *Maximally Diverse Representation:* Dimension (or variance) of features for each class/cluster should be *as large as possible* as long as they stay uncorrelated from the other classes.

Notice that, although the intrinsic structures of each class/cluster may be low-dimensional, they are by no means simply linear in their original representation  $\mathbf{x}$ . Here the subspaces  $\{\mathcal{S}_j\}$  can be viewed as nonlinear *generalized principal components* for  $\mathbf{x}$  [VMS16]. Furthermore, for many clustering or classification tasks (such as object recognition), we consider two samples as *equivalent* if they differ by certain class of domain deformations or augmentations  $\mathcal{T} = \{\tau\}$ . Hence, we are only interested in low-dimensional structures that are *invariant* to such deformations,<sup>10</sup> which are known to have sophisticated geometric and topological structures [WDCB05] and can be difficult to learn in a principled manner even with CNNs [CW16, CGW19]. There are previous attempts to directly enforce subspace structures on features learned by a deep network for supervised [LQMS18] or unsupervised learning [JZL<sup>+</sup>17, ZJH<sup>+</sup>18, PFX<sup>+</sup>17, ZHF18, ZJH<sup>+</sup>19, ZLY<sup>+</sup>19, LQMS18]. However, the *self-expressive* property of subspaces exploited by [JZL<sup>+</sup>17] does not enforce all the desired properties listed above; [LQMS18] uses a nuclear norm based geometric loss to enforce orthogonality between classes, but does not promote diversity in the learned representations, as we will soon see. Figure 1 right illustrates a representation learned by our method on the CIFAR10 dataset. More details can be found in the experimental Section 3.

## 2 Technical Approach and Method

### 2.1 Measure of Compactness for a Representation

Although the above properties are all highly desirable for the latent representation  $\mathbf{z}$ , they are by no means easy to obtain: Are these properties compatible so that we can expect to achieve them all at

<sup>9</sup>There are many reasons why this assumption is plausible: 1. high dimensional data are highly redundant; 2. data that belong to the same class should be similar and correlated to each other; 3. typically we only care about equivalent structures of  $\mathbf{x}$  that are invariant to certain classes of deformation and augmentations.

<sup>10</sup>So  $\mathbf{x} \in \mathcal{M}$  iff  $\tau(\mathbf{x}) \in \mathcal{M}$  for all  $\tau \in \mathcal{T}$ .

once? If so, is there a *simple but principled* objective that can measure the goodness of the resulting representations in terms of all these properties? The key to these questions is to find a principled “measure of compactness” for the distribution of a random variable  $\mathbf{z}$  or from its finite samples  $\mathbf{Z}$ . Such a measure should directly and accurately characterize intrinsic geometric or statistical properties of the distribution, in terms of its intrinsic dimension or volume. Unlike cross-entropy (1) or information bottleneck (2), such a measure should not depend explicitly on class labels so that it can work in all supervised, self-supervised, semi-supervised, and unsupervised settings.

**Low-dimensional degenerate distributions.** In information theory [CT06], the notion of entropy  $H(\mathbf{z})$  is designed to be such a measure.<sup>11</sup> However, entropy is not well-defined for continuous random variables with degenerate distributions.<sup>12</sup> This is unfortunately the case here. To alleviate this difficulty, another related concept in information theory, more specifically in lossy data compression, that measures the “compactness” of a random distribution is the so-called *rate distortion* [CT06]: Given a random variable  $\mathbf{z}$  and a prescribed precision  $\epsilon > 0$ , the rate distortion  $R(\mathbf{z}, \epsilon)$  is the minimal number of binary bits needed to encode  $\mathbf{z}$  such that the expected decoding error<sup>13</sup> is less than  $\epsilon$ . Although this framework has been successful in explaining feature selection in deep networks [MWHK19], the rate distortion of a random variable is difficult, if not impossible to compute, except for simple distributions such as discrete and Gaussian.

**Nonasymptotic rate distortion for finite samples.** When evaluating the lossy coding rate  $R$ , one practical difficulty is that we normally do not know the distribution of  $\mathbf{z}$ . Instead, we have a finite number of samples as learned representations where  $\mathbf{z}_i = f(\mathbf{x}_i, \theta) \in \mathbb{R}^d, i = 1, \dots, m$ , for the given data samples  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$ . Fortunately, [MDHW07] provides a precise estimate on the number of binary bits needed to encode finite samples from a subspace-like distribution. In order to encode the learned representation  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_m]$  up to a precision  $\epsilon$ , the total number of bits needed is given by the following expression<sup>14</sup>:  $\mathcal{L}(\mathbf{Z}, \epsilon) \doteq \left(\frac{m+d}{2}\right) \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}\mathbf{Z}^\top\right)$ . Therefore, the compactness of learned features *as a whole* can be measured in terms of the average coding length per sample (as the sample size  $m$  is large), a.k.a. the *coding rate* subject to the distortion  $\epsilon$ :

$$R(\mathbf{Z}, \epsilon) \doteq \frac{1}{2} \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}\mathbf{Z}^\top\right). \quad (4)$$

**Rate distortion of data with a mixed distribution.** In general, the features  $\mathbf{Z}$  of multi-class data may belong to multiple low-dimensional subspaces. To evaluate the rate distortion of such mixed data *more accurately*, we may partition the data  $\mathbf{Z}$  into multiple subsets:  $\mathbf{Z} = \mathbf{Z}_1 \cup \dots \cup \mathbf{Z}_k$ , with each in one low-dim subspace. So the above coding rate (4) is accurate for each subset. For convenience, let  $\mathbf{\Pi} = \{\mathbf{\Pi}_j \in \mathbb{R}^{m \times m}\}_{j=1}^k$  be a set of diagonal matrices whose diagonal entries encode the membership of the  $m$  samples in the  $k$  classes.<sup>15</sup> Then, according to [MDHW07], with respect to this partition, the average number of bits per sample (the coding rate) is

$$R^c(\mathbf{Z}, \epsilon | \mathbf{\Pi}) \doteq \sum_{j=1}^k \frac{\text{tr}(\mathbf{\Pi}_j)}{2m} \log \det \left(\mathbf{I} + \frac{d}{\text{tr}(\mathbf{\Pi}_j)\epsilon^2} \mathbf{Z}\mathbf{\Pi}_j\mathbf{Z}^\top\right). \quad (5)$$

Notice that when  $\mathbf{Z}$  is given,  $R^c(\mathbf{Z}, \epsilon | \mathbf{\Pi})$  is a concave function of  $\mathbf{\Pi}$ . The function  $\log \det(\cdot)$  in the above expressions has been long known as an effective heuristic for rank minimization problems, with guaranteed convergence to local minimum [FHB03]. As it nicely characterizes the rate distortion of Gaussian or subspace-like distributions,  $\log \det(\cdot)$  can be very effective in clustering or classification of mixed data [MDHW07, WTL+08, KPCC15]. We will soon reveal more desired properties of this function.

## 2.2 Principle of Maximal Coding Rate Reduction

On one hand, for learned features to be discriminative, features of different classes/clusters are preferred to be *maximally incoherent* to each other. Hence they together should span a space of the

<sup>11</sup> given the probability density  $p(\mathbf{z})$  of a random variable,  $H(\mathbf{z}) \doteq - \int p(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z}$ .

<sup>12</sup> The same difficulty resides with evaluating mutual information  $I(\mathbf{x}, \mathbf{z})$  for degenerate distributions.

<sup>13</sup> Say in terms of the  $\ell^2$ -norm, we have  $\mathbb{E}[\|\mathbf{z} - \hat{\mathbf{z}}\|_2] \leq \epsilon$  for the decoded  $\hat{\mathbf{z}}$ .

<sup>14</sup> This formula can be derived either by packing  $\epsilon$ -balls into the space spanned by  $\mathbf{Z}$  or by computing the number of bits needed to quantize the SVD of  $\mathbf{Z}$  subject to the precision, see [MDHW07] for proofs.

<sup>15</sup> That is, the diagonal entry  $\mathbf{\Pi}_j(i, i)$  of  $\mathbf{\Pi}_j$  indicates the probability of sample  $i$  belonging to subset  $j$ . Therefore  $\mathbf{\Pi}$  lies in a simplex:  $\Omega \doteq \{\mathbf{\Pi} | \mathbf{\Pi}_j \geq 0, \mathbf{\Pi}_1 + \dots + \mathbf{\Pi}_k = \mathbf{I}\}$ .

largest possible volume (or dimension) and the coding rate of the whole set  $\mathbf{Z}$  should be as large as possible. On the other hand, learned features of the same class/cluster should be highly correlated and coherent. Hence, each class/cluster should only span a space (or subspace) of a very small volume and the coding rate should be as small as possible. Therefore, a good representation  $\mathbf{Z}$  of  $\mathbf{X}$  is one such that, given a partition  $\mathbf{\Pi}$  of  $\mathbf{Z}$ , achieves a large difference between the coding rate for the whole and that for all the subsets:

$$\Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) \doteq R(\mathbf{Z}, \epsilon) - R^c(\mathbf{Z}, \epsilon | \mathbf{\Pi}). \quad (6)$$

If we choose our feature mapping  $\mathbf{z} = f(\mathbf{x}, \theta)$  to be a deep neural network, the overall process of the feature representation and the resulting rate reduction w.r.t. certain partition  $\mathbf{\Pi}$  can be illustrated by the following diagram:

$$\mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) \xrightarrow{\mathbf{\Pi}, \epsilon} \Delta R(\mathbf{Z}(\theta), \mathbf{\Pi}, \epsilon). \quad (7)$$

Note that  $\Delta R$  is *monotonic* in the scale of the features  $\mathbf{Z}$ . So to make the amount of reduction comparable between different representations,<sup>16</sup> we need to *normalize the scale* of the learned features, either by imposing the Frobenius norm of each class  $\mathbf{Z}_j$  to scale with the number of features in  $\mathbf{Z}_j \in \mathbb{R}^{d \times m_j}$ :  $\|\mathbf{Z}_j\|_F^2 = m_j$  or by normalizing each feature to be on the unit sphere:  $\mathbf{z}_i \in \mathbb{S}^{d-1}$ . This formulation offers a natural justification for the need of “batch normalization” in the practice of training deep neural networks [IS15]. An alternative, arguably simpler, way to normalize the scale of learned representations is to ensure that the mapping of each layer of the network is approximately *isometric* [QYW+20].

Once the representations are comparable, our goal becomes to learn a set of features  $\mathbf{Z}(\theta) = f(\mathbf{X}, \theta)$  and their partition  $\mathbf{\Pi}$  (if not given in advance) such that they maximize the reduction between the coding rate of all features and that of the sum of features w.r.t. their classes:

$$\max_{\theta, \mathbf{\Pi}} \Delta R(\mathbf{Z}(\theta), \mathbf{\Pi}, \epsilon) = R(\mathbf{Z}(\theta), \epsilon) - R^c(\mathbf{Z}(\theta), \epsilon | \mathbf{\Pi}), \quad \text{s.t. } \|\mathbf{Z}_j(\theta)\|_F^2 = m_j, \mathbf{\Pi} \in \Omega. \quad (8)$$

We refer to this as the principle of *maximal coding rate reduction* (MCR<sup>2</sup>), an embodiment of Aristotle’s famous quote: “*the whole is greater than the sum of the parts.*” Note that for the clustering purpose alone, one may only care about the sign of  $\Delta R$  for deciding whether to partition the data or not, which leads to the greedy algorithm in [MDHW07].<sup>17</sup> Here to seek or learn the best representation, we further desire the whole is *maximally* greater than its parts.

**Relationship to information gain.** The maximal coding rate reduction can be viewed as a generalization to *Information Gain* (IG), which aims to maximize the reduction of entropy of a random variable, say  $\mathbf{z}$ , with respect to an observed attribute, say  $\pi$ :  $\max_{\pi} \text{IG}(\mathbf{z}, \pi) \doteq H(\mathbf{z}) - H(\mathbf{z} | \pi)$ , i.e., the *mutual information* between  $\mathbf{z}$  and  $\pi$  [CT06]. Maximal information gain has been widely used in areas such as decision trees [Qui86]. However, MCR<sup>2</sup> is used differently in several ways: 1) One typical setting of MCR<sup>2</sup> is when the data class labels are given, i.e.  $\mathbf{\Pi}$  is known, MCR<sup>2</sup> focuses on learning representations  $\mathbf{z}(\theta)$  rather than fitting labels. 2) In traditional settings of IG, the number of attributes in  $\mathbf{z}$  cannot be so large and their values are discrete (typically binary). Here the “attributes”  $\mathbf{\Pi}$  represent the probability of a multi-class partition for all samples and their values can even be continuous. 3) As mentioned before, entropy  $H(\mathbf{z})$  or mutual information  $I(\mathbf{z}, \pi)$  [HFLM+18] is not well-defined for degenerate continuous distributions whereas the rate distortion  $R(\mathbf{z}, \epsilon)$  is and can be accurately and efficiently computed for (mixed) subspaces, at least.

### 2.3 Properties of the Rate Reduction Function

In theory, the MCR<sup>2</sup> principle (8) benefits from great generalizability and can be applied to representations  $\mathbf{Z}$  of *any* distributions with *any* attributes  $\mathbf{\Pi}$  as long as the rates  $R$  and  $R^c$  for the distributions can be accurately and efficiently evaluated. The optimal representation  $\mathbf{Z}^*$  and partition  $\mathbf{\Pi}^*$  should have some interesting geometric and statistical properties. We here reveal nice properties of the optimal representation with the special case of subspaces, which have many important use cases in

<sup>16</sup>Here different representations can be either representations associated with different network parameters or representations learned after different layers of the same deep network.

<sup>17</sup>Strictly speaking, in the context of clustering *finite* samples, one needs to use the more precise measure of the coding length mentioned earlier, see [MDHW07] for more details.

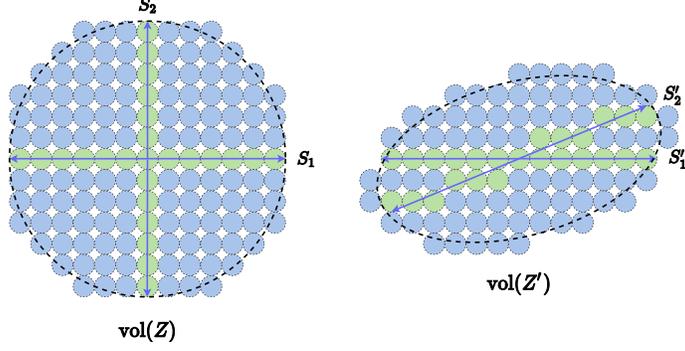


Figure 2: Comparison of two learned representations  $\mathbf{Z}$  and  $\mathbf{Z}'$  via reduced rates:  $R$  is the number of  $\epsilon$ -balls packed in the joint distribution and  $R^c$  is the sum of the numbers for all the subspaces (the green balls).  $\Delta R$  is their difference (the number of blue balls). The MCR<sup>2</sup> principle prefers  $\mathbf{Z}$  (the left one).

machine learning. When the desired representation for  $\mathbf{Z}$  is multiple subspaces, the rates  $R$  and  $R^c$  in (8) are given by (4) and (5), respectively. At the maximal rate reduction, MCR<sup>2</sup> achieves its optimal representations, denoted as  $\mathbf{Z}^* = \mathbf{Z}_1^* \cup \dots \cup \mathbf{Z}_k^* \subset \mathbb{R}^d$  with  $\text{rank}(\mathbf{Z}_j^*) \leq d_j$ . One can show that  $\mathbf{Z}^*$  has the following desired properties (see Appendix A for a formal statement and detailed proofs).

**Theorem 2.1** (Informal Statement). *Suppose  $\mathbf{Z}^* = \mathbf{Z}_1^* \cup \dots \cup \mathbf{Z}_k^*$  is the optimal solution that maximizes the rate reduction (8). We have:*

- **Between-class Discriminative:** *As long as the ambient space is adequately large ( $d \geq \sum_{j=1}^k d_j$ ), the subspaces are all orthogonal to each other, i.e.  $(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* = \mathbf{0}$  for  $i \neq j$ .*
- **Maximally Diverse Representation:** *As long as the coding precision is adequately high, i.e.,  $\epsilon^4 < \min_j \left\{ \frac{m_j d_j^2}{m d_j^2} \right\}$ , each subspace achieves its maximal dimension, i.e.  $\text{rank}(\mathbf{Z}_j^*) = d_j$ . In addition, the largest  $d_j - 1$  singular values of  $\mathbf{Z}_j^*$  are equal.*

In other words, in the case of subspaces, the MCR<sup>2</sup> principle promotes embedding of data into multiple independent subspaces, with features distributed *isotropically* in each subspace (except for possibly one dimension). In addition, among all such discriminative representations, it prefers the one with the highest dimensions in the ambient space. This is substantially different from the objective of information bottleneck (2).

**Comparison to the geometric OLE loss.** To encourage the learned features to be uncorrelated between classes, the work of [LQMS18] has proposed to maximize the difference between the nuclear norm of the whole  $\mathbf{Z}$  and its subsets  $\mathbf{Z}_j$ , called the *orthogonal low-rank embedding* (OLE) loss:  $\max_{\theta} \text{OLE}(\mathbf{Z}(\theta), \mathbf{\Pi}) \doteq \|\mathbf{Z}(\theta)\|_* - \sum_{j=1}^k \|\mathbf{Z}_j(\theta)\|_*$ , added as a regularizer to the cross-entropy loss (1). The nuclear norm  $\|\cdot\|_*$  is a *nonsmooth convex*<sup>18</sup> surrogate for low-rankness, whereas  $\log \det(\cdot)$  is *smooth concave* instead. Unlike the rate reduction  $\Delta R$ , OLE is always *negative* and achieves the maximal value 0 when the subspaces are orthogonal, regardless of their dimensions. So in contrast to  $\Delta R$ , this loss serves as a geometric heuristic and does not promote diverse representations. In fact, OLE typically promotes learning one-dim representations per class, whereas MCR<sup>2</sup> encourages learning subspaces with maximal dimensions (Figure 7 of [LQMS18] versus our Figure 6).

**Relation to contrastive learning.** If samples are *evenly* drawn from  $k$  classes, a randomly chosen pair  $(x_i, x_j)$  is of high probability belonging to difference classes if  $k$  is large.<sup>19</sup> We may view the learned features of two samples together with their their augmentations  $\mathbf{Z}_i$  and  $\mathbf{Z}_j$  as two classes. Then the rate reduction  $\Delta R_{ij} = R(\mathbf{Z}_i \cup \mathbf{Z}_j, \epsilon) - \frac{1}{2}(R(\mathbf{Z}_i, \epsilon) + R(\mathbf{Z}_j, \epsilon))$  gives a “distance” measure for how far the two sample sets are. We may try to further “expand” pairs that likely belong to different classes. From Theorem 2.1, the (averaged) rate reduction  $\Delta R_{ij}$  is maximized when features from different samples are uncorrelated  $\mathbf{Z}_i^\top \mathbf{Z}_j = \mathbf{0}$  (see Figure 2) and features  $\mathbf{Z}_i$  from the same sample are highly correlated. Hence, when applied to sample pairs, MCR<sup>2</sup> naturally conducts the

<sup>18</sup>Nonsmoothness poses additional difficulties in using this loss to learn features via gradient descent.

<sup>19</sup>For example, when  $k \geq 100$ , a random pair is of probability 99% belonging to different classes.

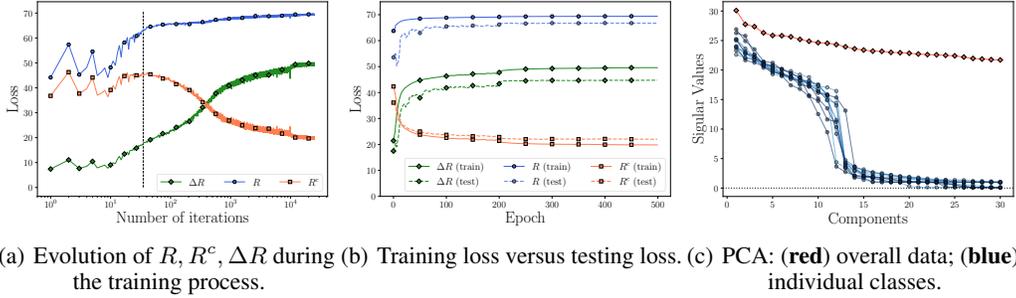


Figure 3: Evolution of the rates of  $MCR^2$  in the training process and principal components of learned features.

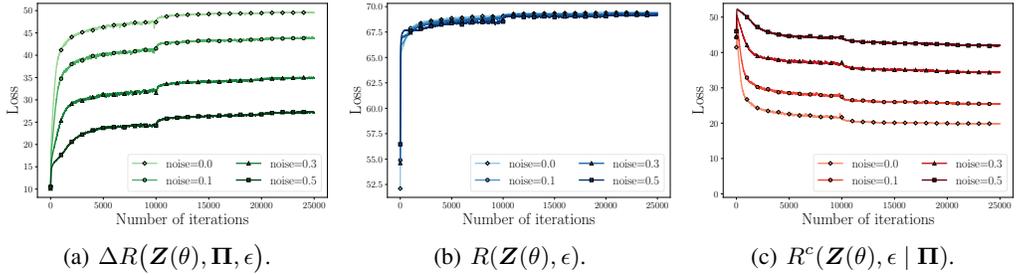


Figure 4: Evolution of rates  $R$ ,  $R^c$ ,  $\Delta R$  of  $MCR^2$  during training with corrupted labels.

so-called *contrastive learning* [HCL06, OLV18, HFW<sup>+</sup>19]. But  $MCR^2$  is *not* limited to expand (or compress) pairs of samples and can uniformly conduct “contrastive learning” for a subset with *any number* of samples as long as we know they likely belong to different (or the same) classes, say by randomly sampling subsets from a large number of classes or with a good clustering method.

### 3 Experiments with Instantiations of $MCR^2$

Our theoretical analysis above shows how the *maximal coding rate reduction* ( $MCR^2$ ) is a principled measure for learning discriminative and diverse representations for mixed data. In this section, we demonstrate experimentally how this principle alone, *without any other heuristics*, is adequate to learning good representations in the supervised, self-supervised, and unsupervised learning settings in a unified fashion. Due to limited space and time, instead of trying to exhaust all its potential and practical implications with extensive engineering, our goal here is only to validate effectiveness of this principle through its most basic usage and fair comparison with existing frameworks. More implementation details and experiments are given in Appendix B. The code can be found in <https://github.com/ryanchankh/mcr2>.

#### 3.1 Supervised Learning of Robust Discriminative Features

**Supervised learning via rate reduction.** When class labels are provided during training, we assign the membership (diagonal) matrix  $\mathbf{\Pi} = \{\mathbf{\Pi}_j\}_{j=1}^k$  as follows: for each sample  $\mathbf{x}_i$  with label  $j$ , set  $\mathbf{\Pi}_j(i, i) = 1$  and  $\mathbf{\Pi}_l(i, i) = 0, \forall l \neq j$ . Then the mapping  $f(\cdot, \theta)$  can be learned by optimizing (8), where  $\mathbf{\Pi}$  remains constant. We apply stochastic gradient descent to optimize  $MCR^2$ , and for each iteration we use mini-batch data  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$  to approximate the  $MCR^2$  loss.

**Evaluation via classification.** As we will see, in the supervised setting, the learned representation has very clear subspace structures. So to evaluate the learned representations, we consider a natural nearest subspace classifier. For each class of learned features  $\mathbf{Z}_j$ , let  $\boldsymbol{\mu}_j \in \mathbb{R}^p$  be its mean and  $\mathbf{U}_j \in \mathbb{R}^{p \times r_j}$  be the first  $r_j$  principal components for  $\mathbf{Z}_j$ , where  $r_j$  is the estimated dimension of class  $j$ . The predicted label of a test data  $\mathbf{x}'$  is given by  $j' = \arg \min_{j \in \{1, \dots, k\}} \|(I - \mathbf{U}_j \mathbf{U}_j^\top)(f(\mathbf{x}', \theta) - \boldsymbol{\mu}_j)\|_2^2$ .

**Experiments on real data.** We consider CIFAR10 dataset [Kri09] and ResNet-18 [HZRS16] for  $f(\cdot, \theta)$ . We replace the last linear layer of ResNet-18 by a two-layer fully connected network with ReLU activation function such that the output dimension is 128. We set the mini-batch size as  $m = 1,000$  and the precision parameter  $\epsilon^2 = 0.5$ . More results can be found in Appendix B.3.2.

Figure 3(a) illustrates how the two rates and their difference (for both training and test data) evolves over epochs of training: After an initial phase,  $R$  gradually increases while  $R^c$  decreases, indicating that features  $\mathbf{Z}$  are expanding as a whole while each class  $\mathbf{Z}_j$  is being compressed. Figure 3(c) shows the distribution of singular values per  $\mathbf{Z}_j$  and Figure 1 (right) shows the angles of features sorted by class. Compared to the geometric loss [LQMS18], our features are *not only orthogonal but also of much higher dimension*. We compare the singular values of representations, both overall data and individual classes, learned by using cross-entropy and MCR<sup>2</sup> in Figure 6 and Figure 7 in Appendix B.3.1. We find that the representations learned by using MCR<sup>2</sup> loss are much more diverse than the ones learned by using cross-entropy loss. In addition, we find that we are able to select diverse images from the same class according to the ‘‘principal’’ components of the learned features (see Figure 8 and Figure 9 in Appendix B.3.1).

**Robustness to corrupted labels.** Because MCR<sup>2</sup> by design encourages richer representations that preserves intrinsic structures from the data  $\mathbf{X}$ , training relies less on class labels than traditional loss such as cross-entropy (CE). To verify this, we train the same network<sup>20</sup> using both CE and MCR<sup>2</sup> with certain ratios of *randomly corrupted* training labels. Figure 4 illustrates the learning process: for different levels of corruption, while the rate for the whole set always converges to the same value, the rates for the classes are inversely proportional to the ratio of corruption, indicating our method only compress samples with valid labels. The classification results are summarized in Table 1. By applying *exact the same* training parameters, MCR<sup>2</sup> is significantly more robust than CE, especially with higher ratio of corrupted labels. This can be an advantage in the settings of self-supervised learning or contrastive learning when the grouping information can be very noisy.

Table 1: Classification results with features learned with labels corrupted at different levels.

	RATIO=0.1	RATIO=0.2	RATIO=0.3	RATIO=0.4	RATIO=0.5
CE TRAINING	90.91%	86.12%	79.15%	72.45%	60.37%
MCR <sup>2</sup> TRAINING	<b>91.16%</b>	<b>89.70%</b>	<b>88.18%</b>	<b>86.66%</b>	<b>84.30%</b>

### 3.2 Self-supervised Learning of Invariant Features

**Learning invariant features via rate reduction.** Motivated by self-supervised learning algorithms [LHB04, KRFL09, OLV18, HFW<sup>+</sup>19, WXYL18], we use the MCR<sup>2</sup> principle to learn representations that are *invariant* to certain class of transformations/augmentations, say  $\mathcal{T}$  with a distribution  $P_{\mathcal{T}}$ . Given a mini-batch of data  $\{\mathbf{x}_j\}_{j=1}^k$ , we augment each sample  $\mathbf{x}_j$  with  $n$  transformations/augmentations  $\{\tau_i(\cdot)\}_{i=1}^n$  randomly drawn from  $P_{\mathcal{T}}$ . We simply label all the augmented samples  $\mathbf{X}_j = [\tau_1(\mathbf{x}_j), \dots, \tau_n(\mathbf{x}_j)]$  of  $\mathbf{x}_j$  as the  $j$ -th class, and  $\mathbf{Z}_j$  the corresponding learned features. Using this self-labeled data, we train our feature mapping  $f(\cdot, \theta)$  the same way as the supervised setting above. For every mini-batch, the total number of samples for training is  $m = kn$ .

**Evaluation via clustering.** To learn invariant features, our formulation itself does *not* require the original samples  $\mathbf{x}_j$  come from a fixed number of classes. For evaluation, we may train on a few classes and observe how the learned features facilitate classification or clustering of the data. A common method to evaluate learned features is to train an additional linear classifier [OLV18, HFW<sup>+</sup>19], with ground truth labels. But for our purpose, because we explicitly verify whether the so-learned invariant features have good subspace structures when the samples come from  $k$  classes, we use an off-the-shelf subspace clustering algorithm EnSC [YLRV16], which is computationally efficient and is provably correct for data with well-structured subspaces. We also use K-Means on the original data  $\mathbf{X}$  as our baseline for comparison. We use normalized mutual information (NMI), clustering accuracy (ACC), and adjusted rand index (ARI) for our evaluation metrics, see Appendix B.3.4 for their detailed definitions.

**Controlling dynamics of expansion and compression.** By directly optimizing the rate reduction  $\Delta R = R - R^c$ , we achieve 0.570 clustering accuracy on CIFAR10 dataset, which is the second best

<sup>20</sup>Both CE and MCR<sup>2</sup> can have better performance by choosing larger models for our mapping.

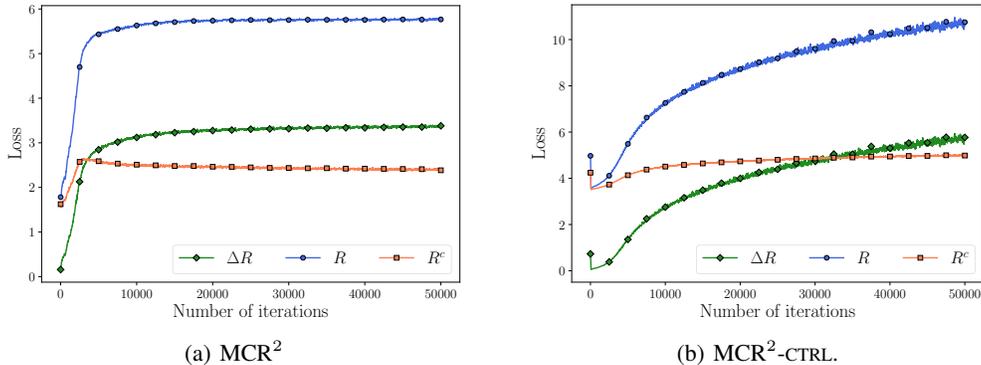


Figure 5: Evolution of the rates of (left)  $MCR^2$  and (right)  $MCR^2\text{-CTRL}$  in the training process in the self-supervised setting on CIFAR10 dataset.

Table 2: Clustering results on CIFAR10, CIFAR100, and STL10 datasets.

DATASET	METRIC	K-MEANS	JULE	RTM	DEC	DAC	DCCM	$MCR^2\text{-CTRL}$
CIFAR10	NMI	0.087	0.192	0.197	0.257	0.395	0.496	<b>0.630</b>
	ACC	0.229	0.272	0.309	0.301	0.521	0.623	<b>0.684</b>
	ARI	0.049	0.138	0.115	0.161	0.305	0.408	<b>0.508</b>
CIFAR100	NMI	0.084	0.103	-	0.136	0.185	0.285	<b>0.362</b>
	ACC	0.130	0.137	-	0.185	0.237	0.327	<b>0.347</b>
	ARI	0.028	0.033	-	0.050	0.087	<b>0.173</b>	0.167
STL10	NMI	0.124	0.182	-	0.276	0.365	0.376	<b>0.446</b>
	ACC	0.192	0.182	-	0.359	0.470	0.482	<b>0.491</b>
	ARI	0.061	0.164	-	0.186	0.256	0.262	<b>0.290</b>

result compared with previous methods. More details can be found in Appendix B.3.3. Empirically, we observe that, without class labels, the overall *coding rate*  $R$  expands quickly and the  $MCR^2$  loss saturates (at a local maximum), see Fig 5(a). Our experience suggests that learning a good representation from unlabeled data might be too ambitious when directly optimizing the original  $\Delta R$ . Nonetheless, from the geometric meaning of  $R$  and  $R^c$ , one can design a different learning strategy by controlling the dynamics of expansion and compression differently during training. For instance, we may re-scale the rate by replacing  $R(\mathbf{Z}, \epsilon)$  with  $\tilde{R}(\mathbf{Z}, \epsilon) \doteq \frac{1}{2\gamma_1} \log \det(\mathbf{I} + \frac{\gamma_2 d}{m\epsilon^2} \mathbf{Z}\mathbf{Z}^\top)$ . With  $\gamma_1 = \gamma_2 = k$ , the learning dynamics change from Fig 5(a) to Fig 5(b): All features are first compressed then gradually expand. We denote the controlled  $MCR^2$  training by  $MCR^2\text{-CTRL}$ .

**Experiments on real data.** Similar to the supervised learning setting, we train *exactly the same* ResNet-18 network on the CIFAR10, CIFAR100, and STL10 [CNL11] datasets. We set the mini-batch size as  $k = 20$ , number of augmentations for each sample as  $n = 50$  and the precision parameter as  $\epsilon^2 = 0.5$ . Table 2 shows the results of the proposed  $MCR^2\text{-CTRL}$  in comparison with methods JULE [YPB16], RTM [NMM19], DEC [XGF16], DAC [CWM<sup>+</sup>17], and DCCM [WLW<sup>+</sup>19] that have achieved the best results on these datasets. Surprisingly, without utilizing any inter-class or inter-sample information and heuristics on the data, the invariant features learned by our method with augmentations alone achieves a better performance over other highly engineered clustering methods. More ablation studies can be found in Appendix B.3.4.

Nevertheless, compared to the representations learned in the supervised setting where the optimal partition  $\mathbf{\Pi}$  in (8) is initialized by correct class information, the representations here learned with self-supervised classes are far from being optimal<sup>21</sup> – they at best correspond to local maxima of the  $MCR^2$  objective (8) when  $\theta$  and  $\mathbf{\Pi}$  are *jointly optimized*. It remains wide open how to design better optimization strategies and dynamics to learn from unlabelled or partially-labelled data better representations (and the associated partitions) close to the global maxima of the  $MCR^2$  objective (8).

<sup>21</sup>We find that the supervised learned representation on CIFAR10 in Section 3.1 can easily achieve a clustering accuracy over 99% on the entire training data.

## 4 Conclusion and Future Work

This work provides rigorous theoretical justifications and clear empirical evidences for why the maximal coding rate reduction (MCR<sup>2</sup>) is a fundamental principle for learning discriminative low-dim representations in almost all learning settings. It unifies and explains existing effective frameworks and heuristics widely practiced in the (deep) learning literature. It remains open *why* MCR<sup>2</sup> is robust to label noises in the supervised setting, *why* self-learned features with MCR<sup>2</sup> alone are effective for clustering, and *how* in future practice instantiations of this principle can be systematically harnessed to further improve clustering or classification tasks.

We believe that MCR<sup>2</sup> gives a principled and practical objective for (deep) learning and can potentially lead to better design operators and architectures of a deep network. A potential direction is to monitor quantitatively the amount of rate reduction  $\Delta R$  gained through every layer of the deep network. By optimizing the rate reduction through the network layers, it is no longer engineered as a “black box.”

On the learning theoretical aspect, although this work has demonstrated only with mixed subspaces, this principle applies to any mixed distributions or structures, for which configurations that achieve maximal rate reduction are of independent theoretical interest. Another interesting note is that the MCR<sup>2</sup> formulation goes beyond the supervised multi-class learning setting often studied through empirical risk minimization (ERM) [DSBDSS15]. It is more related to the expectation maximization (EMX) framework [BDHM<sup>+</sup>17], in which the notion of “compression” plays a crucial role for purely theoretical analysis. We hope this work provides a good connection between machine learning theory and its practice.

## Acknowledgements

Yi would like to thank Professor Yann LeCun of New York University for having a stimulating discussion in his NYU office last November about the search for a proper “energy” function for features to be learned by a deep network [LCH<sup>+</sup>06], during the preparation of a joint proposal. Professor John Wright of Columbia University, who was the leading author of the first two papers on the lossy coding approach to clustering and classification [MDHW07, WTL<sup>+</sup>08], has provided valuable insights and suggestions during germination of this work. We would like to thank Professor Emmanuel Candés of Stanford University for having an online discussion with Yi, during the pandemic, about the rate distortion function for low-dimensional structures. Yi also likes to thank Professor Zhi Ding of UC Davis for discussing the role of rate distortion and lossy coding in communications and information theory and for providing us some pertinent references.

Professor Shankar Sastry of UC Berkeley has always encouraged us to look into fundamental connections between low-dimensional subspaces and deep learning from the perspective of Generalized PCA [VMS16]. Coincidentally, this work was partly motivated to address an inquiry from Professor Ruzena Bajcsy of UC Berkeley earlier this year on how to clarify the role of “latent features” learned in a network in a principled manner. We would also like to thank Professor Jiantao Jiao and Professor Jacob Steinhardt of UC Berkeley for extensive discussions about how to make deep learning robust. During the preparation of this manuscript, Dr. Harry Shum, who collaborated with Yi on lossy coding during his visit to Microsoft Research Asia in 2007 [WTL<sup>+</sup>08], has given excellent suggestions on how to better visualize the learned features, leading to some of the interesting illustrations in this work.

Yaodong would like to thank Zitong Yang and Xili Dai for helpful discussions on the  $\log \det(\cdot)$  function. Ryan would like to thank Yuexiang Zhai for helpful discussions on learning subspace structures. Last but not the least, we are very grateful for Xili Dai and Professor Xiaojun Yuan of UESTC and Professor Hao Chen of UC Davis who have generously provided us their GPU clusters to help us conduct the extensive experiments reported in this paper.

## References

- [ACB17] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pages 214–223, 2017.

- [BDHM<sup>+</sup>17] Shai Ben-David, Pavel Hrubes, Shay Moran, Amir Shpilka, and Amir Yehudayoff. A learning problem that is independent of the set theory ZFC axioms, 2017.
- [BH89] Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- [BV04] Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [CGW19] Taco S Cohen, Mario Geiger, and Maurice Weiler. A general theory of equivariant cnns on homogeneous spaces. In *Advances in Neural Information Processing Systems*, pages 9142–9153, 2019.
- [CNL11] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *International Conference on Artificial Intelligence and Statistics*, pages 215–223, 2011.
- [CT06] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, USA, 2006.
- [CW16] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International Conference on Machine Learning*, pages 2990–2999, 2016.
- [CWM<sup>+</sup>17] Jianlong Chang, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong Pan. Deep adaptive image clustering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5879–5887, 2017.
- [DSBDSS15] Amit Daniely, Sivan Sabato, Shai Ben-David, and Shai Shalev-Shwartz. Multiclass learnability and the ERM principle. *J. Mach. Learn. Res.*, 16(1):2377–2404, January 2015.
- [FHB03] Maryam Fazel, Haitham Hindi, and Stephen P Boyd. Log-det heuristic for matrix rank minimization with applications to hankel and euclidean distance matrices. In *Proceedings of the 2003 American Control Conference, 2003.*, volume 3, pages 2156–2162. IEEE, 2003.
- [GBC16] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [GPAM<sup>+</sup>14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680, 2014.
- [HA85] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of Classification*, 2(1):193–218, 1985.
- [HCL06] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR 2006*, pages 1735–1742, 2006.
- [HFLM<sup>+</sup>18] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018.
- [HFW<sup>+</sup>19] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019.
- [HZRS16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [IS15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [Jol02] Ian T Jolliffe. *Principal Component Analysis*. Springer-Verlag, 2nd edition, 2002.
- [JZL<sup>+</sup>17] Pan Ji, Tong Zhang, Hongdong Li, Mathieu Salzmann, and Ian Reid. Deep subspace clustering networks. In *Advances in Neural Information Processing Systems*, pages 24–33, 2017.

- [KPCC15] Zhao Kang, Chong Peng, Jie Cheng, and Qiang Cheng. Logdet rank minimization with application to subspace clustering. *Computational Intelligence and Neuroscience*, 2015, 2015.
- [Kra91] Mark A Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal*, 37(2):233–243, 1991.
- [KRFL09] Koray Kavukcuoglu, Marc’ Aurelio Ranzato, Rob Fergus, and Yann LeCun. Learning invariant features through topographic filter maps. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1605–1612. IEEE, 2009.
- [Kri09] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- [KTVK18] Artemy Kolchinsky, Brendan D Tracey, and Steven Van Kuyk. Caveats for information bottleneck in deterministic scenarios. *arXiv preprint arXiv:1808.07593*, 2018.
- [LCH<sup>+</sup>06] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and F Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- [LHB04] Yann LeCun, Fu Jie Huang, and Leon Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 2, pages II–104. IEEE, 2004.
- [LPZM20] Ke Li, Shichong Peng, Tianhao Zhang, and Jitendra Malik. Multimodal image synthesis with conditional implicit maximum likelihood estimation. *International Journal of Computer Vision*, 2020.
- [LQMS18] José Lezama, Qiang Qiu, Pablo Musé, and Guillermo Sapiro. OLE: Orthogonal low-rank embedding-a plug and play geometric loss for deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8109–8118, 2018.
- [MDHW07] Yi Ma, Harm Derksen, Wei Hong, and John Wright. Segmentation of multivariate mixed data via lossy data coding and compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(9):1546–1562, 2007.
- [MWHK19] Jan MacDonald, Stephan Wäldchen, Sascha Hauch, and Gitta Kutyniok. A rate-distortion framework for explaining neural network decisions. *CoRR*, abs/1905.11092, 2019.
- [NMM19] Oliver Nina, Jamison Moody, and Clarissa Milligan. A decoder-free approach for unsupervised clustering and manifold learning with random triplet mining. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019.
- [NW06] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, New York, NY, USA, second edition, 2006.
- [OLV18] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [PFX<sup>+</sup>17] Xi Peng, Jiashi Feng, Shijie Xiao, Jiwen Lu, Zhang Yi, and Shuicheng Yan. Deep sparse subspace clustering. *arXiv preprint arXiv:1709.08374*, 2017.
- [PGM<sup>+</sup>19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.
- [Qui86] J. R. Quinlan. Induction of decision trees. *Mach. Learn.*, 1(1):81–106, March 1986.
- [QYW<sup>+</sup>20] Haozhi Qi, Chong You, Xiaolong Wang, Yi Ma, and Jitendra Malik. Deep isometric learning for visual recognition. In *Proceedings of the International Conference on International Conference on Machine Learning*, 2020.
- [RVM<sup>+</sup>11] Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, and Yoshua Bengio. Contractive auto-encoders: Explicit invariance during feature extraction. In *International Conference on Machine Learning*, page 833–840, 2011.

- [SG02] Alexander Strehl and Joydeep Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of Machine Learning Research*, 3(Dec):583–617, 2002.
- [SZ15] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [TZ15] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.
- [VMS16] Rene Vidal, Yi Ma, and S. S. Sastry. *Generalized Principal Component Analysis*. Springer Publishing Company, Incorporated, 1st edition, 2016.
- [WDCB05] Michael B Wakin, David L Donoho, Hyeokho Choi, and Richard G Baraniuk. The multiscale structure of non-differentiable image manifolds. In *Proceedings of SPIE, the International Society for Optical Engineering*, pages 59141B–1, 2005.
- [WLW<sup>+</sup>19] Jianlong Wu, Keyu Long, Fei Wang, Chen Qian, Cheng Li, Zhouchen Lin, and Hongbin Zha. Deep comprehensive correlation mining for image clustering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8150–8159, 2019.
- [WTL<sup>+</sup>08] John Wright, Yangyu Tao, Zhouchen Lin, Yi Ma, and Heung-Yeung Shum. Classification via minimum incremental coding length (micl). In *Advances in Neural Information Processing Systems*, pages 1633–1640, 2008.
- [WXYL18] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018.
- [XGD<sup>+</sup>17] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017.
- [XGF16] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *International Conference on Machine Learning*, pages 478–487, 2016.
- [YLRV16] Chong You, Chun-Guang Li, Daniel P Robinson, and René Vidal. Oracle based active set algorithm for scalable elastic net subspace clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3928–3937, 2016.
- [YPB16] Jianwei Yang, Devi Parikh, and Dhruv Batra. Joint unsupervised learning of deep representations and image clusters. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5147–5156, 2016.
- [ZBH<sup>+</sup>17] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017.
- [ZF14] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833. Springer, 2014.
- [ZHF18] Pan Zhou, Yunqing Hou, and Jiashi Feng. Deep adversarial subspace clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1596–1604, 2018.
- [ZJH<sup>+</sup>18] Tong Zhang, Pan Ji, Mehrtash Harandi, Richard Hartley, and Ian Reid. Scalable deep k-subspace clustering. In *Asian Conference on Computer Vision*, pages 466–481. Springer, 2018.
- [ZJH<sup>+</sup>19] Tong Zhang, Pan Ji, Mehrtash Harandi, Wenbing Huang, and Hongdong Li. Neural collaborative subspace clustering. *arXiv preprint arXiv:1904.10596*, 2019.
- [ZLY<sup>+</sup>19] Junjian Zhang, Chun-Guang Li, Chong You, Xianbiao Qi, Honggang Zhang, Jun Guo, and Zhouchen Lin. Self-supervised convolutional subspace clustering network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5473–5482, 2019.

# Appendices

## A Properties of the Rate Reduction Function

This section is organized as follows. We present background and preliminary results for the  $\log \det(\cdot)$  function and the coding rate function in Section A.1. Then, Section A.2 and A.3 provide technical lemmas for bounding the coding rate and coding rate reduction functions, respectively. Such lemmas are key results for proving our main theoretical results, which are stated informally in Theorem 2.1 and formally in Section A.4. Finally, proof of our main theoretical results is provided in Section A.5.

**Notations** Throughout this section, we use  $\mathbb{S}_{++}^d$ ,  $\mathbb{R}_+$  and  $\mathbb{Z}_{++}$  to denote the set of symmetric positive definite matrices of size  $d \times d$ , nonnegative real numbers and positive integers, respectively.

### A.1 Preliminaries

**Properties of the  $\log \det(\cdot)$  function.**

**Lemma A.1.** *The function  $\log \det(\cdot) : \mathbb{S}_{++}^d \rightarrow \mathbb{R}$  is strictly concave. That is,*

$$\log \det((1 - \alpha)\mathbf{Z}_1 + \alpha\mathbf{Z}_2) \geq (1 - \alpha) \log \det(\mathbf{Z}_1) + \alpha \log \det(\mathbf{Z}_2)$$

for any  $\alpha \in (0, 1)$  and  $\{\mathbf{Z}_1, \mathbf{Z}_2\} \subseteq \mathbb{S}_{++}^d$ , with equality holds if and only if  $\mathbf{Z}_1 = \mathbf{Z}_2$ .

*Proof.* Consider an arbitrary line given by  $\mathbf{Z} = \mathbf{Z}_0 + t\Delta\mathbf{Z}$  where  $\mathbf{Z}_0$  and  $\Delta\mathbf{Z} \neq \mathbf{0}$  are symmetric matrices of size  $d \times d$ . Let  $f(t) \doteq \log \det(\mathbf{Z}_0 + t\Delta\mathbf{Z})$  be a function defined on an interval of values of  $t$  for which  $\mathbf{Z}_0 + t\Delta\mathbf{Z} \in \mathbb{S}_{++}^d$ . Following the same argument as in [BV04], we may assume  $\mathbf{Z}_0 \in \mathbb{S}_{++}^d$  and get

$$f(t) = \log \det \mathbf{Z}_0 + \sum_{i=1}^d \log(1 + t\lambda_i),$$

where  $\{\lambda_i\}_{i=1}^d$  are eigenvalues of  $\mathbf{Z}_0^{-\frac{1}{2}} \Delta\mathbf{Z} \mathbf{Z}_0^{-\frac{1}{2}}$ . The second order derivative of  $f(t)$  is given by

$$f''(t) = - \sum_{i=1}^d \frac{\lambda_i^2}{(1 + t\lambda_i)^2} < 0.$$

Therefore,  $f(t)$  is strictly concave along the line  $\mathbf{Z} = \mathbf{Z}_0 + t\Delta\mathbf{Z}$ . By definition, we conclude that  $\log \det(\cdot)$  is strictly concave.  $\square$

**Properties of the coding rate function.** The following properties, also known as the Sylvester's determinant theorem, for the coding rate function are known in the paper [MDHW07].

**Lemma A.2** (Commutative property [MDHW07]). *For any  $\mathbf{Z} \in \mathbb{R}^{d \times m}$  we have*

$$R(\mathbf{Z}, \epsilon) \doteq \frac{1}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right) = \frac{1}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}^\top \mathbf{Z} \right).$$

**Lemma A.3** (Invariant property [MDHW07]). *For any  $\mathbf{Z} \in \mathbb{R}^{d \times m}$  and any orthogonal matrices  $\mathbf{U} \in \mathbb{R}^{d \times d}$  and  $\mathbf{V} \in \mathbb{R}^{m \times m}$  we have*

$$R(\mathbf{Z}, \epsilon) = R(\mathbf{U} \mathbf{Z} \mathbf{V}^\top, \epsilon).$$

### A.2 Lower and Upper Bounds for Coding Rate

The following result provides an upper and a lower bound on the coding rate of  $\mathbf{Z}$  as a function of the coding rate for its components  $\{\mathbf{Z}_j\}_{j=1}^k$ . The lower bound is tight when all the components  $\{\mathbf{Z}_j\}_{j=1}^k$  have the same covariance (assuming that they have zero mean). The upper bound is tight when the components  $\{\mathbf{Z}_j\}_{j=1}^k$  are pair-wise orthogonal.

**Lemma A.4.** For any  $\{\mathbf{Z}_j \in \mathbb{R}^{d \times m_j}\}_{j=1}^k$  and any  $\epsilon > 0$ , let  $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_k] \in \mathbb{R}^{d \times m}$  with  $m = \sum_{j=1}^k m_j$ . We have

$$\begin{aligned} \sum_{j=1}^k \frac{m_j}{2} \log \det \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right) &\leq \frac{m}{2} \log \det \left( \mathbf{I} + \frac{d}{m \epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right) \\ &\leq \sum_{j=1}^k \frac{m}{2} \log \det \left( \mathbf{I} + \frac{d}{m \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right), \end{aligned} \quad (9)$$

where the first equality holds if and only if

$$\frac{\mathbf{Z}_1 \mathbf{Z}_1^\top}{m_1} = \frac{\mathbf{Z}_2 \mathbf{Z}_2^\top}{m_2} = \dots = \frac{\mathbf{Z}_k \mathbf{Z}_k^\top}{m_k},$$

and the second equality holds if and only if  $\mathbf{Z}_{j_1}^\top \mathbf{Z}_{j_2} = \mathbf{0}$  for all  $1 \leq j_1 < j_2 \leq k$ .

*Proof.* By Lemma A.1,  $\log \det(\cdot)$  is strictly concave. Therefore,

$$\log \det \left( \sum_{j=1}^k \alpha_j \mathbf{S}_j \right) \geq \sum_{j=1}^k \alpha_j \log \det(\mathbf{S}_j), \text{ for all } \{\alpha_j > 0\}_{j=1}^k, \sum_{j=1}^k \alpha_j = 1 \text{ and } \{\mathbf{S}_j \in \mathbb{S}_{++}^d\}_{j=1}^k,$$

where equality holds if and only if  $\mathbf{S}_1 = \mathbf{S}_2 = \dots = \mathbf{S}_k$ . Take  $\alpha_j = \frac{m_j}{m}$  and  $\mathbf{S}_j = \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top$ , we get

$$\frac{m}{2} \log \det \left( \mathbf{I} + \frac{d}{m \epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right) \geq \sum_{j=1}^k \frac{m_j}{2} \log \det \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right),$$

with equality holds if and only if  $\frac{\mathbf{Z}_1 \mathbf{Z}_1^\top}{m_1} = \dots = \frac{\mathbf{Z}_k \mathbf{Z}_k^\top}{m_k}$ . This proves the lower bound in (9).

We now prove the upper bound. By the strict concavity of  $\log \det(\cdot)$ , we have

$$\log \det(\mathbf{Q}) \leq \log \det(\mathbf{S}) + \langle \nabla \log \det(\mathbf{S}), \mathbf{Q} - \mathbf{S} \rangle, \text{ for all } \{\mathbf{Q}, \mathbf{S}\} \subseteq \mathbb{S}_{++}^m,$$

where equality holds if and only if  $\mathbf{Q} = \mathbf{S}$ . Plugging in  $\nabla \log \det(\mathbf{S}) = \mathbf{S}^{-1}$  (see e.g., [BV04]) and  $\mathbf{S}^{-1} = (\mathbf{S}^{-1})^\top$  gives

$$\log \det(\mathbf{Q}) \leq \log \det(\mathbf{S}) + \text{tr}(\mathbf{S}^{-1} \mathbf{Q}) - m. \quad (10)$$

We now take

$$\begin{aligned} \mathbf{Q} &= \mathbf{I} + \frac{d}{m \epsilon^2} \mathbf{Z}^\top \mathbf{Z} = \mathbf{I} + \frac{d}{m \epsilon^2} \begin{bmatrix} \mathbf{Z}_1^\top \mathbf{Z}_1 & \mathbf{Z}_1^\top \mathbf{Z}_2 & \dots & \mathbf{Z}_1^\top \mathbf{Z}_k \\ \mathbf{Z}_2^\top \mathbf{Z}_1 & \mathbf{Z}_2^\top \mathbf{Z}_2 & \dots & \mathbf{Z}_2^\top \mathbf{Z}_k \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{Z}_k^\top \mathbf{Z}_1 & \mathbf{Z}_k^\top \mathbf{Z}_2 & \dots & \mathbf{Z}_k^\top \mathbf{Z}_k \end{bmatrix}, \text{ and} \\ \mathbf{S} &= \mathbf{I} + \frac{d}{m \epsilon^2} \begin{bmatrix} \mathbf{Z}_1^\top \mathbf{Z}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2^\top \mathbf{Z}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{Z}_k^\top \mathbf{Z}_k \end{bmatrix}. \end{aligned} \quad (11)$$

From the property of determinant for block diagonal matrix, we have

$$\log \det(\mathbf{S}) = \sum_{j=1}^k \log \det \left( \mathbf{I} + \frac{d}{m \epsilon^2} \mathbf{Z}_j^\top \mathbf{Z}_j \right). \quad (12)$$

Also, note that

$$\begin{aligned}
& \text{tr}(\mathbf{S}^{-1}\mathbf{Q}) \\
&= \text{tr} \begin{bmatrix} (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_1^\top \mathbf{Z}_1)^{-1} (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_1^\top \mathbf{Z}_1) & \cdots & (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_1^\top \mathbf{Z}_1)^{-1} (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_1^\top \mathbf{Z}_k) \\ \vdots & \ddots & \vdots \\ (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_k^\top \mathbf{Z}_k)^{-1} (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_k^\top \mathbf{Z}_1) & \cdots & (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_k^\top \mathbf{Z}_k)^{-1} (\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_k^\top \mathbf{Z}_k) \end{bmatrix} \\
&= \text{tr} \begin{bmatrix} \mathbf{I} & \cdots & * \\ \vdots & \ddots & \vdots \\ * & \cdots & \mathbf{I} \end{bmatrix} = m, \tag{13}
\end{aligned}$$

where “\*” denotes nonzero quantities that are irrelevant for the purpose of computing the trace. Plugging (12) and (13) back in (10) gives

$$\frac{m}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}^\top \mathbf{Z} \right) \leq \sum_{j=1}^k \frac{m}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^\top \mathbf{Z}_j \right),$$

where the equality holds if and only if  $\mathbf{Q} = \mathbf{S}$ , which by the formulation in (11), holds if and only if  $\mathbf{Z}_{j_1}^\top \mathbf{Z}_{j_2} = \mathbf{0}$  for all  $1 \leq j_1 < j_2 \leq k$ . Further using the result in Lemma A.2 gives

$$\frac{m}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right) \leq \sum_{j=1}^k \frac{m}{2} \log \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right),$$

which produces the upper bound in (9).  $\square$

### A.3 An Upper Bound on Coding Rate Reduction

We may now provide an upper bound on the coding rate reduction  $\Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon)$  (defined in (8)) in terms of its individual components  $\{\mathbf{Z}_j\}_{j=1}^k$ .

**Lemma A.5.** *For any  $\mathbf{Z} \in \mathbb{R}^{d \times m}$ ,  $\mathbf{\Pi} \in \Omega$  and  $\epsilon > 0$ , let  $\mathbf{Z}_j \in \mathbb{R}^{d \times m_j}$  be  $\mathbf{Z} \mathbf{\Pi}_j$  with zero columns removed. We have*

$$\Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) \leq \sum_{j=1}^k \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)} \right), \tag{14}$$

with equality holds if and only if  $\mathbf{Z}_{j_1}^\top \mathbf{Z}_{j_2} = \mathbf{0}$  for all  $1 \leq j_1 < j_2 \leq k$ .

*Proof.* From (4), (5) and (6), we have

$$\begin{aligned}
& \Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) \\
&= R(\mathbf{Z}, \epsilon) - R^c(\mathbf{Z}, \epsilon | \mathbf{\Pi}) \\
&= \frac{1}{2} \log \left( \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right) \right) - \sum_{j=1}^k \left\{ \frac{\text{tr}(\mathbf{\Pi}_j)}{2m} \log \left( \det \left( \mathbf{I} + d \frac{\mathbf{Z} \mathbf{\Pi}_j \mathbf{Z}^\top}{\text{tr}(\mathbf{\Pi}_j) \epsilon^2} \right) \right) \right\} \\
&= \frac{1}{2} \log \left( \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right) \right) - \sum_{j=1}^k \left\{ \frac{m_j}{2m} \log \left( \det \left( \mathbf{I} + d \frac{\mathbf{Z}_j \mathbf{Z}_j^\top}{m_j \epsilon^2} \right) \right) \right\} \\
&\leq \sum_{j=1}^k \frac{1}{2} \log \left( \det \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right) \right) - \sum_{j=1}^k \left\{ \frac{m_j}{2m} \log \left( \det \left( \mathbf{I} + d \frac{\mathbf{Z}_j \mathbf{Z}_j^\top}{m_j \epsilon^2} \right) \right) \right\} \\
&= \sum_{j=1}^k \frac{1}{2m} \log \left( \det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right) \right) - \sum_{j=1}^k \left\{ \frac{1}{2m} \log \left( \det^{m_j} \left( \mathbf{I} + d \frac{\mathbf{Z}_j \mathbf{Z}_j^\top}{m_j \epsilon^2} \right) \right) \right\} \\
&= \sum_{j=1}^k \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)} \right),
\end{aligned}$$

where the inequality follows from the upper bound in Lemma A.4, and that the equality holds if and only if  $\mathbf{Z}_{j_1}^\top \mathbf{Z}_{j_2} = \mathbf{0}$  for all  $1 \leq j_1 < j_2 \leq k$ .  $\square$

#### A.4 Main Results: Properties of Maximal Coding Rate Reduction

We now present our main theoretical results. The following theorem states that for any fixed encoding of the partition  $\mathbf{\Pi}$ , the coding rate reduction is maximized by data  $\mathbf{Z}$  that is maximally discriminative between different classes and is diverse within each of the classes. This result holds provided that the sum of rank for different classes is small relative to the ambient dimension, and that  $\epsilon$  is small.

**Theorem A.6.** *Let  $\mathbf{\Pi} = \{\mathbf{\Pi}_j \in \mathbb{R}^{m \times m}\}_{j=1}^k$  with  $\{\mathbf{\Pi}_j \geq \mathbf{0}\}_{j=1}^k$  and  $\mathbf{\Pi}_1 + \dots + \mathbf{\Pi}_k = \mathbf{I}$  be a given set of diagonal matrices whose diagonal entries encode the membership of the  $m$  samples in the  $k$  classes. Given any  $\epsilon > 0$ ,  $d > 0$  and  $\{d \geq d_j > 0\}_{j=1}^k$ , consider the optimization problem*

$$\begin{aligned} \mathbf{Z}^* \in \arg \max_{\mathbf{Z} \in \mathbb{R}^{d \times m}} \Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) \\ \text{s.t. } \|\mathbf{Z}\mathbf{\Pi}_j\|_F^2 = \text{tr}(\mathbf{\Pi}_j), \text{rank}(\mathbf{Z}\mathbf{\Pi}_j) \leq d_j, \forall j \in \{1, \dots, k\}. \end{aligned} \quad (15)$$

Under the conditions

- (Large ambient dimension)  $d \geq \sum_{j=1}^k d_j$ , and
- (High coding precision)  $\epsilon^4 < \min_{j \in \{1, \dots, k\}} \left\{ \frac{\text{tr}(\mathbf{\Pi}_j) d^2}{m d_j^2} \right\}$ ,

the optimal solution  $\mathbf{Z}^*$  satisfies

- (Between-class discriminative)  $(\mathbf{Z}_{j_1}^*)^\top \mathbf{Z}_{j_2}^* = \mathbf{0}$  for all  $1 \leq j_1 < j_2 \leq k$ , i.e.,  $\mathbf{Z}_{j_1}^*$  and  $\mathbf{Z}_{j_2}^*$  lie in orthogonal subspaces, and
- (Within-class diverse) For each  $j \in \{1, \dots, k\}$ , the rank of  $\mathbf{Z}_j^*$  is equal to  $d_j$  and either all singular values of  $\mathbf{Z}_j^*$  are equal to  $\frac{\text{tr}(\mathbf{\Pi}_j)}{d_j}$ , or the  $d_j - 1$  largest singular values of  $\mathbf{Z}_j^*$  are equal and have value larger than  $\frac{\text{tr}(\mathbf{\Pi}_j)}{d_j}$ ,

where  $\mathbf{Z}_j^* \in \mathbb{R}^{d \times \text{tr}(\mathbf{\Pi}_j)}$  denotes  $\mathbf{Z}^* \mathbf{\Pi}_j$  with zero columns removed.

#### A.5 Proof of Main Results

We start with presenting a lemma that will be used in the proof to Theorem A.6.

**Lemma A.7.** *Given any twice differentiable  $f: \mathbb{R}_+ \rightarrow \mathbb{R}$ , integer  $r \in \mathbb{Z}_{++}$  and  $c \in \mathbb{R}_+$ , consider the optimization problem*

$$\begin{aligned} \max_{\mathbf{x}} \sum_{p=1}^r f(x_p) \\ \text{s.t. } \mathbf{x} = [x_1, \dots, x_r] \in \mathbb{R}_+^r, x_1 \geq x_2 \geq \dots \geq x_r, \text{ and } \sum_{p=1}^r x_p = c. \end{aligned} \quad (16)$$

Let  $\mathbf{x}^*$  be an arbitrary global solution to (16). If the conditions

- $f'(0) < f'(x)$  for all  $x > 0$ ,
- There exists  $x_T > 0$  such that  $f'(x)$  is strictly increasing in  $[0, x_T]$  and strictly decreasing in  $[x_T, \infty)$ ,
- $f''(\frac{c}{r}) < 0$  (equivalently,  $\frac{c}{r} > x_T$ ),

are satisfied, then we have either

- $\mathbf{x}^* = [\frac{c}{r}, \dots, \frac{c}{r}]$ , or
- $\mathbf{x}^* = [x_H, \dots, x_H, x_L]$  for some  $x_H \in (\frac{c}{r}, \frac{c}{r-1})$  and  $x_L > 0$ .

*Proof.* The result holds trivially if  $r = 1$ . Throughout the proof we consider the case where  $r > 1$ .

We consider the optimization problem with the inequality constraint  $x_1 \geq \dots \geq x_r$  in (16) removed:

$$\max_{\mathbf{x}=[x_1, \dots, x_r] \in \mathbb{R}_+^r} \sum_{p=1}^r f(x_p) \quad \text{s.t.} \quad \sum_{p=1}^r x_p = c. \quad (17)$$

We need to show that any global solution  $\mathbf{x}^* = [x_1^*, \dots, x_r^*]$  to (17) is either  $\mathbf{x}^* = [\frac{c}{r}, \dots, \frac{c}{r}]$  or  $\mathbf{x}^* = [x_H, \dots, x_H, x_L] \cdot \mathbf{P}$  for some  $x_H > \frac{c}{r}$ ,  $x_L > 0$  and permutation matrix  $\mathbf{P} \in \mathbb{R}^{r \times r}$ . Let

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \sum_{p=1}^r f(x_p) - \lambda_0 \cdot \left( \sum_{p=1}^r x_p - c \right) - \sum_{p=1}^r \lambda_p x_p$$

be the Lagrangian function for (17) where  $\boldsymbol{\lambda} = [\lambda_0, \lambda_1, \dots, \lambda_r]$  is the Lagrangian multiplier. By the first order optimality conditions (i.e., the Karush–Kuhn–Tucker (KKT) conditions, see, e.g., [NW06, Theorem 12.1]), there exists  $\boldsymbol{\lambda}^* = [\lambda_0^*, \lambda_1^*, \dots, \lambda_r^*]$  such that

$$\sum_{p=1}^r x_p^* = c, \quad (18)$$

$$x_q^* \geq 0, \quad \forall q \in \{1, \dots, r\}, \quad (19)$$

$$\lambda_q^* \geq 0, \quad \forall q \in \{1, \dots, r\}, \quad (20)$$

$$\lambda_q^* \cdot x_q^* = 0, \quad \forall q \in \{1, \dots, r\}, \quad \text{and} \quad (21)$$

$$[f'(x_1^*), \dots, f'(x_r^*)] = [\lambda_0^*, \dots, \lambda_0^*] + [\lambda_1^*, \dots, \lambda_r^*]. \quad (22)$$

By using the KKT conditions, we first show that all entries of  $\mathbf{x}^*$  are strictly positive. To prove by contradiction, suppose that  $\mathbf{x}^*$  has  $r_0$  nonzero entries and  $r - r_0$  zero entries for some  $1 \leq r_0 < r$ . Note that  $r_0 \geq 1$  since an all zero vector  $\mathbf{x}^*$  does not satisfy the equality constraint (18).

Without loss of generality, we may assume that  $x_p^* > 0$  for  $p \leq r_0$  and  $x_p^* = 0$  otherwise. By (21), we have

$$\lambda_1^* = \dots = \lambda_{r_0}^* = 0.$$

Plugging it into (22), we get

$$f'(x_1^*) = \dots = f'(x_{r_0}^*) = \lambda_0^*.$$

From (22) and noting that  $x_{r_0+1} = 0$  we get

$$f'(0) = f'(x_{r_0+1}) = \lambda_0^* + \lambda_{r_0+1}^*.$$

Finally, from (20), we have

$$\lambda_{r_0+1}^* \geq 0.$$

Combining the last three equations above gives  $f'(0) - f'(x_1^*) \geq 0$ , contradicting the assumption that  $f'(0) < f'(x)$  for all  $x > 0$ . This shows that  $r_0 = r$ , i.e., all entries of  $\mathbf{x}^*$  are strictly positive. Using this fact and (21) gives

$$\lambda_p^* = 0 \quad \text{for all } p \in \{1, \dots, r\}.$$

Combining this with (22) gives

$$f'(x_1^*) = \dots = f'(x_r^*) = \lambda_0^*. \quad (23)$$

It follows from the fact that  $f'(x)$  is strictly unimodal that

$$\exists x_H \geq x_L > 0 \quad \text{s.t.} \quad \{x_p^*\}_{p=1}^r \subseteq \{x_L, x_H\}. \quad (24)$$

That is, the set  $\{x_p^*\}_{p=1}^r$  may contain no more than two values. To see why this is true, suppose that there exists three distinct values for  $\{x_p^*\}_{p=1}^r$ . Without loss of generality we may assume that  $0 < x_1^* < x_2^* < x_3^*$ . If  $x_2^* \leq x_T$  (recall  $x_T := \arg \max_{x \geq 0} f'(x)$ ), then by using the fact that  $f'(x)$  is strictly increasing in  $[0, x_T]$ , we must have  $f'(x_1^*) < f'(x_2^*)$  which contradicts (23). A similar contradiction is arrived by considering  $f'(x_2^*)$  and  $f'(x_3^*)$  for the case where  $x_2^* > x_T$ .

There are two possible cases as a consequence of (24). First, if  $x_L = x_H$ , then we have  $x_1^* = \dots = x_r^*$ . By further using (18) we get

$$x_1^* = \dots = x_r^* = \frac{c}{r}.$$

It remains to consider the case where  $x_L < x_H$ . First, by the unimodality of  $f'(x)$ , we must have  $x_L < x_T < x_H$ , therefore

$$f''(x_L) > 0 \text{ and } f''(x_H) < 0. \quad (25)$$

Let  $\ell := |\{p : x_p = x_L\}|$  be the number of entries of  $\mathbf{x}^*$  that are equal to  $x_L$  and  $h := r - \ell$ . We show that it is necessary to have  $\ell = 1$  and  $h = r - 1$ . To prove by contradiction, assume that  $\ell > 1$  and  $h < r - 1$ . Without loss of generality we may assume  $\{x_p^* = x_H\}_{p=1}^h$  and  $\{x_p^* = x_L\}_{p=h+1}^r$ . By (25), we have

$$f''(x_p^*) > 0 \text{ for all } p > h.$$

In particular, by using  $h < r - 1$  we have

$$f''(x_{r-1}^*) > 0 \text{ and } f''(x_r^*) > 0. \quad (26)$$

On the other hand, by using the second order necessary conditions for constraint optimization (see, e.g., [NW06, Theorem 12.5]), the following result holds

$$\begin{aligned} & \mathbf{v}^\top \nabla_{\mathbf{x}\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{v} \leq 0, \text{ for all } \left\{ \mathbf{v} : \left\langle \nabla_{\mathbf{x}} \left( \sum_{p=1}^r x_p^* - c \right), \mathbf{v} \right\rangle = 0 \right\} \\ \iff & \sum_{p=1}^r f''(x_p^*) \cdot v_p^2 \leq 0, \text{ for all } \left\{ \mathbf{v} = [v_1, \dots, v_r] : \sum_{p=1}^r v_p = 0 \right\}. \end{aligned} \quad (27)$$

Take  $\mathbf{v}$  to be such that  $v_1 = \dots = v_{r-2} = 0$  and  $v_{r-1} = -v_r \neq 0$ . Plugging it into (27) gives

$$f''(x_{r-1}^*) + f''(x_r^*) \leq 0,$$

which contradicts (26). Therefore, we may conclude that  $\ell = 1$ . That is,  $\mathbf{x}^*$  is given by

$$\mathbf{x}^* = [x_H, \dots, x_H, x_L], \text{ where } x_H > x_L > 0.$$

By using the condition in (18), we may further show that

$$\begin{aligned} (r-1)x_H + x_L = c & \implies x_H = \frac{c}{r-1} - \frac{c}{x_L} < \frac{x_L}{r-1}, \\ (r-1)x_H + x_L = c & \implies (r-1)x_H + x_H > c \implies x_H > \frac{c}{r}, \end{aligned}$$

which completes our proof.  $\square$

*Proof of Theorem A.6.* Without loss of generality, let  $\mathbf{Z}^* = [\mathbf{Z}_1^*, \dots, \mathbf{Z}_k^*]$  be the optimal solution of problem (15).

To show that  $\mathbf{Z}_j^*, j \in \{1, \dots, k\}$  are pairwise orthogonal, suppose for the purpose of arriving at a contradiction that  $(\mathbf{Z}_{j_1}^*)^\top \mathbf{Z}_{j_2}^* \neq \mathbf{0}$  for some  $1 \leq j_1 < j_2 \leq k$ . By using Lemma A.5, the strict inequality in (14) holds for the optimal solution  $\mathbf{Z}^*$ . That is,

$$\Delta R(\mathbf{Z}^*, \boldsymbol{\Pi}, \epsilon) < \sum_{j=1}^k \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right). \quad (28)$$

On the other hand, since  $\sum_{j=1}^k d_j \leq d$ , there exists  $\{\mathbf{U}'_j \in \mathbb{R}^{d \times d_j}\}_{j=1}^k$  such that the columns of the matrix  $[\mathbf{U}'_1, \dots, \mathbf{U}'_k]$  are orthonormal. Denote  $\mathbf{Z}_j^* = \mathbf{U}'_j \boldsymbol{\Sigma}_j^* (\mathbf{V}_j^*)^\top$  the compact SVD of  $\mathbf{Z}_j^*$ , and let

$$\mathbf{Z}' = [\mathbf{Z}'_1, \dots, \mathbf{Z}'_k], \text{ where } \mathbf{Z}'_j = \mathbf{U}'_j \boldsymbol{\Sigma}_j^* (\mathbf{V}_j^*)^\top.$$

It follows that

$$(\mathbf{Z}'_{j_1})^\top \mathbf{Z}'_{j_2} = \mathbf{V}_{j_1}^* \boldsymbol{\Sigma}_{j_1}^* (\mathbf{U}'_{j_1})^\top \mathbf{U}'_{j_2} \boldsymbol{\Sigma}_{j_2}^* (\mathbf{V}_{j_2}^*)^\top = \mathbf{V}_{j_1}^* \boldsymbol{\Sigma}_{j_1}^* \mathbf{0} \boldsymbol{\Sigma}_{j_2}^* (\mathbf{V}_{j_2}^*)^\top = \mathbf{0} \text{ for all } 1 \leq j_1 < j_2 \leq k.$$

That is, the matrices  $\mathbf{Z}'_1, \dots, \mathbf{Z}'_k$  are pairwise orthogonal. Applying Lemma A.5 for  $\mathbf{Z}'$  gives

$$\begin{aligned} \Delta R(\mathbf{Z}', \boldsymbol{\Pi}, \epsilon) &= \sum_{j=1}^k \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}'_j (\mathbf{Z}'_j)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}'_j (\mathbf{Z}'_j)^\top \right)} \right) \\ &= \sum_{j=1}^k \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right), \end{aligned} \quad (29)$$

where the second equality follows from Lemma A.3. Comparing (28) and (29) gives  $\Delta R(\mathbf{Z}', \mathbf{\Pi}, \epsilon) > \Delta R(\mathbf{Z}^*, \mathbf{\Pi}, \epsilon)$ , which contradicts the optimality of  $\mathbf{Z}^*$ . Therefore, we must have

$$(\mathbf{Z}_{j_1}^*)^\top \mathbf{Z}_{j_2}^* = \mathbf{0} \text{ for all } 1 \leq j_1 < j_2 \leq k.$$

Moreover, from Lemma A.3 we have

$$\Delta R(\mathbf{Z}^*, \mathbf{\Pi}, \epsilon) = \sum_{j=1}^k \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right). \quad (30)$$

We now prove the result concerning the singular values of  $\mathbf{Z}_j^*$ . To start with, we claim that the following result holds:

$$\mathbf{Z}_j^* \in \arg \max_{\mathbf{Z}_j} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)} \right) \text{ s.t. } \|\mathbf{Z}_j\|_F^2 = m_j, \text{ rank}(\mathbf{Z}_j) \leq d_j. \quad (31)$$

To see why (31) holds, suppose that there exists  $\tilde{\mathbf{Z}}_j$  such that  $\|\tilde{\mathbf{Z}}_j\|_F^2 = m_j$ ,  $\text{rank}(\tilde{\mathbf{Z}}_j) \leq d_j$  and

$$\log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \tilde{\mathbf{Z}}_j \tilde{\mathbf{Z}}_j^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \tilde{\mathbf{Z}}_j \tilde{\mathbf{Z}}_j^\top \right)} \right) > \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right). \quad (32)$$

Denote  $\tilde{\mathbf{Z}}_j = \tilde{\mathbf{U}}_j \tilde{\mathbf{\Sigma}}_j \tilde{\mathbf{V}}_j^\top$  the compact SVD of  $\tilde{\mathbf{Z}}_j$  and let

$$\mathbf{Z}' = [\mathbf{Z}_1^*, \dots, \mathbf{Z}_{j-1}^*, \mathbf{Z}'_j, \mathbf{Z}_{j+1}^*, \dots, \mathbf{Z}_k^*], \text{ where } \mathbf{Z}'_j := \mathbf{U}_j^* \tilde{\mathbf{\Sigma}}_j \tilde{\mathbf{V}}_j^\top.$$

Note that  $\|\mathbf{Z}'_j\|_F^2 = m_j$ ,  $\text{rank}(\mathbf{Z}'_j) \leq d_j$  and  $(\mathbf{Z}'_j)^\top \mathbf{Z}'_{j'} = \mathbf{0}$  for all  $j' \neq j$ . It follows that  $\mathbf{Z}'$  is a feasible solution to (15) and that the components of  $\mathbf{Z}'$  are pairwise orthogonal. By using Lemma A.5, Lemma A.3 and (32) we have

$$\begin{aligned} & \Delta R(\mathbf{Z}', \mathbf{\Pi}, \epsilon) \\ &= \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}'_j (\mathbf{Z}'_j)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}'_j (\mathbf{Z}'_j)^\top \right)} \right) + \sum_{j' \neq j} \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}'_{j'} (\mathbf{Z}'_{j'})^\top \right)}{\det^{m_{j'}} \left( \mathbf{I} + \frac{d}{m_{j'} \epsilon^2} \mathbf{Z}'_{j'} (\mathbf{Z}'_{j'})^\top \right)} \right) \\ &= \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \tilde{\mathbf{Z}}_j (\tilde{\mathbf{Z}}_j)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \tilde{\mathbf{Z}}_j (\tilde{\mathbf{Z}}_j)^\top \right)} \right) + \sum_{j' \neq j} \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}'_{j'} (\mathbf{Z}'_{j'})^\top \right)}{\det^{m_{j'}} \left( \mathbf{I} + \frac{d}{m_{j'} \epsilon^2} \mathbf{Z}'_{j'} (\mathbf{Z}'_{j'})^\top \right)} \right) \\ &> \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right) + \sum_{j' \neq j} \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}'_{j'} (\mathbf{Z}'_{j'})^\top \right)}{\det^{m_{j'}} \left( \mathbf{I} + \frac{d}{m_{j'} \epsilon^2} \mathbf{Z}'_{j'} (\mathbf{Z}'_{j'})^\top \right)} \right) \\ &= \sum_{j=1}^k \frac{1}{2m} \log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j^* (\mathbf{Z}_j^*)^\top \right)} \right). \end{aligned}$$

Combining it with (30) shows  $\Delta R(\mathbf{Z}', \mathbf{\Pi}, \epsilon) > \Delta R(\mathbf{Z}^*, \mathbf{\Pi}, \epsilon)$ , contradicting the optimality of  $\mathbf{Z}^*$ . Therefore, the result in (31) holds.

Observe that the optimization problem in (31) depends on  $\mathbf{Z}_j$  only through its singular values. That is, by letting  $\boldsymbol{\sigma}_j := [\sigma_{1,j}, \dots, \sigma_{\min\{m_j, d\}, j}]$  be the singular values of  $\mathbf{Z}_j$ , we have

$$\log \left( \frac{\det^m \left( \mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)}{\det^{m_j} \left( \mathbf{I} + \frac{d}{m_j \epsilon^2} \mathbf{Z}_j \mathbf{Z}_j^\top \right)} \right) = \sum_{p=1}^{\min\{m_j, d\}} \log \left( \frac{\left(1 + \frac{d}{m\epsilon^2} \sigma_{p,j}^2\right)^m}{\left(1 + \frac{d}{m_j \epsilon^2} \sigma_{p,j}^2\right)^{m_j}} \right),$$

also, we have

$$\|\mathbf{Z}_j\|_F^2 = \sum_{p=1}^{\min\{m_j, d\}} \sigma_{p,j}^2 \text{ and } \text{rank}(\mathbf{Z}_j) = \|\boldsymbol{\sigma}_j\|_0.$$

Using these relations, (31) is equivalent to

$$\begin{aligned} & \max_{\sigma_j \in \mathbb{R}_+^{\min\{m_j, d\}}} \sum_{p=1}^{\min\{m_j, d\}} \log \left( \frac{(1 + \frac{d}{m\epsilon^2} \sigma_{p,j}^2)^m}{(1 + \frac{d}{m_j \epsilon^2} \sigma_{p,j}^2)^{m_j}} \right) \\ & \text{s.t. } \sum_{p=1}^{\min\{m_j, d\}} \sigma_{p,j}^2 = m_j, \text{ and } \text{rank}(\mathbf{Z}_j) = \|\sigma_j\|_0 \end{aligned} \quad (33)$$

Let  $\sigma_j^* = [\sigma_{1,j}^*, \dots, \sigma_{\min\{m_j, d\}, j}^*]$  be an optimal solution to (33). Without loss of generality we assume that the entries of  $\sigma_j^*$  are sorted in descending order. It follows that

$$\sigma_{p,j}^* = 0 \text{ for all } p > d_j,$$

and

$$[\sigma_{1,j}^*, \dots, \sigma_{d_j,j}^*] = \arg \max_{\substack{[\sigma_{1,j}, \dots, \sigma_{d_j,j}] \in \mathbb{R}_+^{d_j} \\ \sigma_{1,j} \geq \dots \geq \sigma_{d_j,j}}} \sum_{p=1}^{d_j} \log \left( \frac{(1 + \frac{d}{m\epsilon^2} \sigma_{p,j}^2)^m}{(1 + \frac{d}{m_j \epsilon^2} \sigma_{p,j}^2)^{m_j}} \right) \text{ s.t. } \sum_{p=1}^{d_j} \sigma_{p,j}^2 = m_j. \quad (34)$$

Then we define

$$f(x; d, \epsilon, m_j, m) = \log \left( \frac{(1 + \frac{d}{m\epsilon^2} x)^m}{(1 + \frac{d}{m_j \epsilon^2} x)^{m_j}} \right),$$

and rewrite (34) as

$$\max_{\substack{[x_1, \dots, x_{d_j}] \in \mathbb{R}_+^{d_j} \\ x_1 \geq \dots \geq x_{d_j}}} \sum_{p=1}^{d_j} f(x_p; d, \epsilon, m_j, m) \text{ s.t. } \sum_{p=1}^{d_j} x_p = m_j. \quad (35)$$

We compute the first and second derivative for  $f$  with respect to  $x$ , which are given by

$$\begin{aligned} f'(x; d, \epsilon, m_j, m) &= \frac{d^2 x (m - m_j)}{(dx + m\epsilon^2)(dx + m_j \epsilon^2)}, \\ f''(x; d, \epsilon, m_j, m) &= \frac{d^2 (m - m_j) (m m_j \epsilon^4 - d^2 x^2)}{(dx + m\epsilon^2)^2 (dx + m_j \epsilon^2)^2}. \end{aligned}$$

Note that

- $0 = f'(0) < f'(x)$  for all  $x > 0$ ,
- $f'(x)$  is strictly increasing in  $[0, x_T]$  and strictly decreasing in  $[x_T, \infty)$ , where  $x_T = \epsilon^2 \sqrt{\frac{m}{d} \frac{m_j}{d}}$ , and
- by using the condition  $\epsilon^4 < \frac{m_j}{m} \frac{d^2}{d_j^2}$ , we have  $f''(\frac{m_j}{d_j}) < 0$ .

Therefore, we may apply Lemma A.7 and conclude that the unique optimal solution to (35) is either

- $\mathbf{x}^* = [\frac{m_j}{d_j}, \dots, \frac{m_j}{d_j}]$ , or
- $\mathbf{x}^* = [x_H, \dots, x_H, x_L]$  for some  $x_H \in (\frac{m_j}{d_j}, \frac{m_j}{d_j - 1})$  and  $x_L > 0$ .

Equivalently, we have either

- $[\sigma_{1,j}^*, \dots, \sigma_{d_j,j}^*] = [\sqrt{\frac{m_j}{d_j}}, \dots, \sqrt{\frac{m_j}{d_j}}]$ , or
- $[\sigma_{1,j}^*, \dots, \sigma_{d_j,j}^*] = [\sigma_H, \dots, \sigma_H, \sigma_L]$  for some  $\sigma_H \in (\sqrt{\frac{m_j}{d_j}}, \sqrt{\frac{m_j}{d_j - 1}})$  and  $\sigma_L > 0$ ,

as claimed.  $\square$

## B Additional Simulations and Experiments

### B.1 Simulations - Verifying Diversity Promoting Properties of MCR<sup>2</sup>

As proved in Theorem A.6, the proposed MCR<sup>2</sup> objective promotes within-class diversity. In this section, we use simulated data to verify the diversity promoting property of MCR<sup>2</sup>. As shown in Table 3, we calculate our proposed MCR<sup>2</sup> objective on simulated data. We observe that orthogonal subspaces with *higher* dimension achieve higher MCR<sup>2</sup> value, which is consistent with our theoretical analysis in Theorem A.6.

Table 3: **MCR<sup>2</sup> objective on simulated data.** We evaluate the proposed MCR<sup>2</sup> objective defined in (8), including  $R$ ,  $R^c$ , and  $\Delta R$ , on simulated data. The output dimension  $d$  is set as 512, 256, and 128. We set the batch size as  $m = 1000$  and random assign the label of each sample from 0 to 9, i.e., 10 classes. We generate two types of data: 1) (RANDOM GAUSSIAN) For comparison with data without structures, for each class we generate random vectors sampled from Gaussian distribution (the dimension is set as the output dimension  $d$ ) and normalize each vector to be on the unit sphere. 2) (SUBSPACE) For each class, we generate vectors sampled from its corresponding subspace with dimension  $d_j$  and normalize each vector to be on the unit sphere. We consider the subspaces from different classes are orthogonal/nonorthogonal to each other.

	$R$	$R^c$	$\Delta R$	ORTHOGONAL?	OUTPUT DIMENSION
RANDOM GAUSSIAN	552.70	193.29	360.41	✓	512
SUBSPACE ( $d_j = 50$ )	545.63	108.46	<b>437.17</b>	✓	512
SUBSPACE ( $d_j = 40$ )	487.07	92.71	394.36	✓	512
SUBSPACE ( $d_j = 30$ )	413.08	74.84	338.24	✓	512
SUBSPACE ( $d_j = 20$ )	318.52	54.48	264.04	✓	512
SUBSPACE ( $d_j = 10$ )	195.46	30.97	164.49	✓	512
SUBSPACE ( $d_j = 1$ )	31.18	4.27	26.91	✓	512
<hr/>					
RANDOM GAUSSIAN	292.71	154.13	138.57	✓	256
SUBSPACE ( $d_j = 25$ )	288.65	56.34	<b>232.31</b>	✓	256
SUBSPACE ( $d_j = 20$ )	253.51	47.58	205.92	✓	256
SUBSPACE ( $d_j = 15$ )	211.97	38.04	173.93	✓	256
SUBSPACE ( $d_j = 10$ )	161.87	27.52	134.35	✓	256
SUBSPACE ( $d_j = 5$ )	98.35	15.55	82.79	✓	256
SUBSPACE ( $d_j = 1$ )	27.73	3.92	23.80	✓	256
<hr/>					
RANDOM GAUSSIAN	150.05	110.85	39.19	✓	128
SUBSPACE ( $d_j = 12$ )	144.36	27.72	<b>116.63</b>	✓	128
SUBSPACE ( $d_j = 10$ )	129.12	24.06	105.05	✓	128
SUBSPACE ( $d_j = 8$ )	112.01	20.18	91.83	✓	128
SUBSPACE ( $d_j = 6$ )	92.55	16.04	76.51	✓	128
SUBSPACE ( $d_j = 4$ )	69.57	11.51	58.06	✓	128
SUBSPACE ( $d_j = 2$ )	41.68	6.45	35.23	✓	128
SUBSPACE ( $d_j = 1$ )	24.28	3.57	20.70	✓	128
<hr/>					
SUBSPACE ( $d_j = 50$ )	145.60	75.31	70.29	✗	128
SUBSPACE ( $d_j = 40$ )	142.69	65.68	77.01	✗	128
SUBSPACE ( $d_j = 30$ )	135.42	54.27	81.15	✗	128
SUBSPACE ( $d_j = 20$ )	120.98	40.71	80.27	✗	128
SUBSPACE ( $d_j = 15$ )	111.10	32.89	78.21	✗	128
SUBSPACE ( $d_j = 12$ )	101.94	27.73	74.21	✗	128

### B.2 Implementation Details

**Training Setting.** We mainly use ResNet-18 [HZRS16] in our experiments, where we use 4 residual blocks with layer widths [64, 128, 256, 512]. The implementation of network architectures used in this paper are mainly based on this github repo.<sup>22</sup> For data augmentation in the supervised setting, we apply the RandomCrop and RandomHorizontalFlip. For the supervised setting, we train the models for 500 epochs and use stage-wise learning rate decay every 200 epochs (decay by a factor of 10). For the supervised setting, we train the models for 100 epochs and use stage-wise learning rate decay at 20-th epoch and 40-th epoch (decay by a factor of 10).

<sup>22</sup><https://github.com/kuangliu/pytorch-cifar>

**Evaluation Details.** For the supervised setting, we set the number of principal components for nearest subspace classifier  $r_j = 30$ . We also study the effect of  $r_j$  in Section B.3.2. For the CIFAR100 dataset, we consider 20 superclasses and set the cluster number as 20, which is the same setting as in [CWM<sup>+</sup>17, WXYL18].

**Datasets.** We apply the default datasets in PyTorch, including CIFAR10, CIFAR100, and STL10.

**Augmentations  $\mathcal{T}$  used for the self-supervised setting.** We apply the same data augmentation for CIFAR10 dataset and CIFAR100 dataset and the pseudo-code is as follows.

```
import torchvision.transforms as transforms
TRANSFORM = transforms.Compose([
    transforms.RandomResizedCrop(32),
    transforms.RandomHorizontalFlip(),
    transforms.RandomApply([transforms.ColorJitter(0.4, 0.4, 0.4, 0.1)], p=0.8),
    transforms.RandomGrayscale(p=0.2),
    transforms.ToTensor()])
```

The augmentations we use for STL10 dataset and the pseudo-code is as follows.

```
import torchvision.transforms as transforms
TRANSFORM = transforms.Compose([
    transforms.RandomResizedCrop(96),
    transforms.RandomHorizontalFlip(),
    transforms.RandomApply([transforms.ColorJitter(0.8, 0.8, 0.8, 0.2)], p=0.8),
    transforms.RandomGrayscale(p=0.2),
    GaussianBlur(kernel_size=9),
    transforms.ToTensor()])
```

**Cross-entropy training details.** For CE models presented in Table 1, Figure 6(d)-6(f), and Figure 7, we use the same network architecture, ResNet-18 [HZRS16], for cross-entropy training on CIFAR10, and set the output dimension as 10 for the last layer. We apply SGD, and set learning rate  $lr=0.1$ , momentum  $momentum=0.9$ , and weight decay  $wd=5e-4$ . We set the total number of training epoch as 400, and use stage-wise learning rate decay every 150 epochs (decay by a factor of 10).

### B.3 Additional Experimental Results

#### B.3.1 PCA Results of MCR<sup>2</sup> Training versus Cross-Entropy Training

For comparison, similar to Figure 3(c), we calculate the principle components of representations learned by MCR<sup>2</sup> training and cross-entropy training. For cross-entropy training, we take the output of the second last layer as the learned representation. The results are summarized in Figure 6. We also compare the cosine similarity between learned representations for both MCR<sup>2</sup> training and cross-entropy training, and the results are presented in Figure 7.

As shown in Figure 6, we observe that representations learned by MCR<sup>2</sup> are much more diverse, the dimension of learned features (each class) is around a dozen, and the dimension of the overall features is nearly 120, and the output dimension is 128. In contrast, the dimension of the overall features learned using entropy is slightly greater than 10, which is much smaller than that learned by MCR<sup>2</sup>. From Figure 7, for MCR<sup>2</sup> training, we find that the features of different class are almost orthogonal.

**Visualize representative images selected from CIFAR10 dataset by using MCR<sup>2</sup>.** As mentioned in Section 1, obtaining the properties of desired representation in the proposed MCR<sup>2</sup> principle is equivalent to performing *nonlinear generalized principle components* on the given dataset. As shown in Figure 6(a)-6(c), MCR<sup>2</sup> can indeed learn such diverse and discriminative representations. In order to better interpret the representations learned by MCR<sup>2</sup>, we select images according to their “principal” components (singular vectors using SVD) of the learned features. In Figure 8, we visualize images selected from class-‘Bird’ and class-‘Ship’. For each class, we first compute top-10 singular

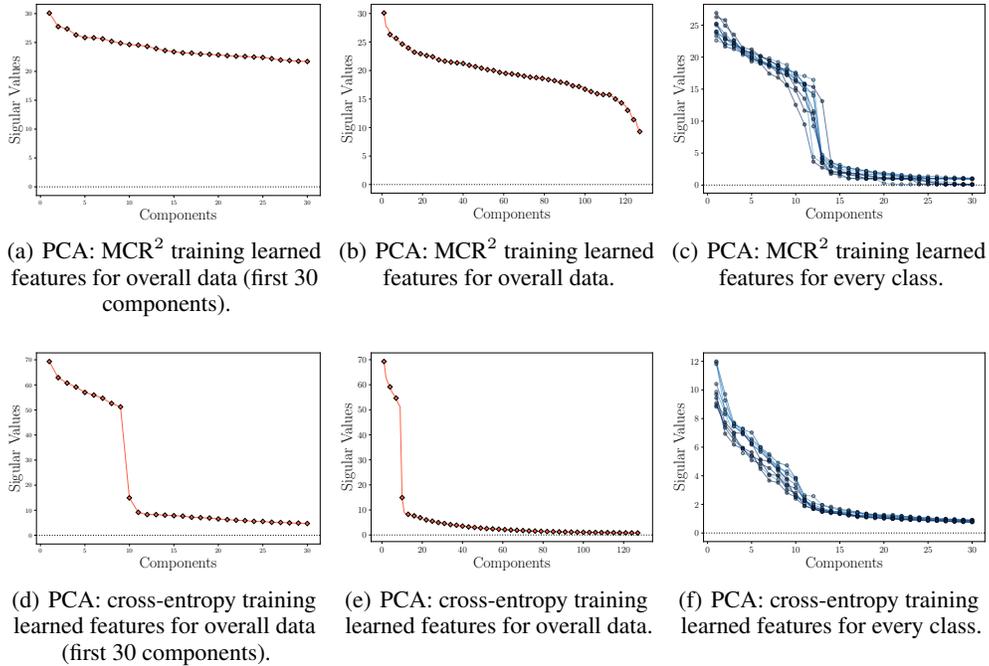


Figure 6: Principal component analysis (PCA) of learned representations for the MCR<sup>2</sup> trained model (**first row**) and the cross-entropy trained model (**second row**).

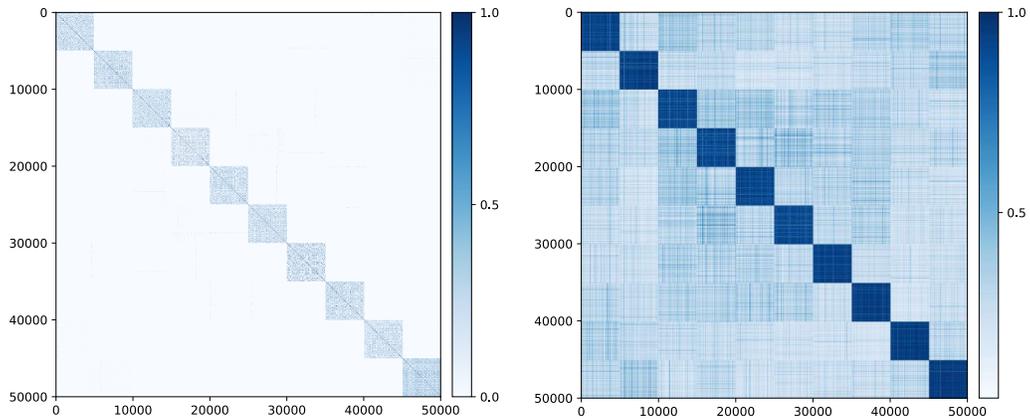


Figure 7: Cosine similarity between learned features by using the MCR<sup>2</sup> objective (**left**) and CE loss (**right**).

vectors of the SVD of the learned features and then for each of the top singular vectors, we display in each row the top-10 images whose corresponding features are closest to the singular vector. As shown in Figure 8, we observe that images in the same row share many common characteristics such as shapes, textures, patterns, and styles, whereas images in different rows are significantly different from each other – suggesting our method captures all the different “modes” of the data even within the same class. Notice that top rows are associated with components with larger singular values, hence they are images that show up more frequently in the dataset.

In Figure 9(a), we visualize the 10 “principal” images selected from CIFAR10 for each of the 10 classes. That is, for each class, we display the 10 images whose corresponding features are most coherent with the top-10 singular vectors. We observe that the selected images are much more diverse and representative than those selected randomly from the dataset (displayed on the CIFAR official website), indicating such principal images can be used as a good “summary” of the dataset.

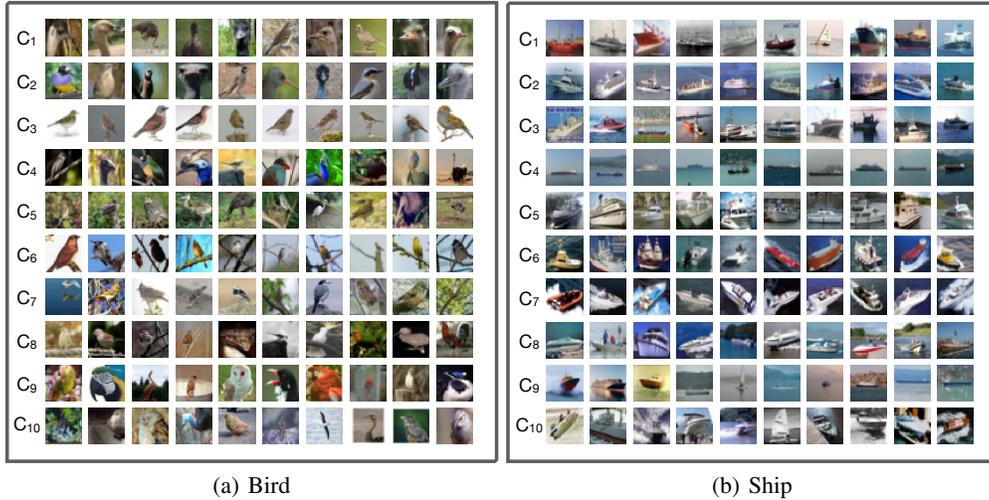


Figure 8: Visualization of principal components learned for class 2-‘Bird’ and class 8-‘Ship’. For each class  $j$ , we first compute the top-10 singular vectors of the SVD of the learned features  $\mathbf{Z}_j$ . Then for the  $l$ -th singular vector of class  $j$ ,  $\mathbf{u}_j^l$ , and for the feature of the  $i$ -th image of class  $j$ ,  $\mathbf{z}_j^i$ , we calculate the absolute value of inner product,  $|\langle \mathbf{z}_j^i, \mathbf{u}_j^l \rangle|$ , then we select the top-10 images according to  $|\langle \mathbf{z}_j^i, \mathbf{u}_j^l \rangle|$  for each singular vector. In the above two figures, each row corresponds to one singular vector (component  $C_l$ ). The rows are sorted based on the magnitude of the associated singular values, from large to small.



(a) 10 representative images from each class based on top-10 principal components of the SVD of learned representations by MCR<sup>2</sup>. (b) Randomly selected 10 images from each class.

Figure 9: Visualization of top-10 “principal” images for each class in the CIFAR10 dataset. (a) For each class- $j$ , we first compute the top-10 singular vectors of the SVD of the learned features  $\mathbf{Z}_j$ . Then for the  $l$ -th singular vector of class  $j$ ,  $\mathbf{u}_j^l$ , and for the feature of the  $i$ -th image of class  $j$ ,  $\mathbf{z}_j^i$ , we calculate the absolute value of inner product,  $|\langle \mathbf{z}_j^i, \mathbf{u}_j^l \rangle|$ , then we select the largest one for each singular vector within class  $j$ . Each row corresponds to one class, and each image corresponds to one singular vector, ordered by the value of the associated singular value. (b) For each class, 10 images are randomly selected in the dataset. These images are the ones displayed in the CIFAR dataset website [Kri09].

### B.3.2 Experimental Results of MCR<sup>2</sup> in the Supervised Learning Setting.

**Training details for mainline experiment.** For the model presented in Figure 1 (Right) and Figure 3, we use ResNet-18 to parameterize  $f(\cdot, \theta)$ , and we set the output dimension  $d = 128$ , precision  $\epsilon^2 = 0.5$ , mini-batch size  $m = 1,000$ . We use SGD in Pytorch [PGM<sup>+</sup>19] as the optimizer, and set the learning rate lr=0.01, weight decay wd=5e-4, and momentum=0.9.

**Experiments for studying the effect of hyperparameters and architectures.** We present the experimental results of MCR<sup>2</sup> training in the supervised setting by using various training hyperparameters and different network architectures. The results are summarized in Table 4. Besides the ResNet architecture, we also consider VGG architecture [SZ15] and ResNext architecture [XGD<sup>+</sup>17]. From Table 4, we find that larger batch size  $m$  can lead to better performance. Also, models with higher output dimension  $d$  require larger training batch size  $m$ .

Table 4: Experiments of MCR<sup>2</sup> in the supervised setting on the CIFAR10 dataset.

ARCH	DIM $d$	PRECISION $\epsilon^2$	BATCHSIZE $m$	LR	ACC	COMMENT
RESNET-18	128	0.5	1,000	0.01	92.20%	MAINLINE, FIG 3
RESNEXT-29	128	0.5	1,000	0.01	92.55%	DIFFERENT ARCHITECTURE
VGG-11	128	0.5	1,000	0.01	90.76%	
RESNET-18	512	0.5	1,000	0.01	88.60%	EFFECT OF OUTPUT DIMENSION
RESNET-18	256	0.5	1,000	0.01	92.10%	
RESNET-18	64	0.5	1,000	0.01	92.21%	
RESNET-18	128	1.0	1,000	0.01	93.06%	EFFECT OF PRECISION
RESNET-18	128	0.4	1,000	0.01	91.93%	
RESNET-18	128	0.2	1,000	0.01	90.06%	
RESNET-18	128	0.5	500	0.01	82.33%	EFFECT OF BATCH SIZE
RESNET-18	128	0.5	2,000	0.01	93.02%	
RESNET-18	128	0.5	4,000	0.01	92.59%	
RESNET-18	512	0.5	2,000	0.01	92.47%	
RESNET-18	512	0.5	4,000	0.01	92.17%	
RESNET-18	128	0.5	1,000	0.05	86.02%	EFFECT OF LR
RESNET-18	128	0.5	1,000	0.005	92.39%	
RESNET-18	128	0.5	1,000	0.001	92.23%	

**Effect of  $r_j$  on classification.** Unless otherwise stated, we set the number of components  $r_j = 30$  for nearest subspace classification. We study the effect of  $r_j$  when used for classification, and the results are summarized in Table 5. We observe that the nearest subspace classification works for a wide range of  $r_j$ .

Table 5: Effect of number of components  $r_j$  for nearest subspace classification in the supervised setting.

NUMBER OF COMPONENTS	$r_j = 10$	$r_j = 20$	$r_j = 30$	$r_j = 40$	$r_j = 50$
MAINLINE (LABEL NOISE RATIO=0.0)	92.68%	92.53%	92.20%	92.32%	92.17%
LABEL NOISE RATIO=0.1	91.71%	91.73%	91.16%	91.83%	91.78%
LABEL NOISE RATIO=0.2	90.68%	90.61%	89.70%	90.62%	90.54%
LABEL NOISE RATIO=0.3	88.24%	87.97%	88.18%	88.15%	88.10%
LABEL NOISE RATIO=0.4	86.49%	86.67%	86.66%	86.71%	86.44%
LABEL NOISE RATIO=0.5	83.90%	84.18%	84.30%	84.18%	83.76%

**Effect of  $\epsilon^2$  on learning from corrupted labels.** To further study the proposed MCR<sup>2</sup> on learning from corrupted labels, we use different precision parameters,  $\epsilon^2 = 0.75, 1.0$ , in addition to the one shown in Table 1. Except for the precision parameter  $\epsilon^2$ , all the other parameters are the same as the mainline experiment (the first row in Table 4). The first row ( $\epsilon^2 = 0.5$ ) in Table 6 is identical to the MCR<sup>2</sup> TRAINING in Table 2. Notice that with slightly different choices in  $\epsilon^2$ , one might even see slightly improved performance over the ones reported in the main body.

Table 6: Effect of Precision  $\epsilon^2$  on classification results with features learned with labels corrupted at different levels by using MCR<sup>2</sup> training.

PRECISION	RATIO=0.1	RATIO=0.2	RATIO=0.3	RATIO=0.4	RATIO=0.5
$\epsilon^2 = 0.5$	91.16%	89.70%	88.18%	86.66%	84.30%
$\epsilon^2 = 0.75$	<b>92.37%</b>	90.82%	<b>89.91%</b>	<b>87.67%</b>	83.69%
$\epsilon^2 = 1.0$	91.93%	<b>91.11%</b>	89.60%	87.09%	<b>84.53%</b>

### B.3.3 Experimental Results of MCR<sup>2</sup> in the Self-supervised Learning Setting

**Training details of MCR<sup>2</sup>-CTRL.** For three datasets (CIFAR10, CIFAR100, and STL10), we use ResNet-18 as in the supervised setting, and we set the output dimension  $d = 128$ , precision  $\epsilon^2 = 0.5$ , mini-batch size  $k = 20$ , number of augmentations  $n = 50$ ,  $\gamma_1 = \gamma_2 = 20$ . We observe that MCR<sup>2</sup>-CTRL can achieve better clustering performance by using smaller  $\gamma_2$ , i.e.,  $\gamma_2 = 15$ , on CIFAR10 and CIFAR100 datasets. We use SGD in Pytorch [PGM<sup>+</sup>19] as the optimizer, and set the learning rate lr=0.1, weight decay wd=5e-4, and momentum=0.9.

**Training dynamic comparison between MCR<sup>2</sup> and MCR<sup>2</sup>-CTRL.** In the self-supervised setting, we compare the training process for MCR<sup>2</sup> and MCR<sup>2</sup>-CTRL in terms of  $R$ ,  $\tilde{R}$ ,  $R^c$ , and  $\Delta R$ . For MCR<sup>2</sup> training, the features first expand (for both  $R$  and  $R^c$ ) then compress (for ). For MCR<sup>2</sup>-CTRL, both  $\tilde{R}$  and  $R^c$  first compress then  $\tilde{R}$  expands quickly and  $R^c$  remains small, as we have seen in Figure 5 in the main body.

**Clustering results comparison.** We compare the clustering performance between MCR<sup>2</sup> and MCR<sup>2</sup>-CTRL in terms of NMI, ACC, and ARI. The clustering results are summarized in Table 7. We find that MCR<sup>2</sup>-CTRL can achieve better performance for clustering.

Table 7: Clustering comparison between MCR<sup>2</sup> and MCR<sup>2</sup>-CTRL on CIFAR10 dataset.

	NMI	ACC	ARI
MCR <sup>2</sup>	0.544	0.570	0.399
MCR <sup>2</sup> -CTRL	0.630	0.684	0.508

### B.3.4 Clustering Metrics and More Results

We first introduce the definitions of normalized mutual information (NMI) [SG02], clustering accuracy (ACC), and adjusted rand index (ARI) [HA85].

**Normalized mutual information (NMI).** Suppose  $Y$  is the ground truth partition and  $C$  is the prediction partition. The NMI metric is defined as

$$\text{NMI}(Y, C) = \frac{\sum_{i=1}^k \sum_{j=1}^s |Y_i \cap C_j| \log \left( \frac{m |Y_i \cap C_j|}{|Y_i| |C_j|} \right)}{\sqrt{\left( \sum_{i=1}^k |Y_i| \log \left( \frac{|Y_i|}{m} \right) \right) \left( \sum_{j=1}^s |C_j| \log \left( \frac{|C_j|}{m} \right) \right)}},$$

where  $Y_i$  is the  $i$ -th cluster in  $Y$  and  $C_j$  is the  $j$ -th cluster in  $C$ , and  $m$  is the total number of samples.

**Clustering accuracy (ACC).** Given  $m$  samples,  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$ . For the  $i$ -th sample  $\mathbf{x}_i$ , let  $\mathbf{y}_i$  be its ground truth label, and let  $c_i$  be its cluster label. The ACC metric is defined as

$$\text{ACC}(Y, C) = \max_{\sigma \in S} \frac{\sum_{i=1}^m \mathbf{1}\{\mathbf{y}_i = \sigma(c_i)\}}{m},$$

where  $S$  is the set includes all the one-to-one mappings from cluster to label, and  $Y = [\mathbf{y}_1, \dots, \mathbf{y}_m]$ ,  $C = [c_1, \dots, c_m]$ .

**Adjusted rand index (ARI).** Suppose there are  $m$  samples, and let  $Y$  and  $C$  be two clustering of these samples, where  $Y = \{Y_1, \dots, Y_r\}$  and  $C = \{C_1, \dots, C_s\}$ . Let  $m_{ij}$  denote the number of the intersection between  $Y_i$  and  $C_j$ , i.e.,  $m_{ij} = |Y_i \cap C_j|$ . The ARI metric is defined as

$$\text{ARI} = \frac{\sum_{i,j} m_{ij} \binom{m_{ij}}{2} - \left( \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right) / \binom{m}{2}}{\frac{1}{2} \left( \sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right) - \left( \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right) / \binom{m}{2}},$$

where  $a_i = \sum_j m_{ij}$  and  $b_j = \sum_i m_{ij}$ .

**More experiments on the effect of hyperparameters of  $\text{MCR}^2\text{-CTRL}$ .** We provide more experimental results of  $\text{MCR}^2\text{-CTRL}$  training in the self-supervised setting by varying training hyperparameters on the STL10 dataset. The results are summarized in Table 8. Notice that the choice of hyperparameters only has small effect on the performance with the  $\text{MCR}^2\text{-CTRL}$  objective. We may hypothesize that, in order to further improve the performance, one has to seek other, potentially better, control of optimization dynamics or strategies. We leave those for future investigation.

Table 8: Experiments of  $\text{MCR}^2\text{-CTRL}$  in the self-supervised setting on STL10 dataset.

ARCH	PRECISION $\epsilon^2$	LEARNING RATE LR	NMI	ACC	ARI
RESNET-18	0.5	0.1	0.446	0.491	0.290
RESNET-18	0.75	0.1	0.450	0.484	0.288
RESNET-18	0.25	0.1	0.447	0.489	0.293
RESNET-18	0.5	0.2	0.477	0.473	0.295
RESNET-18	0.5	0.05	0.444	0.496	0.293
RESNET-18	0.25	0.05	0.454	0.489	0.294