

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Symmetric Encoding-decoding Framework without Backpropagation

Pengyuan Zhai ¹ ¹ Harvard University; pzhai@g.harvard.edu

Abstract: We propose a forward-only multi-layer encoding-decoding framework based on the principle of Maximal Coding Rate Reduction (MCR²), an information-theoretic metric that describes a statistical distance between two sets of feature vectors up to the second moment. The encoder directly transforms data vectors themselves via gradient ascent to maximize the MCR² distance between different classes in the feature space, resulting in class-wise mutually orthogonal subspace representations. The decoder follows a process symmetric to the encoder, and transforms the subspace feature vectors via gradient descent to minimize the MCR² distance between the reconstructed data and the original data. We show that the encoder transforms data to linear discriminative representations without breaking the higher-order manifolds, preserving higher-moment correlation/information. It is exactly due to this information-preserving property that the decoder is able to reconstruct the data with high fidelity.

Keywords: Backpropagation-free Network, Information Theory, Maximal Coding Rate Reduction

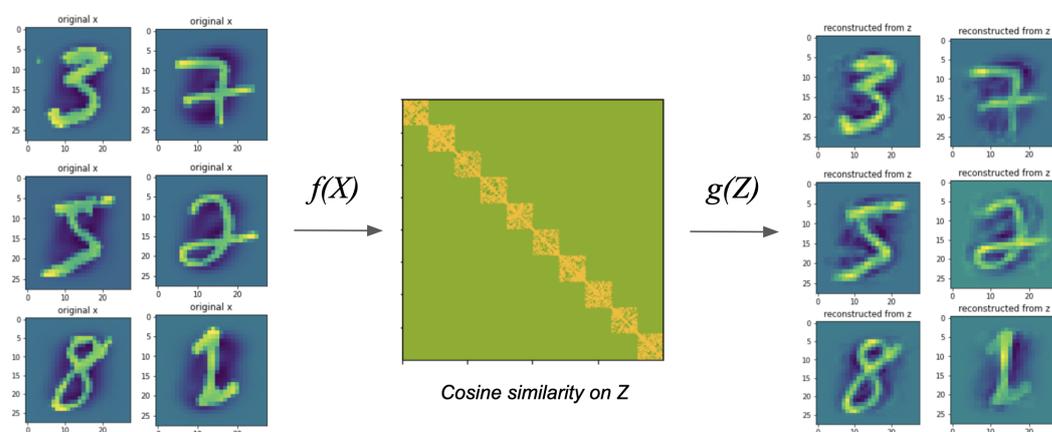


Figure 1. We propose a non-backpropagation encoding-decoding pipeline based on the Maximal Coding Rate Reduction (MCR²) principle. Given data vectors $\{X^{(c)}\}_C$ from C classes, the forward transformation corresponds to a gradient-based transformations on data vectors $\{X^{(c)}\}_C$ themselves, which simultaneously maps multi-class data vectors into linear subspaces $\{Z^{(c)}\}_C$ respective to each class. The decoding process follows a multi-layer pipeline that mirrors the structure of the encoding process, reconstructing $\{Z^{(c)}\}_C$ into $\{\tilde{X}^{(c)}\}_C$ with high visual fidelity with the original $\{X^{(c)}\}_C$.

0. Introduction and Related Work

One important goal of modern machine learning is to learn and model complex distributions of real-world data. A typical learning paradigm is to establish an (often parameterized) mapping between the distribution of the real data, i.e., $x \in \mathbb{R}^D$, and a more compact representation, i.e., $z \in \mathbb{R}^d$:

$$\text{encode: } f(\cdot, \theta_f) : x \in \mathbb{R}^D \mapsto z \in \mathbb{R}^d \quad \text{and/or decode: } g(\cdot, \theta_g) : z \in \mathbb{R}^d \mapsto x \in \mathbb{R}^D, \quad (1)$$

where z has simpler or more discriminative structures. The learnt feature z is thus easier to manipulate for both generative or discriminative purposes.

0.1. Auto-Encoding

One typical approach to encoding-decoding is the ‘‘Auto Encoder’’ (AE) [1,2], which learns an encoding mapping f from x to z and a decoding mapping g from z back to x :

$$\mathbf{X} \xrightarrow{f(x, \theta_f)} \mathbf{Z} \xrightarrow{g(z, \theta_g)} \hat{\mathbf{X}}. \quad (2)$$

To make the learned distribution $\hat{p}(z)$ close to a target $p(z)$, we minimize some distance metric, e.g., the KL-divergence: $\min \mathcal{D}_{KL}(\hat{p}, p)$. Between two arbitrary degenerate distributions in high-dimensional spaces, however, KL-divergence is often intractable. In practice, one instead minimizes instead certain bounds of the distance, e.g., the variational bound in *variational auto-encoding* (VAE) [3,4]. The exact distribution/structure of $\hat{p}(z|\mathbf{X}, \theta)$ is thus unclear.

0.2. Coding Rate Reduction

Recently, [5] [6] proposed a novel objective for learning a linear discriminative representation (LDR) for multi-class data, called ‘‘coding rate reduction’’. Unlike conventional discriminative methods only trying to predict class labels, the LDR learns the structure of the data, and is more suitable for both discriminative and generative purposes. [6] maps distributions of the input data on *multiple* nonlinear submanifolds to multiple distinctive linear subspaces. Intuitively, this approach generalizes nonlinear PCA [1] to settings where one simultaneously applies *multiple nonlinear PCAs* to data on multiple nonlinear submanifolds.

In this work, we extend the work of [6] and show that the *coding rate reduction* principle is also suitable for encoding-decoding purposes, and that one can build a pipeline that learns LDRs in the feature space and also reconstructs high-fidelity data in a forward-only fashion without back-propagation. The manifold structures can be learned directly from data itself, free of explicit parameterization.

1. Problem Formulation

In this section, we first explain the objective formulations for the forward (encoding) mapping and the backward (decoding) mapping, and then cover the details of the network structure. The main measurement of ‘‘closeness’’ between data distributions is the coding rate reduction metric [5]. Our goal is to come up with a encoding-decoding scheme:

$$\mathbf{X} \xrightarrow{f(x)} \mathbf{Z} \xrightarrow{g(z)} \hat{\mathbf{X}}, \quad (3)$$

such that the learnt features \mathbf{Z} of different classes are maximally discriminative, and the reconstructed data $\hat{\mathbf{X}}$ is the ‘‘closest’’ to the real data \mathbf{X} by the coding rate reduction metric.

1.1. Encoding mapping via maximal coding rate reduction

The (labeled) training data \mathbf{X} is the concatenation of data of all k classes: $\mathbf{X} \doteq [\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^k] \in \mathbb{R}^{D \times m}$, where $\mathbf{X}^j \in \mathbb{R}^{D \times m_j}$ denotes the j -th class’s data matrix, D is the dimension of each data sample, and m_j is the number of samples of class j . Following the notation in [6], we use the matrix $\Pi^j(i, i) = 1$ to denote the membership of sample i belonging to class j (and $\Pi^j = 0$ otherwise). For any general encoding mapping $f(\cdot, \theta_f) : x \mapsto z$ from \mathbf{X} to an optimal representation $\mathbf{Z} \in \mathbb{R}^{d \times m}$:

$$\mathbf{X} \xrightarrow{f(x, \theta_f)} \mathbf{Z}, \quad (4)$$

we wish to make each class's features more discriminative of others by maximize the coding rate reduction objective [5]:

$$\max_{\mathbf{Z}} \Delta R(\mathbf{Z} | \mathbf{\Pi}, \epsilon) \doteq \underbrace{\frac{1}{2} \log \det \left(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^* \right)}_{R(\mathbf{Z} | \epsilon)} - \underbrace{\sum_{j=1}^k \frac{\gamma_j}{2} \log \det \left(\mathbf{I} + \alpha_j \mathbf{Z} \mathbf{\Pi}^j \mathbf{Z}^* \right)}_{R_c(\mathbf{Z} | \mathbf{\Pi}, \epsilon)}, \quad (5)$$

where $\alpha = \frac{d}{m\epsilon^2}$, $\alpha_j = \frac{d}{\text{tr}(\mathbf{\Pi}^j)\epsilon^2}$, $\gamma_j = \frac{\text{tr}(\mathbf{\Pi}^j)}{m}$ for $j = 1, \dots, k$, and \mathbf{Z}^* denotes the transpose of \mathbf{Z} . For simplicity we denote $\Delta R(\mathbf{Z} | \mathbf{\Pi}, \epsilon)$ as $\Delta R(\mathbf{Z})$, assuming $\mathbf{\Pi}, \epsilon$ are known and fixed. The first term $R(\mathbf{Z} | \epsilon)$ is the coding rate of the whole set of features \mathbf{Z} (coded as a Gaussian source) with precision ϵ ; the second term $R_c(\mathbf{Z} | \mathbf{\Pi}, \epsilon)$ is the weighted average coding rate of the k subsets of feature, each coded as a Gaussian.

For more readability, Eq. 5 is equivalently:

$$\max_{\mathbf{Z}} \Delta R(\mathbf{Z} | \mathbf{\Pi}, \epsilon) \doteq \underbrace{\frac{1}{2} \log \det \left(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^* \right)}_{R(\mathbf{Z} | \epsilon)} - \underbrace{\sum_{j=1}^k \frac{\gamma_j}{2} \log \det \left(\mathbf{I} + \alpha_j \mathbf{Z}^j \mathbf{Z}^{j*} \right)}_{R_c(\mathbf{Z} | \mathbf{\Pi}, \epsilon)}, \quad (6)$$

given known class partitions of the feature set: $\mathbf{Z} = [\mathbf{Z}^1, \mathbf{Z}^2, \dots, \mathbf{Z}^k] \in \mathbb{R}^{d \times m}$.

1.2. Decoding via coding rate matching.

The pair-wise coding rate reduction metric for two identical sets of data samples is equal to the minimal value of zero. To match the distribution of the reconstructed data to that of the real data, for any general decoding mapping $g(\cdot, \theta_g) : \mathbf{z} \mapsto \hat{\mathbf{x}}$ from the encoded \mathbf{Z} features to an optimal reconstructed $\hat{\mathbf{X}}$:

$$\mathbf{Z} \xrightarrow{g(\cdot, \theta_g)} \hat{\mathbf{X}}, \quad (7)$$

we propose the following objective:

$$\min_{\hat{\mathbf{X}}} \sum_{j=1}^k \underbrace{\Delta R \left([\mathbf{X}^j, \hat{\mathbf{X}}^j] \right)}_{\text{coding rate matching}} = \sum_{j=1}^k \left\{ \underbrace{\frac{1}{2} \log \det \left(\mathbf{I} + \alpha (\mathbf{X}^j \mathbf{X}^{j*} + \hat{\mathbf{X}}^j \hat{\mathbf{X}}^{j*}) \right)}_{R(\mathbf{Z} | \epsilon)} - \underbrace{\frac{1}{4} \log \det \left(\mathbf{I} + 2\alpha \mathbf{X}^j \mathbf{X}^{j*} \right) - \frac{1}{4} \log \det \left(\mathbf{I} + 2\alpha \hat{\mathbf{X}}^j \hat{\mathbf{X}}^{j*} \right)}_{R_c(\mathbf{Z} | \mathbf{\Pi}, \epsilon)} \right\}, \quad (8)$$

where $\alpha = \frac{d}{2m_j\epsilon^2}$, if we assume that the number of data samples in \mathbf{X}^j and $\hat{\mathbf{X}}^j$ are equal.

Eq. 8 minimizes the pairwise coding rate reduction of the real data \mathbf{X}^j and the reconstructed data $\hat{\mathbf{X}}^j$ of class j . If the reconstructed data is exactly identical to the real data for all classes, the objective evaluates to zero.

2. Forward-only Encoding-decoding Pipeline

The encoding and decoding objectives (Eq. 6 and Eq. 8) are suitable for any general parameterized functions, which can be readily modeled by deep neural networks. In this work, we propose a non-parameterized, forward-only pipeline that does not require training through back-propagation.

There are two key factors that make such pipeline viable (details will be discussed in later sections):

1. **Differentiable gradients.** Eq. 6 and Eq. 8 are differentiable with \mathbf{Z} and $\hat{\mathbf{X}}$ respectively. Given any initial configuration of \mathbf{Z} or $\hat{\mathbf{X}}$, one can directly apply gradient ascent/descent steps to transform the data features to increase/decrease the objective function, following a normalization step after each update.
2. **Structure-preserving transformation.** According to our numerical simulations, both the encoding and decoding mappings (through gradient ascent/descent) preserve the local topological relationships between data points. The encoding and decoding transformations resembles a Ricci flow, which preserves local structures of the manifold [7]. Additionally, due to the normalization (by Frobenius norm) after each gradient update, the Riemannian metric of the manifold is constant, providing a fair comparison of coding rate reduction metrics at different layers.

2.0.1. Encoding via Gradient-based Transformation

Consider the encoding objective in Eq.6, if we directly optimize the objective $\Delta R(\mathbf{Z})$ as a function of the (normalized) data feature \mathbf{Z} . The gradient ascent step is:

$$\mathbf{Z}_{\ell+1} \propto \mathbf{Z}_{\ell} + \eta \cdot \left. \frac{\partial \Delta R}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_{\ell}} \text{ subject to } \|\mathbf{Z}_{\ell+1}\|_F = 1. \quad (9)$$

This means that at each layer, we apply a gradient ascent transformation directly on data features themselves to incrementally improve the coding rate reduction metric ΔR .

The each part of the gradient $\left. \frac{\partial \Delta R}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_{\ell}}$ is computed as (refer to [8]):

$$\begin{aligned} \left. \frac{1}{2} \frac{\partial \log \det(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*)}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_{\ell}} &= \underbrace{\alpha (\mathbf{I} + \alpha \mathbf{Z}_{\ell} \mathbf{Z}_{\ell}^*)^{-1}}_{\mathbf{E}_{\ell} \in \mathbb{R}^{d \times d}} \mathbf{Z}_{\ell} \in \mathbb{R}^{d \times m}, \\ \left. \frac{1}{2} \frac{\partial (\gamma_j \log \det(\mathbf{I} + \alpha_j \mathbf{Z}^j \mathbf{Z}^{j*}))}{\partial \mathbf{Z}^j} \right|_{\mathbf{Z}_{\ell}^j} &= \gamma_j \alpha_j \underbrace{(\mathbf{I} + \alpha_j \mathbf{Z}_{\ell}^j \mathbf{Z}_{\ell}^{j*})^{-1}}_{\mathbf{C}_{\ell}^j \in \mathbb{R}^{d \times d}} \mathbf{Z}_{\ell}^j \in \mathbb{R}^{d \times m_j}, \end{aligned} \quad (10)$$

where $\alpha = \frac{d}{m\epsilon^2}$. Thus the total gradient of the encoding objective with respect to \mathbf{Z}_{ℓ}^j of class j is:

$$\left. \frac{\partial \Delta R}{\partial \mathbf{Z}^j} \right|_{\mathbf{Z}_{\ell}^j} = \underbrace{\mathbf{E}_{\ell}}_{\text{Expansion}} \mathbf{Z}_{\ell}^j - \sum_{j=1}^k \gamma_j \underbrace{\mathbf{C}_{\ell}^j}_{\text{Compression}} \mathbf{Z}_{\ell}^j \in \mathbb{R}^{n \times m}. \quad (11)$$

2.0.2. Symmetric Decoding

Each of the k terms in the summation of Eq. 8 has the following gradient with respect to the reconstructed data $\hat{\mathbf{X}}^j$ of class j (with $\alpha = \frac{d}{2m_j\epsilon^2}$):

$$\left. \frac{\partial \Delta R([\mathbf{X}^j, \hat{\mathbf{X}}^j])}{\partial \hat{\mathbf{X}}^j} \right|_{\hat{\mathbf{X}}_{\ell}^j} = \underbrace{\alpha (\mathbf{I} + \alpha (\mathbf{X}_{\ell}^j \mathbf{X}_{\ell}^{j*} + \hat{\mathbf{X}}_{\ell}^j \hat{\mathbf{X}}_{\ell}^{j*}))^{-1}}_{\mathbf{E}_{\ell}^j} \hat{\mathbf{X}}_{\ell}^j - \frac{1}{2} \underbrace{2\alpha (\mathbf{I} + 2\alpha \hat{\mathbf{X}}_{\ell}^j \hat{\mathbf{X}}_{\ell}^{j*})^{-1}}_{\mathbf{C}_{\ell}^j} \hat{\mathbf{X}}_{\ell}^j. \quad (12)$$

Thus, the total gradient of the decoding objective with respect to $\hat{\mathbf{X}}^j$ is

$$\left. \frac{\partial \Delta R([\mathbf{X}^j, \hat{\mathbf{X}}^j])}{\partial \hat{\mathbf{X}}^j} \right|_{\hat{\mathbf{X}}_{\ell}^j} = \left(\mathbf{E}_{\ell}^j - \frac{1}{2} \mathbf{C}_{\ell}^j \right) \hat{\mathbf{X}}_{\ell}^j. \quad (13)$$

3. Pipeline Architecture 80

With the knowledge of the gradients, we now build the layer-by-layer pipeline for forward-only encoder f and forward-only decoder g :

$$\mathbf{X} \xrightarrow{f(x)} \mathbf{Z} \xrightarrow{g(z)} \hat{\mathbf{X}}. \quad (14)$$

3.1. Encoding Pipeline 81

Assuming a fixed number of L_f layers, define the encoding mapping function as $f(\mathbf{X}) \doteq f_{L_f} \circ f_{L_f-1} \circ \dots \circ f_2 \circ f_1(\mathbf{X})$, where $\mathbf{Z}_\ell = f_\ell(\mathbf{Z}_{\ell-1})$ is the transformation function acting on the input ($\mathbf{Z}_{\ell-1}$) to the ℓ -th layer:

$$\begin{aligned} \mathbf{Z}_\ell^j &= f_\ell(\mathbf{Z}_{\ell-1}^j) \doteq FN \left(\mathbf{Z}_\ell^j + \eta \frac{\partial \Delta R}{\partial \mathbf{Z}^j} \Big|_{\mathbf{Z}_\ell^j} \right) = FN \left(\mathbf{Z}_\ell^j + \eta \left(\mathbf{E}_\ell \mathbf{Z}_\ell^j - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{Z}_\ell^j \right) \right) \\ &= FN \left(\mathbf{Z}_\ell^j + \eta \left(\mathbf{E}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \right) \mathbf{Z}_\ell^j \right) = FN \left(\left[\mathbf{I} + \eta \left(\mathbf{E}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \right) \right] \mathbf{Z}_\ell^j \right), \end{aligned} \quad (15)$$

where $FN(\cdot)$ is the Frobenius normalization function:

$$FN(\mathbf{Z}) = \frac{\mathbf{Z}}{\|\mathbf{Z}\|_F} \quad (16)$$

3.2. Decoding pipeline. 82

Assuming a fixed number of L_f layers, define the decoding mapping function as $g(\mathbf{Z}_L) \doteq g_{L_g} \circ g_{L_g-1} \circ \dots \circ g_2 \circ g_1(\mathbf{Z}_L) = \hat{\mathbf{X}}$, where $\hat{\mathbf{X}}$ denotes the decoded features in the data space (referred to as the "reconstructed data"). We use a similar procedure:

$$\begin{aligned} \hat{\mathbf{X}}_\ell^j &= f_\ell(\hat{\mathbf{X}}_{\ell-1}^j) \doteq FN \left(\hat{\mathbf{X}}_\ell^j + \eta \frac{\partial \Delta R_g}{\partial \hat{\mathbf{X}}^j} \Big|_{\hat{\mathbf{X}}_\ell^j} \right) = FN \left(\hat{\mathbf{X}}_\ell^j + \eta \left(\mathbf{E}_\ell^j - \frac{1}{2} \mathbf{C}_\ell^j \right) \hat{\mathbf{X}}_\ell^j \right) \\ &= FN \left(\left[\mathbf{I} + \eta \left(\mathbf{E}_\ell^j - \frac{1}{2} \mathbf{C}_\ell^j \right) \right] \hat{\mathbf{X}}_\ell^j \right). \end{aligned} \quad (17)$$

We make the following remarks: 83

1. For encoding (f), we have that $\mathbf{Z}_0 = \mathbf{X}$, and the last layer outputs \mathbf{Z}_L . The deepest-layer features \mathbf{Z}_L is then fed into the decoder (g) to produce $\hat{\mathbf{X}}$. 84
2. The structure of the pipeline bears resemblance to that of a neural network with residual connections (by the appearance of the identity matrix \mathbf{I}). The normalization steps can be seen as a non-linear function applied on the piece-wise linearly transformed features (piece-wise linear by each class). 85
86
87
88
89

4. Numerical Simulations 90

In this section, we present results from the experiments to demonstrate: 91

1. How the encoder transforms the data \mathbf{X} to a linear discriminative representation (LDR) \mathbf{Z} , while preserving the local manifold structure; 92
93
2. How the decoder transforms/unfolds \mathbf{Z} to $\hat{\mathbf{X}}$ such that the principal components of $\hat{\mathbf{X}}^j$ align well with the those of \mathbf{X}^j at the optimal configuration; 94
95
3. The importance of preserving the local manifold structures during encoding in order to assure successful decoding; 96
97
4. The performance of our pipeline on the MNIST dataset. 98

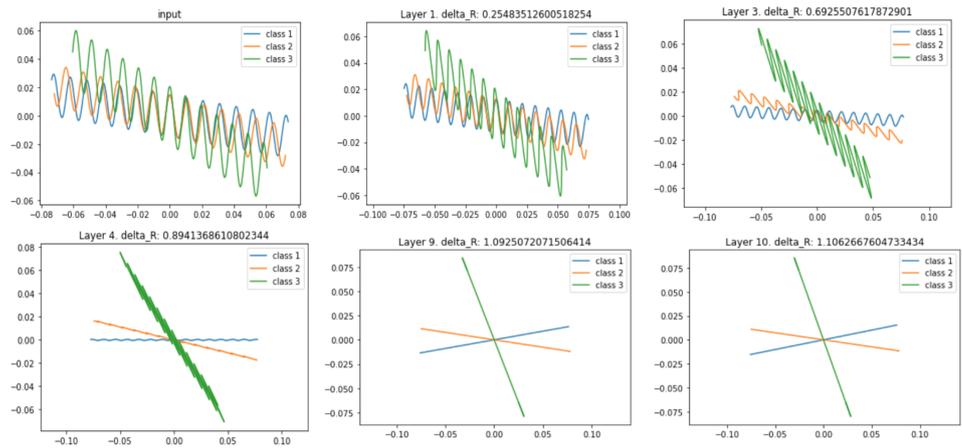


Figure 2. Visualization of the encoding process on 2-dimensional data vectors from 3 classes. At deeper layers, data vectors from each class are compressed to a linear subspace. Later we show that the higher-order “micro” manifolds are preserved in this process, allowing reconstruction/decoding with high-fidelity.

4.1. Experiment 1: Encoding while preserving manifold structures

For each of the two classes, we generate 500 2-dimensional (x, y) pairs as the training data with the relationship:

$$y = \frac{(2 \cdot \cos(x \cos a - x \sin a) - x \sin a)}{0.3 \cos a}, \quad (18)$$

where a controls the rotation around the origin.

We use a 10-layer encoder with $\epsilon = 0.01$ and the learning rate $\eta_f = 0.1$ and plot the features output by each layer. Two key observations are: i. the dominant principal components of each class becomes farther apart in deeper layers, in the last layer, the manifolds are compressed to LDRs; ii. the local structure of the manifolds are preserved (notice the zig-zag shapes of the manifolds). Although the manifolds are hard to visualize in deeper layers, we demonstrate in the next experiment that the manifolds structures are indeed preserved (Fig. 2).

4.2. Experiment 2: Decoding by manifold unfolding and component alignment

We first use a 30-layer encoder with $\epsilon = 0.01$ and the learning rate $\eta_f = 0.01$ and plot the features output by the last layer (Fig. 3), which will be the starting \mathbf{Z} which we will decode from. We transform this last layer \mathbf{Z} using the decoding pipeline (Section 3) and reconstruct $\tilde{\mathbf{X}}$ to match the distribution of \mathbf{X} . The decoder has 10 layers with the learning rate $\eta_g = 0.1$ (Fig. 4)

Notice that the decoder does two things: i. it “unfolds” the manifolds from a compressive state to a more expanded state; ii. it aligns the principal components of $\tilde{\mathbf{X}}$ to those of \mathbf{X} by gradually rotating the manifolds.

4.3. Experiment 3: Matching two different manifolds

In this experiment we transform a \mathbf{Z} to match \mathbf{X} with a different manifold structure. Namely we feed into the decoder a Gaussian-like \mathbf{Z} to match \mathbf{X} , which has the sinusoidal manifolds. We use a 20-layer decoder with learning rate $\eta_g = 0.1$:

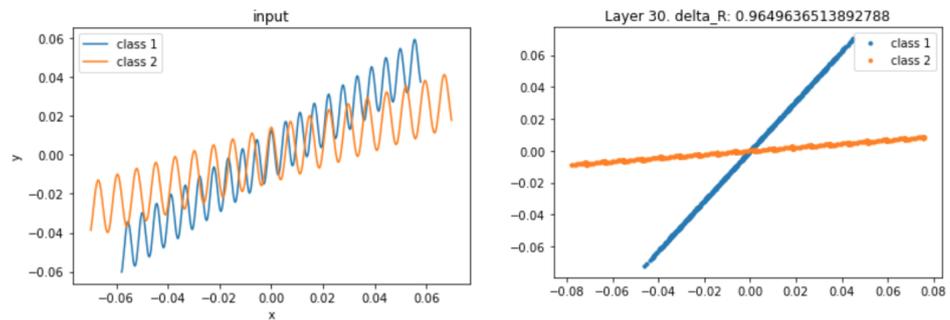


Figure 3. To prepare for the decoding, we first prepare the encoding process with two classes. We then transform the compressed data vectors (right-hand side) into the original manifold as shown below (Fig. 4).

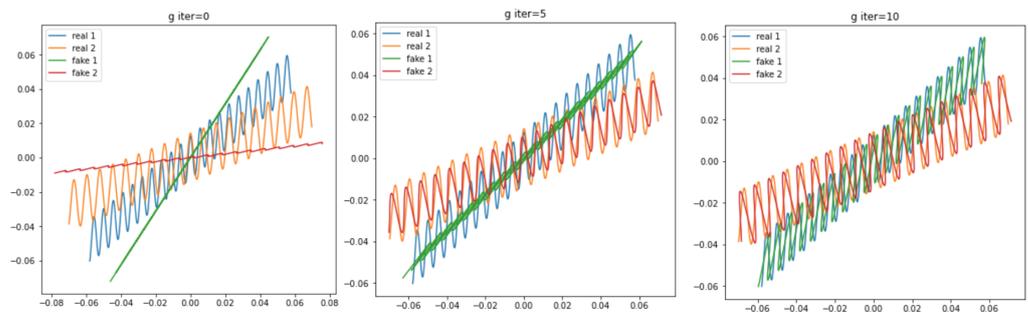


Figure 4. Visualization of the decoding process. The green (compressed class 1) and red (compressed class 2) linear subspaces are the compressed data vectors from Fig. 3, juxtaposed on the blue (original class 1) and orange (original data 2) data points from the true distribution. We transform the compressed data vectors based on the decoding pipeline described in Section 3 until convergence. The reconstructed data vectors (red and green) match the true data distribution with high fidelity, as the “micro-manifolds” in the compressed state are able to expand back to the original shape.

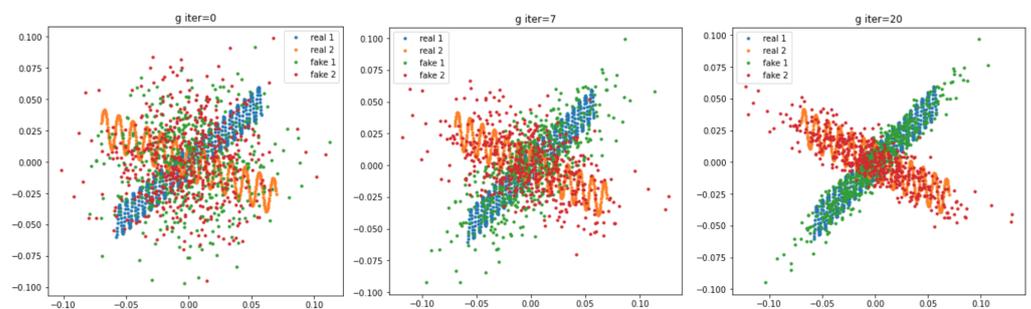


Figure 5. Visualization of a decoding process, where the target X are the blue and orange “spiral” distributions, and the input Z are two random Gaussian distributions (green and red). The decoding progress matches the second moment of the reconstructed \tilde{X} and the real X in the data space.

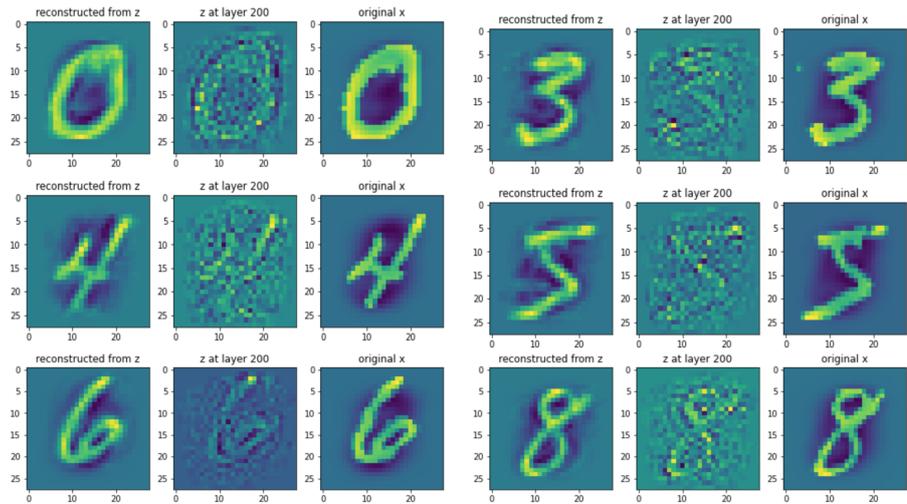


Figure 6. The original data X (right column), the compressed data in respective linear subspaces Z (middle column) and the reconstructed data \hat{X} (left column) for a selection of MNIST digits.

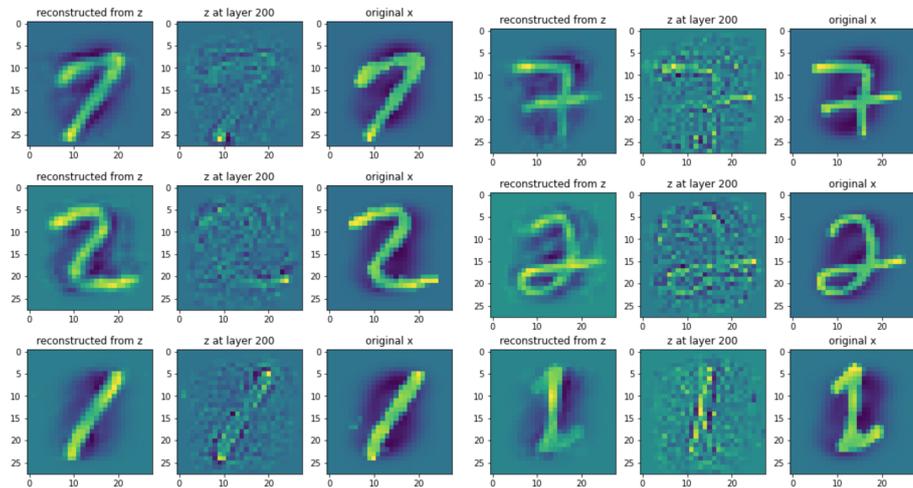


Figure 7. For different styles of digits, the decoder is able to reconstruct with high visual fidelity. For both sub-panels, right columns are the real input X , middle columns are the transformed images Z that lie on respective linear subspaces, and the left columns are the reconstructed images \hat{X} from Z .

Notice that the principal components of \hat{X} gradually align with those of X , but due to the intrinsic differences of the manifolds (Gaussian vs sinusoidal), the reconstructed \hat{X} does not learn the local manifold structures if only using the coding rate reduction objective.

4.4. Experiment 4: MNIST

We perform encoding-decoding on randomly selected MNIST digits [9]. For each of the 10 classes, we randomly select 200 images, resize each $28p \times 28p$ image into a 784-dimensional vector, and then center the data. With $\epsilon = 0.1$ we use a 200-layer encoder with learning rate $\eta_f = 0.003$ to arrive at the deepest $Z_L \in R^{784 \times 2000}$ (we keep the number of features to 784). From Z_L , we reconstruct \hat{X} using a 200-layer decoder with learning rate $\eta_g = 0.001$. We present some selected samples to showcase the qualitative decoding performance (Fig 6, 7).

For each digit, the left column is the reconstructed image, the middle is the deepest features Z_L plotted as an image, and the right is the original input. The reconstructed \hat{X} almost exactly replicate the original input, given a seemingly low-quality Z_L . We also notice that different modes of the same digits are also captured by the decoder (Fig. 7).

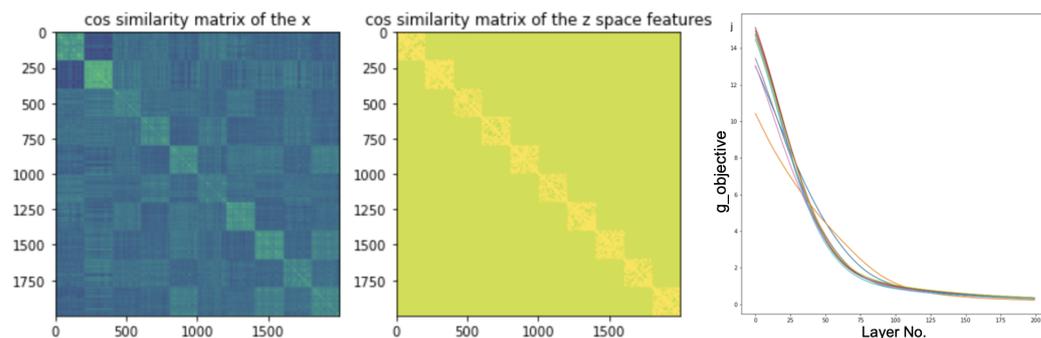


Figure 8. Left: the cosine similarity between any two pair of MNIST images (grouped by the 10 classes). In the original data space, there is no distinct similarity pattern. However, when we transform the original image vectors X into Z , and visualize the cosine similarity of vectors from Z (middle), we see high intra-class similarities as showcased by the block-diagonal pattern.

To investigate whether Z^j for class j is truly a linear discriminative representation, we plot the cosine similarity matrix of the the deepest learnt features Z_L and the cosine similarity matrix of the original data X , where the data points are grouped by classes (Fig. 8). It is clear that in the Z space, the learnt data features Z^j for class j has high cosine similarity with other data points in the same class, while far away from other classes (shown by the 10 bright diagonal blocks on the left plot), meanwhile the unencoded X is not a good linear discriminative representation (right plot).

Additionally, we show the stability of the decoder but plotting the decoder objective value (Eq. 8) at different layer of g , which shows very nice convergent properties for all classes (represented by each curve, Fig. 8 right).

5. Conclusions and Future Work

We proposed a non-parametric, forward-only, multi-layer encoding-decoding framework based on the principle of maximal coding rate reduction [5], an information-theoretic metric that measures the distance between two sets of feature vectors which are often degenerate with high dimensions. By designing the encoding and decoding objectives based on the *coding rate reduction* principle, we are able to derive gradients for both the encoding and decoding mappings with respect to data/features themselves. The encoder directly applies gradient ascent on data itself to maximize the distances between different classes, while the decoder transforms the learnt features by directly applying gradient descent to minimize the distance between the reconstructed and the original data. Through numerical simulations, we show that the encoder transforms data to linear discriminative representations, and the decoder reconstructs the data with high fidelity. Future work includes investigating the exact connections to the Ricci Flow, and under what conditions this pipeline will converge/diverge. More experiments shall also be done to investigate the performance on larger, more complex datasets such as the Cifar [10] datasets and the ImageNet[11] dataset.

References

1. Kramer, M.A. Nonlinear principal component analysis using autoassociative neural networks. *Aiche Journal* **1991**, *37*, 233–243.
2. Hinton, G.E.; Zemel, R.S. Autoencoders, Minimum Description Length and Helmholtz Free Energy. In Proceedings of the Proceedings of the 6th International Conference on Neural Information Processing Systems; Morgan Kaufmann Publishers Inc.: San Francisco, CA, USA, 1993; NIPS'93, p. 3–10.
3. Kingma, D.P.; Welling, M. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114* **2013**.
4. Zhao, S.; Song, J.; Ermon, S. InfoVAE: Information maximizing variational autoencoders. *arXiv preprint arXiv:1706.02262* **2017**.
5. Yu, Y.; Chan, K.H.R.; You, C.; Song, C.; Ma, Y. Learning Diverse and Discriminative Representations via the Principle of Maximal Coding Rate Reduction. In Proceedings of the Advances in neural information processing systems, 2020.
6. Chan, K.H.R.; Yu, Y.; You, C.; Qi, H.; Wright, J.; Ma, Y. ReduNet: A White-box Deep Network from the Principle of Maximizing Rate Reduction. *CoRR* **2021**, *abs/2105.10446*, [2105.10446].

-
7. Li, Y.; Lu, R. Applying Ricci Flow to High Dimensional Manifold Learning, 2017, [[arXiv:cs.LG/1703.10675](https://arxiv.org/abs/cs/1703.10675)]. 176
 8. Chan, K.H.R.; Yu, Y.; You, C.; Qi, H.; Wright, J.; Ma, Y. Deep Networks from the Principle of Rate Reduction, 2020, [[arXiv:cs.LG/2010.14765](https://arxiv.org/abs/cs/2010.14765)]. 177
 9. Deng, L. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine* **2012**, 29, 141–142. 179
 10. Krizhevsky, A. Learning Multiple Layers of Features from Tiny Images. 2009. 180
 11. Deng, J.; Dong, W.; Socher, R.; Li, L.J.; Li, K.; Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In Proceedings of the 2009 IEEE conference on computer vision and pattern recognition. Ieee, 2009, pp. 248–255. 182