

# Incremental Learning via Rate Reduction

*Kyung Eun Baek*  
*Yi Ma*

Electrical Engineering and Computer Sciences  
University of California, Berkeley

Technical Report No. UCB/EECS-2021-130

<http://www2.eecs.berkeley.edu/Pubs/TechRpts/2021/EECS-2021-130.html>

May 14, 2021



Copyright © 2021, by the author(s).  
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

---

# Incremental Learning via Rate Reduction

by Christina Baek

## Research Project

Submitted to the Department of Electrical Engineering and Computer Sciences,  
University of California at Berkeley, in partial satisfaction of the requirements for the  
degree of **Master of Science, Plan II**.

Approval for the Report and Comprehensive Examination:

### Committee:

*Yi Ma*

---

Professor Yi Ma  
Research Advisor

May 13, 2021

---

(Date)

\* \* \* \* \*

*Anant Sahai*

---

Professor Anant Sahai  
Second Reader

May 14, 2021

---

(Date)

Incremental Learning via Rate Reduction

by

Christina Baek

A thesis submitted in partial satisfaction of the

requirements for the degree of

Master of Science

in

Electrical Engineering and Computer Science

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor Yi Ma, Chair

Professor Anant Sahai

Spring 2021

Abstract

Incremental Learning via Rate Reduction

by

Christina Baek

Master of Science in Electrical Engineering and Computer Science

University of California, Berkeley

Professor Yi Ma, Chair

Current deep learning architectures suffer from catastrophic forgetting, a failure to retain knowledge of previously learned classes when incrementally trained on new classes. The fundamental roadblock faced by deep learning methods is that the models are optimized as “black boxes,” making it difficult to properly adjust the model parameters to preserve knowledge about previously seen data. To overcome the problem of catastrophic forgetting, we propose utilizing an alternative “white box” architecture derived from the principle of rate reduction, where each layer of the network is explicitly computed without back propagation. Under this paradigm, we demonstrate that, given a pretrained network and new data classes, our approach can provably construct a new network that emulates joint training with all past and new classes. Finally, our experiments show that our proposed learning algorithm observes significantly less decay in classification performance, outperforming state of the art methods on MNIST and CIFAR-10 by a large margin and justifying the use of “white box” algorithms for incremental learning even for sufficiently complex image data.

# Contents

<b>Contents</b>	<b>i</b>
<b>List of Figures</b>	<b>ii</b>
<b>List of Tables</b>	<b>iii</b>
<b>1 Interpretable Network from Rate Reduction</b>	<b>1</b>
1.1 Principle of Maximal Coding Rate Reduction . . . . .	1
1.2 Rate Reduction Network (ReduNet) . . . . .	3
<b>2 ReduNet for Interpretable Incremental Learning</b>	<b>5</b>
2.1 Introduction . . . . .	5
2.2 Related Work . . . . .	6
2.3 Incremental Learning with ReduNet . . . . .	8
2.4 Experiments . . . . .	12
2.5 Conclusions and Future Work . . . . .	16
<b>3 Low Rank Approximation for Efficient ReduNet</b>	<b>17</b>
3.1 Low rank approximation . . . . .	17
3.2 Block Diagonal Approximation by Eigenspace . . . . .	20
3.3 Conclusion . . . . .	25
<b>Bibliography</b>	<b>26</b>

# List of Figures

1.1	$MCR^2$ learns a map $f(\mathbf{x}, \theta)$ such that the class representation $\mathbf{Z}_j = f(\mathbf{X}_j, \theta)$ lies on maximally uncorrelated subspaces $\{S_j\}$ . . . . .	2
1.2	ReduNet Architecture in which we here adopt a slightly different normalization than [2] but is more suitable for the incremental learning as we will see in our derivation. . . . .	4
2.1	The joint network can be derived using simply $\mathbf{Z}_{0,t}\mathbf{Z}_{0,t}^\top$ . We do not need the task $t$ data $\mathbf{Z}_{0,t}$ directly. . . . .	10
2.2	Incremental learning results (accuracy) on MNIST and CIFAR-10. Both datasets have 5 incremental batches. We also provide the upper bound (UB) given by joint training a model utilizing the same architecture as the baseline methods. In <i>solid</i> lines are regularization-based methods and in <i>dashed</i> are exemplar-based methods, which saves 200 samples from previous tasks. Note that the decay in the performance in ReduNet is simply because classification is harder to accomplish with more classes, not because of catastrophic forgetting. . . . .	15
3.1	The first 100 eigenvalues of each class distribution of MNIST and CIFAR10. MNIST is $32 \times 32 = 784$ dimensional and CIFAR10 is $28 \times 28 \times 3 = 3072$ dimensional. Note that eigenvalues nearly decay to 0 by the 20th dimension for both datasets. . . . .	18
3.2	We pass through the training data of MNIST and CIFAR10 through the original ReduNet and compressed ReduNet, and compare the MCR loss $\Delta R(\mathbf{Z}_l)$ at each layer $l$ . We observe the loss curve for compressed ReduNet to be better than that of true ReduNet . . . . .	20
3.3	In MNIST, we observe that only a few eigenvectors are necessary to approximate $C_j$ of the first layer for small tolerance values. . . . .	21
3.4	ReduNet trained on MNIST with 500 samples. Notice that the number of eigenvectors chosen is proportional to the compression loss . . . . .	24

# List of Tables

2.1 Test Accuracy (%) on Task 1 After Each Training Session on MNIST and CIFAR-10. 14



## Acknowledgments

This work was a joint effort with students Ziyang Wu and Chong You, advised by Professor Yi Ma. I would like to thank Professor Yi Ma for the invaluable opportunities he has given me this year to learn and grow as a student and researcher. I learned tremendously from him and graduate students Simon, Yaodong, Chong, Xili, and others, whose dedication and compassion inspired me on a daily basis to keep working towards a better understanding of this field. I dedicate my thesis to my family, from whom I always receive unwavering support.

# Chapter 1

## Interpretable Network from Rate Reduction

### 1.1 Principle of Maximal Coding Rate Reduction

Given a set of training data  $\{\mathbf{x}_i\}$  and their corresponding labels  $\{\mathbf{y}_i\}$ , classical deep learning aims to learn a nonlinear mapping  $h(\cdot) : \mathbf{x} \rightarrow \mathbf{y}$ , implemented as a series of simple linear and nonlinear maps, that minimizes the cross-entropy loss. One popular way to interpret the role of multiple layers is to consider the output of each intermediate layer as a latent representation space. Then, the beginning layers aim to learn a latent representation  $\mathbf{z} = f(\mathbf{x}, \theta) \in \mathbb{R}^d$  that best facilitates the later layers  $\mathbf{y} = g(\mathbf{z})$  for the downstream classification task.

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

As a concrete example, in image recognition tasks,  $f(\cdot)$  is a convolutional backbone that encodes an image  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  into a vector representation  $\mathbf{z} = f(\mathbf{x}, \theta) \in \mathbb{R}^d$  and  $g(\mathbf{z}) = \mathbf{w} \cdot \mathbf{z}$  is a linear classifier where  $\mathbf{w} \in \mathbb{R}^{k \times d}$  and  $k$  is the number of classes. Recent work [15] shows that this direct label fitting leads to a phenomena called *neural collapse* in deep networks, where within-class variability and structural information are completely suppressed. Namely, at the final hidden layer, the variance of each latent class distribution often converges to 0, such that for each class, the training inputs map to a single point. Therefore, it is unclear to what extent the feature representation captures any intrinsic structure of the data.

To address the aforementioned problem, a recent work by Yu *et al.* [25] presented a framework for learning useful and geometrically meaningful representations by maximizing the coding rate reduction (i.e., MCR<sup>2</sup>). Given  $m$  training samples of  $d$  dimension  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{d \times m}$  that belong to  $k$  classes, let  $\mathbf{Z} = [f(\mathbf{x}_1, \theta), \dots, f(\mathbf{x}_m, \theta)] \in \mathbb{R}^{d \times m}$  be the latent representation. Let  $\mathbf{\Pi} = \{\mathbf{\Pi}^j\}_{j=1}^k$  be the membership of the data in the  $k$  classes, where each  $\mathbf{\Pi}^j \in \mathbb{R}^{m \times m}$  is a diagonal matrix such that  $\mathbf{\Pi}^j(i, i)$  is the probability of  $\mathbf{x}_i$  belonging to class  $j$ . Then, MCR<sup>2</sup> aims to learn a feature representation  $\mathbf{Z}$  by maximizing

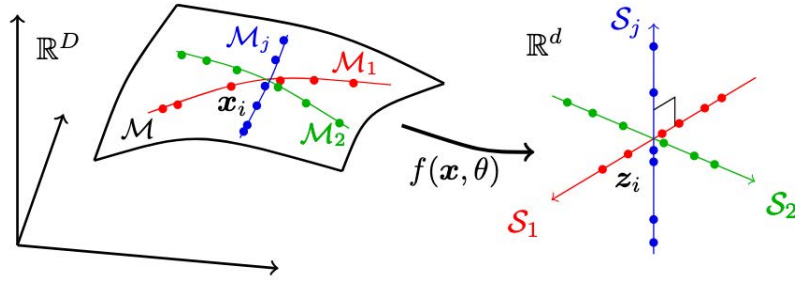


Figure 1.1:  $MCR^2$  learns a map  $f(\mathbf{x}, \theta)$  such that the class representation  $\mathbf{Z}_j = f(\mathbf{X}_j, \theta)$  lies on maximally uncorrelated subspaces  $\{S_j\}$ .

the following *rate reduction*:

$$\Delta R(\mathbf{Z}) = R(\mathbf{Z}) - R_c(\mathbf{Z}, \mathbf{\Pi}), \quad (1.1)$$

where

$$R(\mathbf{Z}) = \frac{1}{2} \log \det (\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^\top), \quad \text{and} \quad (1.2)$$

$$R_c(\mathbf{Z}, \mathbf{\Pi}) = \sum_{j=1}^k \frac{\gamma_j}{2} \log \det (\mathbf{I} + \alpha_j \mathbf{Z} \mathbf{\Pi}^j \mathbf{Z}^\top). \quad (1.3)$$

subject to the constraint that  $\mathbf{Z}$  is properly normalized, *i.e.*, have the Frobenius norm of class features  $\mathbf{Z}^j = \mathbf{Z} \mathbf{\Pi}^j$  scale with the number of samples in class  $j$ :  $\|\mathbf{Z}^j\|_F^2 = m_j = \text{tr}(\mathbf{\Pi}^j)$ . Above, we denote  $\alpha = d/(m\epsilon^2)$ ,  $\alpha_j = d/(\text{tr}(\mathbf{\Pi}^j)\epsilon^2)$ ,  $\gamma_j = \text{tr}(\mathbf{\Pi}^j)/m$ , and  $\epsilon > 0$  is a prescribed quantization error.

This coding rate measure utilizes local  $\epsilon$ -ball packing to estimate the coding rate of the latent distribution from finite samples [13].  $R(\mathbf{Z})$ , known as the *expansion* term, represents the total coding length of all features  $\mathbf{Z}$  while  $R_c(\mathbf{Z}, \mathbf{\Pi})$ , named *compression* term, measures the sum of coding lengths of each latent class distribution. They are called expansion and compression terms respectively, since to maximize  $\Delta R$ , the first coding rate term is maximized and the second coding rate term is minimized.

The features learned by  $MCR^2$  have precise statistical and geometric interpretations formalized in the following theorem.

**Theorem 1.1.1.** (Theorem 2.1 in [25]) Suppose  $\mathbf{Z}^* = \mathbf{Z}_1^* \cup \dots \cup \mathbf{Z}_k^*$  is the optimal solution that maximizes the rate reduction (1.1) with the extra constraint that  $\text{rank}(\mathbf{Z}_j) \leq d_j$ ,  $\forall j \in \{1, \dots, k\}$ . As long as the ambient space is adequately large ( $d \geq \sum_{j=1}^k d_j$ ), we have:

1. Between-class Discriminative: The subspaces are all orthogonal to each other, *i.e.*  $(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* = 0$  for  $i \neq j$ .

2. Maximally Diverse Representation: As long as the coding precision is adequately high, i.e.,  $\varepsilon^4 < \min_j \left\{ \frac{m_j d^2}{m d_j^2} \right\}$ , 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 [25], it is demonstrated empirically that maximizing  $\Delta R(\mathbf{Z})$  enforces each latent class distribution to be a low-dimensional subspace-like distribution of approximately  $\frac{d}{k}$  dimension with class balance *i.e.* among all such discriminative representations, it prefers the one that spans the whole ambient space. The data points are distributed *isotropically* in each subspace except for possibly one dimension. In addition, these class distributions are orthogonal to each other (See Figure 1.1). Thus, by maximizing the coding rate difference, the features become *between-class* discriminative, whilst maintaining *intra-class* diversity. We refer interested readers to [25] for detailed proofs and empirical results.

## 1.2 Rate Reduction Network (ReduNet)

While an existing neural network architecture (such as ResNet) can be used for feature learning with MCR<sup>2</sup>, a follow-up work [2] showed that a novel architecture can be explicitly constructed without back-propagation via emulating the projected gradient ascent scheme for maximizing  $\Delta R(\mathbf{Z})$ . This produces a “white box” network, called ReduNet, which has precise statistical and geometric interpretations. We review the construction of ReduNet as follows.

Let  $\mathbf{Z}$  be initialized as the training data, i.e.,  $\mathbf{Z}_0 = \mathbf{X}$ . Then, the projected gradient ascent step for optimizing the rate reduction  $\Delta R(\mathbf{Z})$  in (1.1) is given by

$$\begin{aligned} \mathbf{Z}_{\ell+1} &\propto \mathbf{Z}_\ell + \eta \left( \frac{\partial \Delta R}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} \right) \\ &= \mathbf{Z}_\ell + \eta \left( \mathbf{E}_\ell \mathbf{Z}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{Z}_\ell^j \right) \\ \text{s.t. } \quad &\|\mathbf{Z}_{\ell+1}^j\|_F^2 = \text{tr}(\mathbf{\Pi}^j) = m_j \quad \forall j \in \{1, \dots, k\}, \end{aligned} \tag{1.4}$$

where we use  $\mathbf{Z}_\ell^j = \mathbf{Z}_\ell \mathbf{\Pi}^j \in \mathbb{R}^{d \times m}$  to denote the feature matrix associated with the  $j$ -th class at the  $\ell$ -th iteration, and  $\eta > 0$  is the learning rate. The matrices  $\mathbf{E}_\ell$  and  $\mathbf{C}_\ell^j$  are obtained by evaluating the derivative  $\frac{\partial \Delta R}{\partial \mathbf{Z}}$  at  $\mathbf{Z}_\ell$ , given by

$$\mathbf{E}_\ell = \alpha \left( \mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^\top \right)^{-1}, \tag{1.5}$$

$$\mathbf{C}_\ell^j = \alpha_j \left( \mathbf{I} + \alpha_j \mathbf{Z}_\ell^j \mathbf{Z}_\ell^{j\top} \right)^{-1}. \tag{1.6}$$

Observe that  $\mathbf{E}_\ell \in \mathbb{R}^{d \times d}$  is applied to all features  $\mathbf{Z}_\ell$  and it expands the coding length of the entire data. Meanwhile,  $\mathbf{C}_\ell^j \in \mathbb{R}^{d \times d}$  is applied to features from class  $j$ , i.e.,  $\mathbf{Z}_\ell^j$ , and it compresses the coding lengths of the  $j$ -th class.

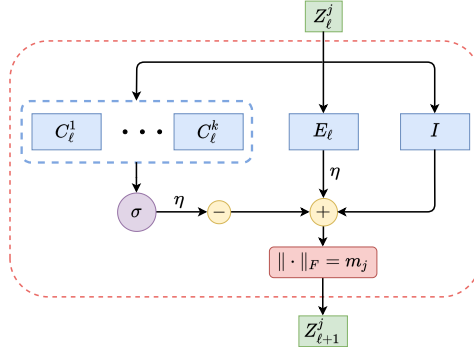


Figure 1.2: ReduNet Architecture in which we here adopt a slightly different normalization than [2] but is more suitable for the incremental learning as we will see in our derivation.

Each projected gradient step can be interpreted as one layer of a neural network composed of matrix multiplication and subtraction operators, with  $\mathbf{E}_\ell$  and  $\mathbf{C}_\ell^j$  being parameters associated with the  $\ell$ -th layer computed using training data (See Figure 1.2). Then, given a test sample  $\mathbf{x} \in \mathbb{R}^d$ , its feature can be computed by setting  $\mathbf{z}_0 = \mathbf{x}$  and iteratively carrying out the following incremental transform

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_\ell + \eta \left( \mathbf{E}_\ell \mathbf{z}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{z}_\ell \pi^j(\mathbf{z}_\ell) \right). \quad (1.7)$$

Notice that the increment depends on  $\pi^j(\mathbf{z}_\ell)$ , the membership of the feature  $\mathbf{z}_\ell$ , which is unknown for the test data. Therefore, [2] presented a method that replaces  $\pi^j(\mathbf{z}_\ell)$  in (1.7) by the following estimated membership

$$\hat{\pi}_\ell^j(\mathbf{z}) = \frac{\exp(-\lambda k \|\mathbf{C}_\ell^j \mathbf{z}\|)}{\sum_{j=1}^k \exp(-\lambda k \|\mathbf{C}_\ell^j \mathbf{z}\|)} \in [0, 1], \quad (1.8)$$

where  $\lambda > 0$  is a confidence parameter. This leads to a nonlinear operator  $\sigma(\mathbf{C}_\ell^1 \mathbf{z}_\ell, \dots, \mathbf{C}_\ell^k \mathbf{z}_\ell) \doteq \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{z}_\ell \hat{\pi}_\ell^j$  that, after being plugged into (1.7), produces a nonlinear layer as summarized in Figure 1.2. Stacking multiple such layers produces a multi-layer neural network for extracting discriminative features. Then, a nearest subspace classifier as we will discuss in Section 2.4 can classify the data. Note that each layer is interpretable and computed explicitly.

Note that when membership is known, ReduNet is simply a linear transformation for each class data. Additionally, note that  $\mathbf{E}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j$  cannot move the data  $\mathbf{Z}_\ell$  towards a direction that is not spanned by  $\mathbf{Z}_\ell$ . Thus, the final representation achieved by each class does not necessarily span all  $\frac{d}{k}$  dimensions like the globally optimal configuration (Theorem 1.1.1) if the initial data lies on a lower dimensional subspace.

## Chapter 2

# ReduNet for Interpretable Incremental Learning

The work in this chapter was published at the Conference on Computer Vision and Pattern Recognition (CVPR) 2021 by the author, Ziyang Wu, Chong You, and Yi Ma [23].

### 2.1 Introduction

Humans are capable of acquiring new information continuously while retaining previously obtained knowledge. This seemingly natural capability, however, is extremely difficult for deep neural networks (DNNs) to achieve. *Incremental learning* (IL), also known as continual learning or life-long learning, thus studies the design of machine learning systems that can assimilate new information without forgetting past knowledge.

In incremental learning, models go through rounds of training sessions to accumulate knowledge for a particular objective (*e.g.* classification). Specifically, under *class* incremental learning (class-IL), an agent has access to training data from a subset of the classes, known as a *task*, at each training session and is evaluated on all seen classes at inference time. The overarching goal is to precisely fine-tune a model trained on previously seen tasks to additionally classify new classes of data. However, due to the absence of old data, such models often suffer from *catastrophic forgetting* [14], which refers to a drastic drop in performance after training incrementally on different tasks.

In the last few years, a flurry of continual learning algorithms have been proposed for DNNs, aiming to alleviate the effect of catastrophic forgetting. These methods can be roughly partitioned into three categories: 1) regularization-based methods that often involve knowledge distillation [12, 6, 19, 26], 2) exemplar-based methods that keep partial copies of data from previously learned tasks [16, 1, 22], and 3) modified architectures that attempt to utilize network components specialized for different tasks [17, 19, 11]. In practice, these algorithms exhibit varying performance across different datasets and their ability to mitigate catastrophic forgetting is inadequate. Factors including domain shift [18] across tasks and

imbalance of new and past classes [22] are part of the reason.

The fundamental roadblock in deep continual learning is that DNNs are trained and optimized in a “black box” fashion. Each model contains millions of mathematical operations and its complexity prevents humans from following the mapping from data input to prediction. Given our current limited understanding of network parameters, it is difficult, if not impossible, to precisely control the parameters of a pre-trained model such that the decision boundary learned fits to new data without losing its understanding of old data.

In this work, we take a drastically different approach to incremental learning. We avoid “black box” architectures entirely, and instead utilize a recently proposed “white box” DNN architecture derived from the principle of rate reduction [2]. Termed ReduNet, each layer of this DNN can be explicitly computed in a forward-propagation fashion and each parameter has precise statistical interpretations. The so-constructed network is intrinsically suitable for incremental learning because the second-order statistics of any previously-seen training data is preserved in the network parameters to be leveraged for future tasks.

We propose a new incremental learning algorithm utilizing ReduNet to demonstrate the power and scalability of designing more interpretable networks for continual learning. Specifically, we prove that a ReduNet trained incrementally can be constructed to be equivalent to one obtained by joint training, where all data, both new and old, is assumed to be available at training time. Finally, we observe that ReduNet performs significantly better on MNIST [9] and CIFAR-10 [7] in comparison to current continual DNN approaches.

## 2.2 Related Work

Since the early success of deep learning in classification tasks such as object recognition, attention has lately shifted to the problem of incremental learning in hopes of designing deep learning systems that are capable of continuously adapting to data from non-stationary and changing distributions.

Incremental learning can refer to different problem settings and most studies focus on three widely accepted scenarios [20]. Most of the earlier works [12, 17, 6, 19] study the *task* incremental (task-IL) setting, where a model, after trained on multiple tasks, must be able to classify on data belonging to all the classes it has seen so far. However, the model is additionally provided a task-ID indicating the task or subset of classes each datapoint belongs to. Models trained under this setting are thus required to distinguish among typically only a small number of classes. Recent works [24, 27] explore the more difficult *class* incremental (class-IL) setting, where task-ID is withheld at inference time. This setting is considerably more difficult since without the task-ID, each datapoint could potentially belong to any of the classes the model has seen so far. The other setting, known as *domain* incremental learning (domain-IL) differs from the previous two settings in that each task consists of all the classes the model needs to learn. Instead, a task-dependent transformation is applied to the data. For example, each task could contain the same training data rotated by differing

degrees and the model must learn to classify images of all possible rotations without access to the task-ID.

Deep continual learning literature from the last few years can be roughly partitioned into three categories as follows:

**Regularization-based** methods usually attempt to preserve some part of the network parameters deemed important for previously learned tasks. Knowledge distillation [4] is a popular technique utilized to preserve knowledge obtained in the past. Learning without Forgetting (LwF) [12], for example, attempts to prevent the model parameters from large drifts during the training of the current task by employing cross-entropy loss regularized by a distillation loss. Alternatively, elastic weight consolidation (EWC) [6] attempts to curtail learning on weights based on their importance to previously seen tasks. This is done by imposing a quadratic penalty term that encourages weights to move along directions with low Fisher information. Schwarz *et al.* [19] later proposed an online variant (oEWC) that reduces the cost of estimating the Fisher information matrix. Similarly, Zenke *et al.* [26] limits the changes of important parameters in the network by using an easy-to-compute surrogate loss during training.

**Exemplar-based** methods typically use a memory buffer to store a small set of data from previous tasks in order to alleviate catastrophic forgetting. The data stored is used along with the data from the current task to jointly train the model. Rebuffi *et al.* [16] proposed iCaRL which uses a herding algorithm to decide which samples from each class to store during each training session. This technique is combined with regularization with a distillation loss to further encourage knowledge retention [16]. A recent work by Wu *et al.* [22] achieved further improvements by correcting the bias towards new classes due to data imbalance, which they empirically show causes degradation in performance for large-scale incremental learning settings. This is accomplished by appending a bias-correction layer at the end of the network. Another increasingly popular approach is to train a generative adversarial network (GAN) [5, 21] on previously seen classes and use the generated synthetic data to facilitate training on future tasks.

**Architecture-based** methods either involve designing specific components in the architecture to retain knowledge of previously seen data or appending new parameters or entire networks when encountering new classes of data. Progressive Neural Network (PNN) [17], for example, instantiates a new network for each task with lateral connection between networks in order to overcome forgetting. This results in the number of networks to grow linearly with respect to the number of tasks as training progresses. Progress & Compress (P & C) [19] utilizes one network component to learn the new task, then distills knowledge into a main component that aggregates knowledge from previously encountered data. Li *et al.* [11], proposes a neural architecture search method that utilizes a separate network to learn whether to reuse, adapt, or add certain building blocks of the main classification network for each task encountered.

Our work studies the more difficult class-IL scenario and does not involve regularization or storing any exemplars. Our method thus can be characterized as an architecture-based approach. However, our method differs with the aforementioned works in several important aspects. First, we use a “white box” architecture that is computed exactly in a feed-forward



manner. Moreover, the network, when trained under class-IL scenario, can be shown to perform equivalently to one obtained from joint training while most existing works [11, 19, 17] based on modified architectures target the less challenging task-IL setting. We discuss the differences in more detail later in Section 2.3, after we have introduced our method properly.

## 2.3 Incremental Learning with ReduNet

In this paper, we tackle the task of *class* incremental learning, formalized as follows. Suppose we have a stream of tasks  $\mathbf{D}^1, \mathbf{D}^2, \dots, \mathbf{D}^t, \dots$ , where each task  $\mathbf{D}^t$  consists of data from  $k_t$  classes, i.e.,  $\mathbf{D}^t = \{\mathbf{X}^{(t-1)k_t+1}, \dots, \mathbf{X}^{t \cdot k_t}\}$  where  $\mathbf{X}^j$  is a set of points in class  $j$ . The classes in different tasks are assumed to be mutually exclusive. Furthermore, it is assumed that the tasks arrive in an online setting, meaning that at timestep  $t$  when data  $\mathbf{D}^t$  arrives, the data associated with old tasks  $\{\mathbf{D}_i, i < t\}$  becomes unavailable. Therefore, the objective is to design a learning system that can adapt the model from the old tasks so as to correctly classify on all tasks hitherto, i.e.,  $\mathbf{D}^1, \dots, \mathbf{D}^t$ . In addition, we assume that we are not given the information on the task a test data belongs to, making this problem significantly more challenging than task-IL.

In this section, we show that ReduNet can perfectly adapt to a new task without forgetting old tasks. Specifically, we present an algorithm to adapt the ReduNet constructed from data  $\{\mathbf{D}_i, i < t\}$  by using only the data in  $\mathbf{D}^t$ , so that the updated ReduNet is *exactly the same* as the ReduNet constructed as if data from all tasks  $\{\mathbf{D}_i, i \leq t\}$  were available.

### Derivation of Incrementally-Trained ReduNet

Without loss of generality, we consider the simple case with two tasks  $t$  and  $t'$  where  $t$  is treated as the old task and  $t'$  is treated as the new task. Assume that  $t$  and  $t'$  contain  $m_t, m_{t'}$  training samples and  $k_t, k_{t'}$  distinct classes, respectively. We denote such training data by  $\mathbf{Z}_{0,t} \in \mathbb{R}^{d \times m_t}$  (for task  $t$ ) and  $\mathbf{Z}_{0,t'} \in \mathbb{R}^{d \times m_{t'}}$  (for task  $t'$ ), and assume that they have been normalized by Frobenius norm as described in (1.4). For ease of notation, we label the classes as  $\{1, \dots, k_t\}$  for task  $t$  and  $\{k_t + 1, \dots, k_t + k_{t'}\}$  for task  $t'$ .

Let  $\Theta_t$  be the ReduNet of depth  $L$  trained on task  $t$  as described in Section 1.2. Given the new task  $\mathbf{Z}_{0,t'}$ , our objective is to train a network  $\Theta_{t \rightarrow t'}$  that adapts  $\Theta_t$  to have good performance for both tasks  $t$  and  $t'$ . Next, we show that a network  $\Theta_{t \rightarrow t'}$  can be constructed from  $\Theta_t$  and  $\mathbf{Z}_{0,t'}$  such that it is *equivalent* to  $\Theta$  obtained from training on  $\mathbf{Z}_0 = [\mathbf{Z}_{0,t} | \mathbf{Z}_{0,t'}] \in \mathbb{R}^{d \times m}$  where  $m = m_t + m_{t'}$ .

To start, consider the initial expansion term  $\mathbf{E}_0 \in \mathbb{R}^{d \times d}$  and compression terms  $\mathbf{C}_0^j$  at layer 0 of the joint network  $\Theta$  given by

$$\begin{aligned} \mathbf{E}_0 &= \alpha (\mathbf{I} + \alpha \mathbf{Z}_0 \mathbf{Z}_0^\top)^{-1} \\ &= \alpha (\mathbf{I} + \alpha (\mathbf{Z}_{0,t} \mathbf{Z}_{0,t}^\top + \mathbf{Z}_{0,t'} \mathbf{Z}_{0,t'}^\top))^{-1}, \end{aligned} \tag{2.1}$$

and

$$\mathbf{C}_0^j = \begin{cases} \alpha_j \left( \mathbf{I} + \alpha_j \mathbf{Z}_{0,t}^j \mathbf{Z}_{0,t}^{j\top} \right)^{-1}, & \text{if } j \leq k_t, \\ \alpha_j \left( \mathbf{I} + \alpha_j \mathbf{Z}_{0,t'}^j \mathbf{Z}_{0,t'}^{j\top} \right)^{-1}, & \text{else} \end{cases} \quad (2.2)$$

where  $\alpha = d/(m\epsilon^2)$  and  $\alpha_j = d/(\text{tr}(\mathbf{\Pi}^j)\epsilon^2)$ .

Note that the term  $\mathbf{Z}_{0,t'}^j \mathbf{Z}_{0,t'}^{j\top}$  can be directly computed from input data  $\mathbf{Z}_{0,t'}^j$ . On the other hand, the term  $\mathbf{Z}_{0,t}^j \mathbf{Z}_{0,t}^{j\top}$  cannot be directly computed from input data as  $\mathbf{Z}_{0,t}^j$  is from the old task, which is no longer available under the IL setup. Our key observation is that  $\mathbf{Z}_{0,t}^j \mathbf{Z}_{0,t}^{j\top}$  can be computed from the network  $\Theta_t$ . Specifically, by denoting the compression matrices of  $\Theta_t$  as  $\{\mathbf{C}_{\ell,t}^j\}$  for  $\ell \in \{0, \dots, L-1\}$ , we have

$$\mathbf{Z}_{0,t}^j \mathbf{Z}_{0,t}^{j\top} = \left( (\mathbf{C}_{0,t}^j / \alpha_j)^{-1} - \mathbf{I} \right) / \alpha_j. \quad (2.3)$$

Next, we show by induction that one can recursively compute  $\mathbf{E}_\ell$  and  $\{\mathbf{C}_\ell^j\}$  of  $\Theta$  for  $\ell > 0$  from (2.1) and (2.2). To construct layer 1 of  $\Theta$ , we observe that the output features of class  $j$  at layer 0 *before* normalization is as follows.

$$\mathbf{P}_0^j = \left( \mathbf{Z}_0 + \eta \mathbf{E}_0 \mathbf{Z}_0 - \eta \sum_{i=1}^k \gamma_i \mathbf{C}_0^i \mathbf{Z}_0 \right) \mathbf{\Pi}^j \quad (2.4)$$

$$= \mathbf{Z}_0^j + \eta \mathbf{E}_0 \mathbf{Z}_0^j - \eta \gamma_j \mathbf{C}_0^j \mathbf{Z}_0^j \quad (2.5)$$

$$= \underbrace{\left( \mathbf{I} + \eta \mathbf{E}_0 - \eta \gamma_j \mathbf{C}_0^j \right)}_{\mathbf{L}_0^j \in \mathbb{R}^{d \times d}} \mathbf{Z}_0^j. \quad (2.6)$$

Notice the term  $\mathbf{L}_0^j$  only depends on quantities already obtained at layer 0. To compute  $\mathbf{E}_1$  and  $\mathbf{C}_1^j$ , we need the covariance matrix of  $\mathbf{P}_0^j$ , which we observe to be

$$\mathbf{T}_1^j = \mathbf{P}_0^j \mathbf{P}_0^{j\top} = \mathbf{L}_0^j \mathbf{Z}_0^j \mathbf{Z}_0^{j\top} \mathbf{L}_0^{j\top}. \quad (2.7)$$

Notice that  $\mathbf{T}_1^j$  can be expressed with known quantities of  $\mathbf{L}_0^j$  and  $\mathbf{Z}_{0,t}^j \mathbf{Z}_{0,t}^{j\top}$  if  $j \leq k_t$  or  $\mathbf{Z}_{0,t'}^j \mathbf{Z}_{0,t'}^{j\top}$  if  $j > k_t$ . The remaining step would be to re-scale  $\mathbf{T}_1^j$  as the updated representation  $\mathbf{P}_0^j$  needs to be normalized to get  $\mathbf{Z}_1^j$ . Recall that we adopt the normalization scheme that imposes the Frobenius norm of each class  $\mathbf{Z}_j$  to scale with  $m_j$ :

$$\|\mathbf{Z}_1^j\|_F^2 = m_j \iff \text{tr}(\mathbf{Z}_1^j \mathbf{Z}_1^{j\top}) = m_j. \quad (2.8)$$

The re-scaling factor is then easy to calculate:

$$\mathbf{Z}_1^j \mathbf{Z}_1^{j\top} = \frac{m_j}{\text{tr}(\mathbf{T}_1^j)} \mathbf{T}_1^j. \quad (2.9)$$

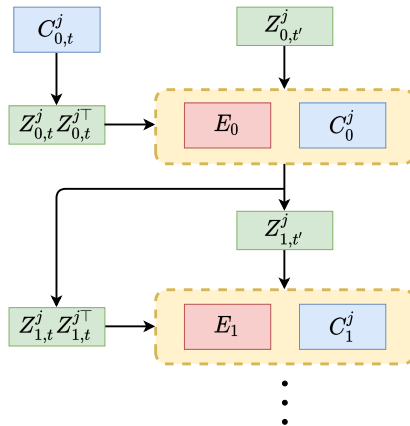


Figure 2.1: The joint network can be derived using simply  $\mathbf{Z}_{0,t}\mathbf{Z}_{0,t}^\top$ . We do not need the task  $t$  data  $\mathbf{Z}_{0,t}$  directly.

From above, we see that we can obtain the correct value of the covariance matrix  $\mathbf{Z}_1^j\mathbf{Z}_1^{j\top}$ , from which we can derive  $\mathbf{E}_1$  and  $\mathbf{C}_1^j$  for layer 1 of the joint network  $\Theta$  and obtain  $\mathbf{Z}_{2,t}^j$ . With these values, we can compute  $\mathbf{T}_2^j$ .

By the same logic, we can recursively update  $\mathbf{E}_\ell$  and  $\mathbf{C}_\ell^j$  for all  $\ell > 1$ . Specifically, once we have obtained  $\mathbf{Z}_{\ell-1}^j\mathbf{Z}_{\ell-1}^{j\top}$  and  $\mathbf{L}_{\ell-1}^j$ , it is straightforward to compute  $\mathbf{T}_\ell^j = \mathbf{L}_{\ell-1}^j\mathbf{Z}_{\ell-1}^j\mathbf{Z}_{\ell-1}^{j\top}\mathbf{L}_{\ell-1}^{j\top}$  and therefore obtain

$$\mathbf{Z}_\ell^j\mathbf{Z}_\ell^{j\top} = \frac{m_j}{\text{tr}(\mathbf{T}_\ell^j)}\mathbf{T}_\ell^j. \quad (2.10)$$

Note that we never need to access  $\mathbf{Z}_{0,t} \in \mathbb{R}^{d \times m_t}$  directly. Instead, we iteratively update the covariance matrix  $\mathbf{Z}_{\ell-1}^j\mathbf{Z}_{\ell-1}^{j\top} \in \mathbb{R}^{d \times d}$  for each class  $j$  using the procedure described. This concludes our induction and Algorithm 1 describes the entire training process for incremental learning on two tasks. The procedure is illustrated in Figure 2.1. This procedure can be naturally extended to settings with more than two tasks.

## Comparison to Existing Methods

Incremental learning with ReduNet offers several nice properties: 1) Each parameter of the network has an explicit purpose, computed precisely to emulate the gradient ascent on the feature representation. 2) It does not require a memory buffer which is often needed in many state-of-the-art methods [16, 22, 1]. 3) It can be proven to behave like a network reconstructed from joint training, thus eliminating the problem of catastrophic forgetting.

Note that many existing works without relying on exemplars [12, 6, 19, 26] regularize the original weights of the model at each training session, effectively freezing certain parts of the

---

**Algorithm 1** Incremental Learning with ReduNet
 

---

- 1: **Input:** Network  $\Theta_t$  with parameters  $\mathbf{E}_{\ell,t}$  and  $\{\mathbf{C}_{\ell,t}^j\}$ , data  $\mathbf{Z}_{0,t'}^j \forall j \in \{k_t+1, \dots, k_t+k_{t'}\}$ .
  - 2: Compute  $\Sigma_{0,t}^j = \mathbf{Z}_{0,t}^j \mathbf{Z}_{0,t}^{j\top}$ ,  $\forall j \in \{1, \dots, k_t\}$  by (2.3).
  - 3: **for**  $\ell = 0, 1, 2, \dots, L-1$  **do**
  - 4:      $\Sigma_{\ell,t} = \sum_{j=1}^{k_t} \Sigma_{\ell,t}^j$ ,
  - 5:      $\Sigma_{\ell,t'} = \sum_{j=k_t+1}^{k_t+k_{t'}} \mathbf{Z}_{\ell,t'}^j \mathbf{Z}_{\ell,t'}^{j\top}$ ,
  - 6:      $\mathbf{E}_\ell = \alpha (\mathbf{I} + \alpha (\Sigma_{\ell,t} + \Sigma_{\ell,t'}))^{-1}$ ,
  - 7:      $\mathbf{L}_\ell^j = \mathbf{I} + \eta \mathbf{E}_\ell - \eta \gamma_j \mathbf{C}_\ell^j \quad \forall j \in \{1, \dots, k_t\}$ ,
  - 8:     **for**  $j = 1, 2, \dots, k_t$  **do**
  - 9:          $\mathbf{C}_\ell^j = \alpha_j (\mathbf{I} + \alpha_j \Sigma_\ell^j)^{-1}$ ,
  - 10:          $\mathbf{T}_{\ell+1,t}^j = \mathbf{L}_\ell^j \Sigma_{\ell,t}^j \mathbf{L}_\ell^{j\top}$ ,
  - 11:          $\Sigma_{\ell+1,t}^j = \frac{m_j}{\text{tr}(\mathbf{T}_{\ell+1,t}^j)} \mathbf{T}_{\ell+1,t}^j$ .
  - 12:     **end for**
  - 13:      $\mathbf{Z}_{\ell+1,t'} \propto \mathbf{Z}_{\ell,t'} + \eta \mathbf{E}_\ell \mathbf{Z}_{\ell,t'} - \eta \sum_{i=k_t+1}^{k_t+k_{t'}} \gamma_i \mathbf{C}_\ell^i \mathbf{Z}_{\ell,t'}^i$
  - 14:     s.t.  $\|\mathbf{Z}_{\ell+1,t'}^j\|_F^2 = m_j$ .
  - 15: **end for**
  - 16: **Output:** Network  $\Theta$  with parameters  $\mathbf{E}_\ell$  and  $\{\mathbf{C}_\ell^j\}$ .
- 

network. Different tasks, however, tend to depend on different parts of the network, which eventually leads to conflicts on which parameters to regularize as the number of tasks to learn increases. These methods, as we see later in Figure 2.2, empirically perform sub-optimally in the class-IL setting. This in fact reveals the fundamental limitation that underlies in many incremental learning methods: a lack of understanding of how individual weights impact the learned representation of data points. ReduNet, on the other hand, sidesteps this problem by utilizing a fully interpretable architecture.

One notable property of ReduNet, at its current form, is that its width grows linearly with the number of classes as a new compression term  $\mathbf{C}_\ell^j$  is appended to each layer whenever we see a new class. On the surface, this makes ReduNet similar to some architecture-based methods [11, 17] that dynamically expand the capacity of the network. However, there exists a major difference. ReduNet is naturally suited for class-IL scenario, whilst the aforementioned works do not address class-IL directly. Instead they only directly address task-IL, which they accomplish by optimizing a sub-network *per task*. These networks, which are designed to accomplish each task individually, fail to properly share information between the sub-networks to discriminate between classes of different tasks. ReduNet accomplishes class-IL by not only appending the class compression terms  $\mathbf{C}_\ell^j$  to the network, but also modifying the expansion term  $\mathbf{E}_\ell$  to share information about classes of all previous tasks.

For class-IL, such methods that also append new parameters to the architecture fail to completely address the problem of catastrophic forgetting. One can see why with a simple

example. Consider an ensemble learning technique where for each class  $j$ , we train an all-versus-one model that predicts whether a data point belongs to class  $j$  or not. At each task, we can feed the available data points into each model, labeled as 1 if it belongs in that class or 0 otherwise. However, by optimizing such “black box” models by back-propagation, we again arrive at the problem of catastrophic forgetting. Specifically, the model only sees training points of its own class only for one task or training session. For the remaining tasks, all data points that it must train will be of label 0, which prevents standard gradient descent from correctly learning the desired all-versus-one decision boundary, and there is no clear way to precisely address this optimization problem.

Although it is natural to expect the network to expand as the number of classes increases, it remains interesting to see if the growth of certain variations of the ReduNet can be sublinear instead of linear in the number of classes.

## 2.4 Experiments

We evaluate the proposed method on MNIST and CIFAR-10 datasets in a class-IL scenario and compare the results with existing methods. In short, for both MNIST and CIFAR-10, the 10 classes are split into 5 incremental batches or tasks of 2 classes each. After training on each task, we evaluate the model’s performance on test data from all classes the model has seen so far. The same setting is applied to all other methods we compared to.

### Datasets

We compare the incremental learning performance of ReduNet on the following two standard datasets.

**MNIST** [10]. MNIST contains 70,000 greyscale images of handwritten digits 0-9, where each image is of size  $28 \times 28$ . The dataset is split into training and testing sets, where the training set contains 60,000 images and the testing dataset contains 10,000 images.

**CIFAR-10** [8]. CIFAR-10 contains 60,000 RGB images of 10 object classes, where each image is of size  $32 \times 32$ . Each class has 5,000 training images and 1,000 testing images. We normalize the input data by dividing the pixel values by 255, and subtracting the mean image of the training set.

### Implementation Details

We implement ReduNet for each training dataset in the following manner.

**ReduNet on MNIST.** To construct a ReduNet on MNIST, we first flatten the input image and represent it by a vector of dimension 784. Then, with a precision  $\epsilon = 0.5$  in the MCR<sup>2</sup> objective (1.1), we apply 200 iterations of projected gradient iterations to compute  $\mathbf{E}_\ell$  and  $\mathbf{C}_\ell^j$  matrices for each iteration  $\ell$ . The learning rate is set to  $\eta = 0.5 \times 0.933^\ell$  at the  $\ell$ -th iteration. These matrices are the parameters of the constructed ReduNet. Given a test

data, its feature can be extracted with the incremental transform in (1.7) with estimated labels computed as in (1.8) with parameter  $\lambda = 1$ . At each training session, we update the ReduNet by the procedure described in Algorithm 1.

We note that hyper-parameter tuning in ReduNet does not require a training/validation splitting as in regular supervised learning methods. The hyper-parameters described above for ReduNet are chosen based on the training data. This is achieved by evaluating the estimated label through (1.8) on the training data, and comparing such labels with ground truth labels. Then, the model parameter  $\epsilon$ , learning rate  $\eta$  and the softmax confidence parameter  $\lambda$  are chosen as those that gives the highest accuracy with the estimated labels (at the final layer).

**ReduNet on CIFAR-10.** We apply 5 random Gaussian kernels with stride 1, size  $3 \times 3$  on the input RGB images.<sup>1</sup> This lifts each image to a multi-channel signal of size  $32 \times 32 \times 5$ , which is subsequently flattened to be a  $\mathbb{R}^{5,120}$  dimensional vector. Subsequently, we construct a 50-layer ReduNet with all other hyper-parameters the same as those for MNIST. All hyperparameters stated above, including the depth of the network, were chosen such that the  $\Delta R$  loss has sufficiently converged.

**Comparing Methods.** We compare our approach to the following state of the art algorithms: iCaRL [16], LwF [12], oEWC [19], SI [26] and DER [1]. For these algorithms, we utilize the same benchmark and training protocol as Buzzega *et al.* [1]. For MNIST, we employ a fully-connected network with two hidden layers comprised of 100 ReLU units. For CIFAR-10, we rely on ResNet18 without pre-training [3]. All the networks were trained by stochastic gradient descent. For MNIST, we train on one epoch per task. For CIFAR-10, we train on 100 epochs per task. The number of epochs were chosen based on the complexity of the dataset. For each algorithm, batch size, learning rate, and specific hyperparameters for each algorithm were selected by performing a grid-search using 10% of the training data as a validation set and selecting the hyperparameter that achieves the highest final accuracy. The optimal hyperparameters utilized for the benchmark experiments are reported in [1].

The performance of state of the art algorithms utilizing a replay buffer highly depends on the number of exemplars, or samples from previous tasks, it is allowed to retain. We test on two exemplar-based algorithm, iCaRL and DER. For both MNIST and CIFAR-10, we set the total number of exemplars to 200.

## Nearest Subspace Classification

By the principle of maximal rate reduction, the ReduNet  $f(\mathbf{X}, \theta)$  extracts features such that each class lies in a low-dimensional linear subspace and different subspaces are orthogonal. As suggested by the original MCR<sup>2</sup> work [2], we utilize a nearest subspace classifier to classify the test data featurized to maximize  $\Delta R$ . Given a test sample  $\mathbf{z}_{test} = f(\mathbf{x}_{test}, \theta)$ , the label

---

<sup>1</sup>This choice is limited by our current computational resources. Although this choice is not adequate to achieve top classification performance, it is adequate to verify the advantages of our method in the incremental setting.

Algorithm	MNIST					CIFAR-10				
	Task 1	Task 2	Task 3	Task 4	Task 5	Task 1	Task 2	Task 3	Task 4	Task 5
LwF	0.999	0.009	0.0	0.0	0.0	0.979	0.0	0.0	0.0	0.0
oEWC	<b>1.0</b>	0.004	0.0	0.0	0.0	0.981	0.0	0.0	0.0	0.0
SI	0.997	0.004	0.001	0.0	0.0	<b>0.989</b>	0.0	0.0	0.0	0.0
iCaRL (200 Exemplars)	0.999	0.806	0.708	0.612	0.596	0.964	0.720	0.427	0.362	0.313
DER (200 Exemplars)	0.999	0.967	0.941	0.883	0.735	0.985	<b>0.816</b>	0.608	0.404	0.292
ReduNet(Ours)	0.999	<b>0.994</b>	<b>0.993</b>	<b>0.989</b>	<b>0.987</b>	0.875	0.754	<b>0.714</b>	<b>0.642</b>	<b>0.562</b>
Upper Bound (UB)	0.999	0.995	0.990	0.988	0.982	0.989	0.971	0.957	0.963	0.920

Table 2.1: Test Accuracy (%) on Task 1 After Each Training Session on MNIST and CIFAR-10.

predicted by a nearest subspace classifier is

$$y = \arg \min_{y \in \{1, \dots, k\}} \|(\mathbf{I} - \mathbf{U}^y \mathbf{U}^{y\top}) \mathbf{z}_{test}\|_2^2, \quad (2.11)$$

where  $\mathbf{U}^y$  is a matrix containing the top  $x$  principal components of the covariance of the training data passed through ReduNet, *i.e.*  $\mathbf{Z}_{train} \mathbf{Z}_{train}^\top$  for  $\mathbf{Z}_{train} = f(\mathbf{X}_{train}, \theta)$ .

Since we do not have access to  $\mathbf{Z}_{train}$  during evaluation, we instead collect the  $\mathbf{C}^j$  matrices at the very last layer  $L$  and extract the covariance matrix  $\mathbf{\Sigma}_L^j$  to be further processed by SVD. For MNIST, we utilize the top 28 principal components. For CIFAR-10, we utilize the top 15 principal components.

## Results and Analysis

In this section, we evaluate the class-IL performance of incremental ReduNet against three regularization-based methods (oEWC, SI, LwF) and two replay-based methods leveraging 200 exemplars (iCaRL, DER) on MNIST and CIFAR-10. We also provide the upper bound (UB) achieved by joint training a model utilizing the same architecture as the baseline methods. After the model is trained on each task, performance is evaluated by computing the accuracy on test data from all classes the model has seen so far. To observe the degree of forgetting, we record the model’s performance on Task 1 after training on each subsequent task. For both MNIST and CIFAR-10, we observe a substantial performance increase by utilizing incremental ReduNet as shown in Figure 2.2. Additionally, we observe ReduNet shows significantly less forgetting (see Table 2.1).

On MNIST, we observe a 3% decay in accuracy across the tasks on ReduNet versus a 20-80% decay on benchmark methods (Figure 2.2). We measure *decay* as the difference in average accuracy between the first and last task. ReduNet retains a classification accuracy of 96%. This is of no surprise since MNIST is relatively linearly separable, allowing second-order information about the data to be sufficient for ReduNet to correctly classify the digits. We observe that even for a very simple task as MNIST, competing continual learning

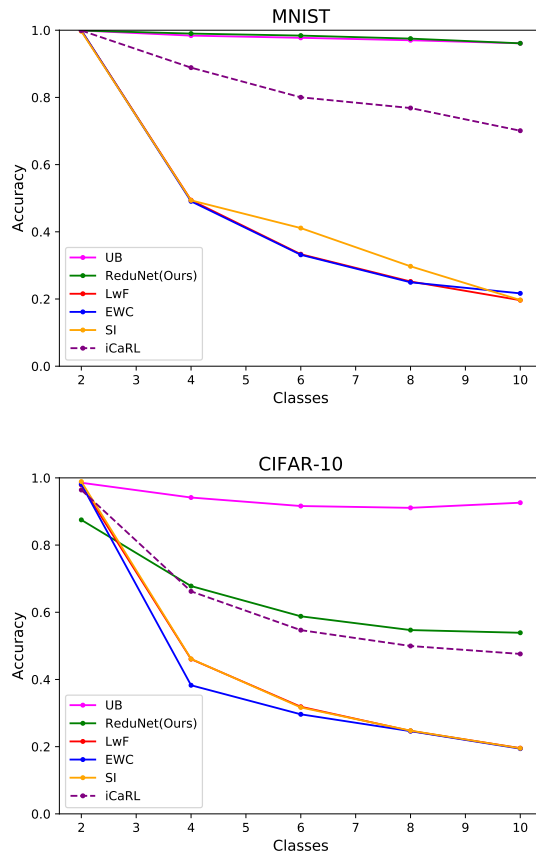


Figure 2.2: Incremental learning results (accuracy) on MNIST and CIFAR-10. Both datasets have 5 incremental batches. We also provide the upper bound (UB) given by joint training a model utilizing the same architecture as the baseline methods. In *solid* lines are regularization-based methods and in *dashed* are exemplar-based methods, which saves 200 samples from previous tasks. Note that the decay in the performance in ReduNet is simply because classification is harder to accomplish with more classes, not because of catastrophic forgetting.

algorithms fail spectacularly due to catastrophic forgetting. In fact, as shown in Table 2.1, models trained by regularization-based methods retain 0% accuracy for classes 0 and 1 after incrementally training on digits 2 to 5 (up to Task 3). The drastic decay in performance of benchmark methods is expected and replicated often in class-IL literature [1, 24]. Note that ReduNet observes no catastrophic forgetting and the decay in performance is due to the fact that classification is increasingly harder to accomplish with more classes.

Similar improvement in performance utilizing ReduNet is also observed on CIFAR-10, a more complex image dataset. We observe a 42-80% decay in accuracy for benchmark methods, whereas incremental ReduNet observes a 34% decrease (in Figure 2.2). The



algorithm that achieves the closest performance to ReduNet is iCaRL and DER, exemplar-based methods that require access to 200 previously observed exemplars. Certainly, as can be seen by the 88% accuracy on Task 1 of CIFAR-10, ReduNet at its current basic form (only using 5 randomly initialized kernels without back-propagation) is not able to reach the same classification accuracy as ResNet-18 for complex image classification tasks. It is thus not surprising that DER, based on more established network architectures, exceeds ReduNet in terms of average accuracy. However, as shown in Table 2.1, ReduNet decays gracefully and significantly outperforms other methods in terms of forgetting, retaining over 55% accuracy on Task 1 versus less than 30% by DER.

We note that ReduNet is currently a slower training framework given its current naive implementation using CuPy. Utilizing a single NVIDIA TITAN V GPU, each task training session took approximately 1500 seconds for MNIST and 9200 seconds for CIFAR10. On the other hand, joint training a model by back-propagation for each task took 23 and 2500 seconds for MNIST and CIFAR10, respectively.

## 2.5 Conclusions and Future Work

We demonstrated through an incremental version of the recently proposed ReduNet, the promise of leveraging interpretable network design for continual learning. The proposed network has shown significant performance increases in both synthetic and complex real data, even without utilizing any fine-tuning with back-propagation. It has clearly shown that if knowledge of past learned tasks are properly utilized, catastrophic forgetting needs not to happen as new tasks continue to be learned.

We purpose of this work is not to push the state of art classification accuracy or efficiency on any single large-scale real-world dataset. Rather we want to use the simplest experiments to show beyond doubt the remarkable effectiveness and great potential of this new framework. Using CIFAR-10 as an example, simply utilizing a relatively small set of 5 random lifting kernels was already sufficient for a decent incremental classification performance. We believe that to achieve better performance for more complex tasks and datasets, judicious design or learning of more convolution kernels would be needed. This leaves plenty of room for further improvements.

This work also opens up a few promising new extensions. As we have mentioned earlier, the current framework requires the width of the network to grow linearly in the number of classes. It would be interesting to see if some of the filters can be shared among old/new classes so that the growth can be sublinear. In the next chapter, we will discuss several approaches to decrease the memory footprint of incremental ReduNet.

To a large extent, the rate reduction gives a unified measure for learning discriminative representations in supervised, semi-supervised, and unsupervised settings. We believe our method can be easily extended to cases when some of the new data do not have class information.

## Chapter 3

# Low Rank Approximation for Efficient ReduNet

As we've observed in the previous chapter, incremental ReduNet tends to have a large memory footprint that increases linearly in the number of classes  $k$ . Namely, given  $d$  dimensional data, each layer requires one to save  $k + 1$   $d \times d$  weight matrices:  $\mathbf{E}_{l,t}, \{\mathbf{C}_{l,t}^j\}$ . In this chapter, we discuss naive, yet effective methods at reducing this memory footprint.

### 3.1 Low rank approximation

Though real world data such as images are often high dimensional, we often expect the data from each class to lie on low dimensional curved manifolds. For example, for MNIST and CIFAR-10, the singular values of the sample covariance matrix of each class is very small across most dimensions (Figure 3.1). Thus, the matrices are approximately low-rank.

We may utilize this understanding to find a substantially compressed version of ReduNet with equal performance. Namely, note that given that the distribution of class  $j$ ,  $\mathbf{Z}_j \mathbf{Z}_j^\top$ , is low rank,  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$  is also low rank where  $\mathbf{C}_j$  is the corresponding compression matrix.

$$\begin{aligned} \frac{1}{\alpha_j} \mathbf{C}_j &= (\mathbf{I} + c_j \mathbf{Z}_j \mathbf{Z}_j^\top)^{-1} = \mathbf{U}_j \text{diag} \left( \frac{1}{1 + c_j \lambda_{ji}} \right) \mathbf{U}_j^\top \\ \Rightarrow \mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j &= \mathbf{U}_j \text{diag} \left( 1 - \frac{1}{1 + c_j \lambda_{ji}} \right) \mathbf{U}_j^\top \end{aligned} \quad (3.1)$$

We may find a low rank approximation of the compression matrix by optimizing over the following objective.

$$\begin{aligned} \min_{\mathbf{D}_j} & \left\| \left( \mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j \right) - \mathbf{D}_j \right\|_F^2 \\ & \text{subject to } \text{rank}(\mathbf{D}_j) < \frac{d}{k} \end{aligned} \quad (3.2)$$

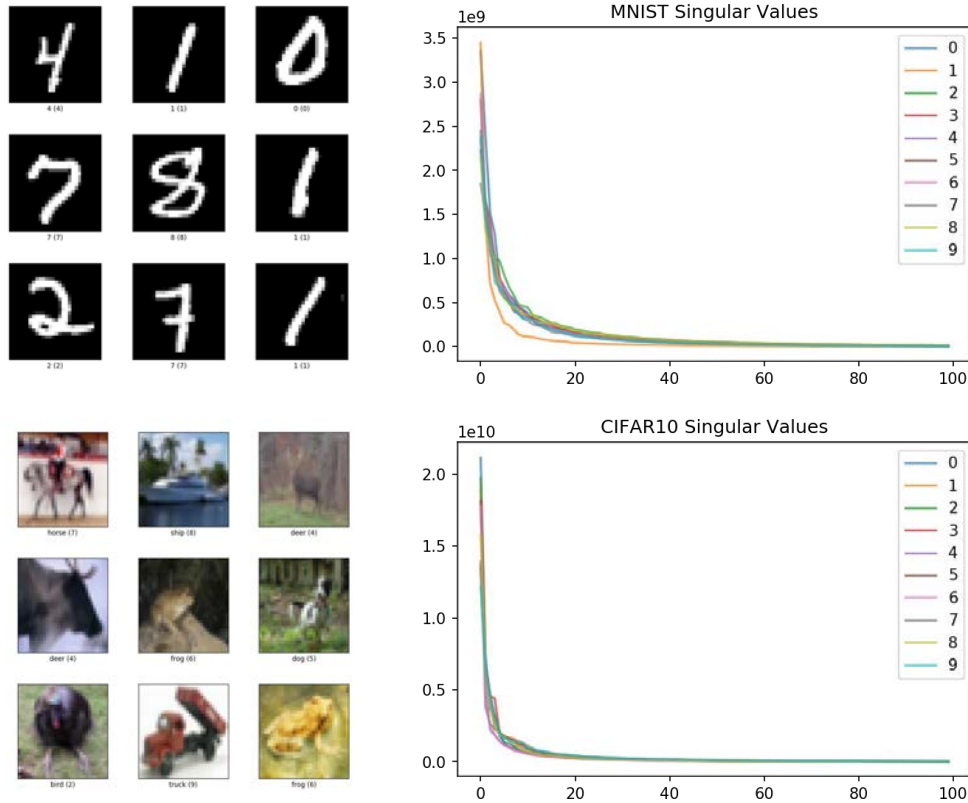


Figure 3.1: The first 100 eigenvalues of each class distribution of MNIST and CIFAR10. MNIST is  $32 \times 32 = 784$  dimensional and CIFAR10 is  $28 \times 28 \times 3 = 3072$  dimensional. Note that eigenvalues nearly decay to 0 by the 20th dimension for both datasets.

where  $d$  is the ambient dimension and  $k$  is the number of classes. The rank constraint is chosen based off the fact that with class balance, the globally optimal configuration  $\mathbf{Z}_j^*$  spans at most a  $\frac{d}{k}$  dimensional subspace.

The above objective has a closed form solution given by  $\mathbf{D}_j = \mathbf{U}_j[:\frac{d}{k}] \boldsymbol{\lambda}_j[:\frac{d}{k}] \mathbf{U}_j[:\frac{d}{k}]^\top$  where  $\mathbf{U}_j[:\frac{d}{k}]$  and  $\boldsymbol{\lambda}_j[:\frac{d}{k}]$  are the top  $\frac{d}{k}$  eigenvectors of  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$  and their corresponding eigenvalues, respectively.

To be more flexible with the rank constraint, we may alternatively approximate  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$  up to a certain tolerance

$$\frac{\left\| \left( \mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j \right) - \mathbf{D}_j \right\|_F^2}{\left\| \left( \mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j \right) \right\|_F^2} \leq \tau \quad (3.3)$$

and greedily choose the the eigenvectors of  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$  until the approximation is below the tolerance. Specifically,

---

**Algorithm 2** Greedy Low Rank Approximation
 

---

```

1: Input:  $\varepsilon, \mathbf{Z}_j$ 
2:  $r = 1, error = \infty$ 
3:  $\mathbf{U}, \boldsymbol{\lambda}, \mathbf{V} = \text{SVD}(\mathbf{Z}_j \mathbf{Z}_j^T)$ 
4: while  $error > \varepsilon$  do
5:    $\mathbf{D}_j = \mathbf{U}[:, r] \text{diag}(\boldsymbol{\lambda}[:, r]) \mathbf{V}[:, r]^T$ 
6:    $error = \frac{\left\| (\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j) - \mathbf{D}_j \right\|_F^2}{\left\| (\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j) \right\|_F^2}$ 
7:    $r += 1$ 
8: end while
9: Output:  $\mathbf{U}[:, r], \boldsymbol{\lambda}[:, r], \mathbf{V}[:, r]$ 

```

---

From  $\mathbf{D}_j$ , we may approximately recover  $\mathbf{C}_j$  by doing  $\mathbf{C}_j = \alpha_j (\mathbf{I} - \mathbf{D}_j)$ . By saving  $r \leq \frac{d}{k}$  top eigenvalues and corresponding eigenvectors of  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$  the memory required to save each layer of ReduNet *at least* reduces to

$$\text{Original: } (k + 1) \times d \times d \implies \text{Compressed: } 2 \times d \times d + d$$

The extra  $d$  term comes from saving the eigenvalues. Note that the memory of this compressed network is now independent on the number of classes.

## Experiments

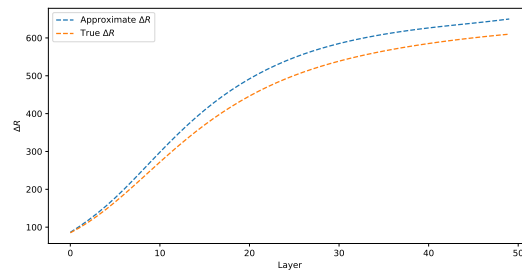
### MNIST

We run 50-layer ReduNet on MNIST with 1000 training samples from each of the 10 classes. We set learning rate to be  $\eta = 0.1$  and ball packing error to be  $\varepsilon^2 = 0.5$ . Each image is  $28 \times 28 = 784$ , so the dimension of  $\mathbf{E}$  and  $\mathbf{C}_j$  is  $784 \times 784$ . We derive the compressed network by saving the  $\frac{d}{k}$ -rank approximation of the compression matrices  $\{\mathbf{C}_l^j\}$  at each layer  $l$ .

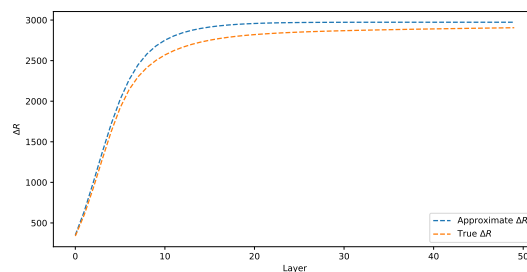
### CIFAR10

For CIFAR10, we run a 50-layer ReduNet with 1000 training samples from each of the 10 classes. We set learning rate to be  $\eta = 0.1$  and ball packing error to be  $\varepsilon^2 = 0.5$ . Each image is  $32 \times 32 \times 3 = 3072$  dimensions, so the dimension of  $\mathbf{E}$  and  $\mathbf{C}_j$  is  $3072 \times 3072$ . We derive the compressed network by saving the  $\frac{d}{k}$ -rank approximation of the compression matrices  $\{\mathbf{C}_l^j\}$  at each layer  $l$ .

We compare the losses and feature representations between original ReduNet and approximate



(a) MNIST



(b) CIFAR10

Figure 3.2: We pass through the training data of MNIST and CIFAR10 through the original ReduNet and compressed ReduNet, and compare the MCR loss  $\Delta R(\mathbf{Z}_l)$  at each layer  $l$ . We observe the loss curve for compressed ReduNet to be better than that of true ReduNet

ReduNet. Recall that the class distributions of MNIST are already quite low rank and separable (Figure 3.1). Since the class distributions of CIFAR10 and MNIST are low dimensional, we observe that the low rank approximations of the compression matrices  $\mathbf{C}_l^j$  are close to the original matrices. Furthermore, compressed ReduNet achieved better loss curves than true ReduNet. By compressing ReduNet up to some tolerance, we were able to reduce the memory footprint of the network even further 3.3.

## 3.2 Block Diagonal Approximation by Eigenspace

We may further compress the model from our knowledge that the eigenvectors of each class distribution  $\mathbf{Z}_j \mathbf{Z}_j^\top$  must be some linear combination of the eigenvectors of the joint distribution  $\mathbf{Z} \mathbf{Z}^\top = \sum_j \mathbf{Z}_j \mathbf{Z}_j^\top$ . In fact, at a critical point  $\mathbf{Z}^*$  of the loss  $\Delta R$ , the eigenvectors of  $\mathbf{Z} \mathbf{Z}^\top$  is precisely equal to the eigenvectors of  $\mathbf{Z}_j \mathbf{Z}_j^\top$ . This statement is formalized in the following theorem.

**Theorem 3.2.1.** *If  $\mathbf{Z}^*$  is a local maxima of  $\Delta R(\mathbf{Z})$ , then each eigenvector with positive*

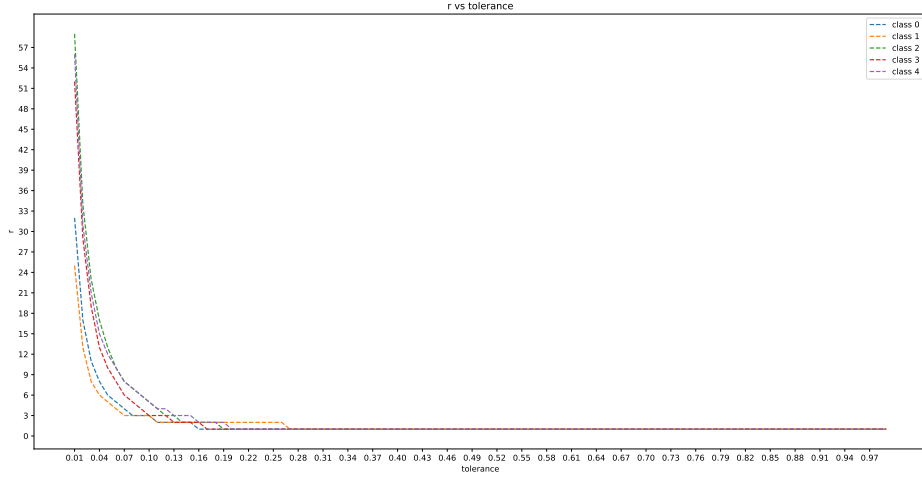


Figure 3.3: In MNIST, we observe that only a few eigenvectors are necessary to approximate  $C_j$  of the first layer for small tolerance values.

corresponding eigenvalue of the class subspace  $\mathbf{Z}_j^* \mathbf{Z}_j^{*\top}$  must be an eigenvector of the whole distribution  $\mathbf{Z}^* \mathbf{Z}^{*\top}$ .

*Proof.* Recall MCR<sup>2</sup> objective:

$$\max \Delta R(\mathbf{Z}) = R(\mathbf{Z}, \epsilon) - R^c(\mathbf{Z}, \epsilon | \mathbf{\Pi}) \quad \text{s.t.} \quad \|\mathbf{Z}_j\|_F^2 = m_j \quad \forall j \in \{1, \dots, k\} \quad (3.4)$$

By first order necessary conditions, note that  $\mathbf{Z}$  is a stationary point if for each class distribution  $\nabla_{\mathbf{Z}_j} \Delta R = a_j \mathbf{Z}_j$  for some constant  $a_j \in \mathbb{R}$ . Recall that the gradient is

$$\frac{\partial \Delta R}{\partial \mathbf{Z}_j} = \left( \mathbf{E} - \gamma_j \mathbf{C}_j^j \right) \mathbf{Z}_j \quad (3.5)$$

So at a stationary point, each sample  $\mathbf{z}_i^j$  in  $\mathbf{Z}_j$  for  $i \in 1, \dots, m_j$  satisfies

$$\left( \mathbf{E} - \gamma_j \mathbf{C}_j \right) \mathbf{z}_i^j = a_j \mathbf{z}_i^j.$$

Furthermore, note that each eigenvector  $\mathbf{w}_l^j$  of  $\mathbf{Z}_j \mathbf{Z}_j^\top$  for  $0 \leq l \leq d$  with corresponding eigenvalue  $\lambda_l^j > 0$ , is some linear combination of class samples

$$\mathbf{w}_l^j = \sum_{i=1}^{m_j} \alpha_i \mathbf{z}_i^j.$$

Thus,

$$\begin{aligned}
& (\mathbf{E} - \gamma_j \mathbf{C}_j) \mathbf{u}_l^j \\
&= \sum_{i=1}^{m_j} \alpha_i (\mathbf{E} - \gamma_j \mathbf{C}_j) \mathbf{z}_i^j \\
&= a_j \sum_{i=1}^{m_j} \alpha_i \mathbf{z}_i^j \\
&= a_j \mathbf{u}_l^j
\end{aligned} \tag{3.6}$$

Since  $\mathbf{C}_j \mathbf{u}_l^j = \alpha_j (1 + c_j \lambda_l^j)^{-1} \mathbf{u}_l^j$  and  $(\mathbf{E} - \gamma_j \mathbf{C}_j) \mathbf{u}_l^j = a_j \mathbf{u}_l^j$ , we may conclude that

$$\begin{aligned}
\mathbf{E} \mathbf{u}_l^j &= a_j \mathbf{u}_l^j + \gamma_j \mathbf{C}_j \mathbf{u}_l^j \\
&= (a_j + \gamma_j \alpha_j (1 + c_j \lambda_l^j)^{-1}) \mathbf{u}_l^j.
\end{aligned} \tag{3.7}$$

Therefore, the eigenvectors (with positive eigenvalue) of each class subspace  $\mathbf{Z}_j \mathbf{Z}_j^\top$  must also be an eigenvector of the whole distribution  $\mathbf{Z} \mathbf{Z}^\top$ .  $\square$

Thus, we may simply utilize the eigenvectors of  $\mathbf{I} - \frac{1}{\alpha} \mathbf{E} = \mathbf{U} (\mathbf{I} - (\mathbf{I} + \alpha \mathbf{\Sigma})^{-1}) \mathbf{U}^\top$  as a dictionary for the model, and use them to approximate the eigenvectors of  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C} = \mathbf{U}_j (\mathbf{I} - (\mathbf{I} + \alpha_j \mathbf{\Sigma}_j)^{-1}) \mathbf{U}_j^\top$ . Additionally, at a critical point, we know that the eigenvectors of  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$  with positive corresponding eigenvalues is equal to approximately  $\frac{d}{k}$  of the eigenvalues of  $\mathbf{I} - \frac{1}{\alpha}$ . Formally, we solve the following objective

$$\begin{aligned}
& \min_{\mathbf{D}_j} \left\| (\mathbf{I} - \mathbf{C}_j) - \mathbf{U} \mathbf{D}_j \mathbf{U}^\top \right\|_F^2 \\
& \text{subject to } \text{rank}(\mathbf{D}_j) < \frac{d}{k}
\end{aligned} \tag{3.8}$$

Similar to low-rank approximation, we may be more flexible with the rank constraint by approximating  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$  up to a certain tolerance

$$\frac{\left\| (\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j) - \mathbf{U} \mathbf{D}_j \mathbf{U}^\top \right\|_F^2}{\left\| (\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j) \right\|_F^2} \leq \tau. \tag{3.9}$$

We may solve the objective up to some tolerance using the following greedy algorithm (Algorithm 3).

Essentially, the algorithm projects each eigenvector  $u \in \mathbf{U}$  of  $\mathbf{I} - \mathbf{E}$  onto each class subspace  $\|(\mathbf{I} - \mathbf{C}_j)u\|_2$  and selects the top  $r \leq \frac{d}{k}$  eigenvectors in  $\mathbf{I} - \frac{1}{\alpha} \mathbf{E}$  with largest

**Algorithm 3** Greedy Block Diagonal Approximation

- 
- 1: **Input:**  $\varepsilon, \mathbf{Z}_j, \mathbf{Z}$
  - 2:  $r = 1, error = \infty$
  - 3:  $\mathbf{U}, \boldsymbol{\lambda}, \mathbf{V} = \text{SVD}(\mathbf{Z}\mathbf{Z}^T)$
  - 4: Sort columns  $u_i$  of  $\mathbf{U}$  by  $\left\| \left( \mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j \right) u_i \right\|_2^2$  in descending order
  - 5: **while**  $error > \varepsilon$  **do**
  - 6:      $\mathbf{D}_j = \min_{\mathbf{D}_j} \left\| \left( \mathbf{I} - \mathbf{C}_j \right) - \mathbf{U}[:, r] \mathbf{D}_j \mathbf{U}[:, r]^T \right\|_F^2 = \mathbf{U}[:, r]^T \left( \mathbf{I} - \mathbf{C}_j \right) \mathbf{U}[:, r]$
  - 7:      $error = \frac{\left\| \left( \mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j \right) - \mathbf{U} \mathbf{D}_j \mathbf{U}^T \right\|_F^2}{\left\| \left( \mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j \right) \right\|_F^2}$
  - 8:      $r += 1$
  - 9: **end while**
  - 10:  $\boldsymbol{\pi}_j =$  Indices corresponding to columns of  $\mathbf{U}$  chosen
  - 11: **Output:**  $\mathbf{D}_j, \boldsymbol{\pi}_j$
- 

projection to approximate  $\mathbf{I} - \frac{1}{\alpha_j} \mathbf{C}_j$ . By saving  $\mathbf{D}_j \in \mathbb{R}^{r \times r}$  and  $\boldsymbol{\pi} \in \mathbb{R}^r$  for  $r \leq \frac{d}{k}$  instead of  $\mathbf{C}_j$ , the memory *at least* reduces to

**Original:**  $(k+1) \times d \times d \implies$  **Compressed:**  $k \times r \times r + d \times d + k \times d = \left(1 + \frac{1}{k}\right) \times d \times d + d$

Additionally, we may reduce the computation needed during forward passing of test data at the label approximation layers. Specifically, once  $\mathbf{D}_j$  is computed, we can save computation by no longer requiring passing each  $\mathbf{C}_j$  with data  $\mathbf{Z}$  every time. Concretely, given  $\mathbf{A} = \mathbf{U}^T \mathbf{Z}$ , we can obtain  $\mathbf{E}\mathbf{Z} = \mathbf{U}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T \mathbf{Z} = \mathbf{U}\boldsymbol{\Sigma}^{-1}\mathbf{A}$ .

Also, we can compute  $\mathbf{C}_j \mathbf{Z} = \left( \mathbf{I} - \left( \mathbf{I} - \mathbf{C}_j \right)_{approx} \right) \mathbf{Z} = \mathbf{Z} - \mathbf{U}[:, \boldsymbol{\pi}_j] \mathbf{D}_j^T \mathbf{A}[\boldsymbol{\pi}_j, :]$  where  $\mathbf{U}[:, \boldsymbol{\pi}_j]$  is selecting columns of  $\mathbf{U}$  and  $\mathbf{A}[\boldsymbol{\pi}_j, :]$  is selecting rows of  $\mathbf{A}$  given indices  $\boldsymbol{\pi}_j$ . The expensive operation of multiplying a  $d \times d$  matrix with the data matrix of dimension  $d \times m$  only needs to occur once for each ReduNet layer the data passes through.

Namely, the computational complexity reduces to

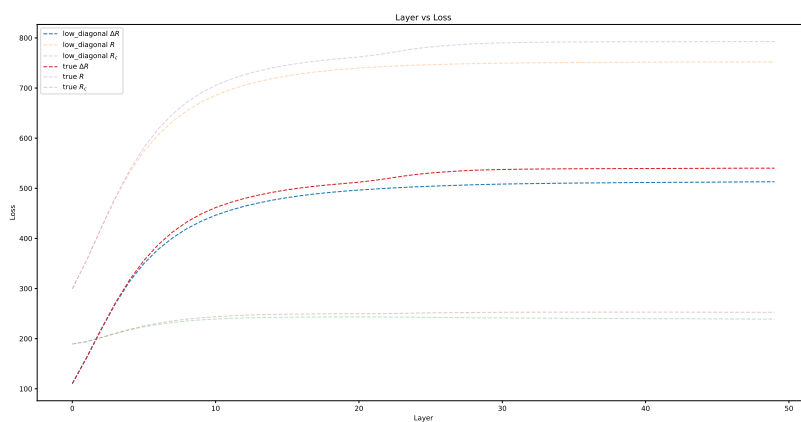
**Original:**  $O((k+1)d^2n) \implies$  **Compressed:**  $\underbrace{O(2d^2n)}_{\mathbf{A}+\mathbf{E}\mathbf{Z}} + \underbrace{O(kdrn) + O(kdr^2)}_{\mathbf{C}_j\mathbf{Z}} = O(3d^2n) + O(kdr^2)$

## Experiments

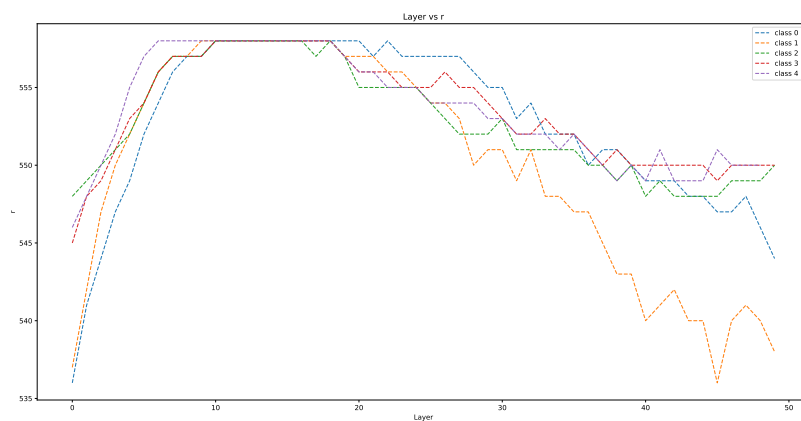
We train a 50 layer ReduNet on 500 samples of 5 classes MNIST. We train the model with tolerance  $\tau = 0.01$ , learning rate  $\eta = 0.1$ , and ball packing error  $\varepsilon^2 = 0.1$ . We observe the following loss curves and number of eigenvectors chosen to approximate the class distributions at each layer (3.4).

Recall that the dimension of vectorized MNIST images is 784. With a tolerance of  $\tau = 0.01$ , each class chose approximately 550 vectors to approximate  $\mathbf{C}_j$  at each layer.





(a) We plot the training loss over iterations using approximated  $C_j$ 's and true  $C_j$ 's.



(b) Number of eigenvectors  $r$  selected over iterations to approximate  $C_j$ .

Figure 3.4: ReduNet trained on MNIST with 500 samples. Notice that the number of eigenvectors chosen is proportional to the compression loss

Additionally, we see that with further training, the number of eigenvectors  $r$  needed to approximate each class’s compression matrix decreases. This is as expected, since at a local minima of the  $MCR^2$  objective, the eigenvector of  $\mathbf{Z}\mathbf{Z}^\top$  become closer to the eigenvectors of  $\mathbf{Z}_j\mathbf{Z}_j^\top$ ,  $\forall j$ . However 550 vectors is no where close to  $\frac{d}{k} = \frac{784}{5}$  which implies that that it is difficult to approximate  $\mathbf{Z}_j$  using a subset of eigenvectors of  $\mathbf{Z}$ , especially at the initial layers. Additionally, this suggests it takes a while for  $\mathbf{Z}$  to converge to an optimal solution.

### 3.3 Conclusion

From our experiments, we conclude that low rank approximation of the covariance matrices is an efficient way to reduce memory of ReduNet. By setting  $\tau = 0.1$  we found a significant reduction in memory footprint in both our experiments in MNIST and CIFAR10 3.3. For future work, it would be interesting to see if we can identify low rank memory reduced networks for the Convolutional ReduNet [2], where  $\mathbf{E}$  and  $\mathbf{C}_j$  are circulant matrices. Additionally,

# Bibliography

- [1] P. Buzzega et al. “Dark Experience for General Continual Learning: a Strong, Simple Baseline”. In: *Adv. Neural Inform. Process. Syst.* (2020).
- [2] Ryan Chan et al. “Deep Networks from the Principle of Rate Reduction”. In: *arXiv preprint arXiv:2010.14765* (2020).
- [3] Kaiming He et al. “Deep residual learning for image recognition”. In: *IEEE Conf. Comput. Vis. Pattern Recog.* (2016).
- [4] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. “Distilling the knowledge in a neural network”. In: *arXiv preprint arXiv:1503.02531* (2015).
- [5] Ronald Kemker and Christopher Kanan. “Fearnnet: Brain-inspired model for incremental learning”. In: *arXiv preprint arXiv:1711.10563* (2017).
- [6] James Kirkpatrick et al. “Overcoming catastrophic forgetting in neural networks”. In: *PNAS* (2017).
- [7] Alex Krizhevsky. “Learning multiple layers of features from tiny images”. In: *Technical report* (2009).
- [8] Alex Krizhevsky, Geoffrey Hinton, et al. “Learning multiple layers of features from tiny images”. In: (2009).
- [9] Yann LeCun. “The MNIST database of handwritten digits”. In: <http://yann.lecun.com/exdb/mnist/> (1998).
- [10] Yann LeCun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.
- [11] Xilai Li et al. “Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting”. In: *arXiv preprint arXiv:1904.00310* (2019).
- [12] Zhizhong Li and Derek Hoiem. “Learning without Forgetting”. In: *Eur. Conf. Comput. Vis.* (2016).
- [13] Yi Ma et al. “Segmentation of multivariate mixed data via lossy data coding and compression”. In: *IEEE transactions on pattern analysis and machine intelligence* 29.9 (2007), pp. 1546–1562.

- [14] Michael McCloskey and Neal J Cohen. “Catastrophic interference in connectionist networks: The sequential learning problem”. In: *Psychology of learning and motivation*. Vol. 24. Elsevier, 1989, pp. 109–165.
- [15] Vardan Papayan, X.Y. Han, and David Donoho. “Prevalence of Neural Collapse during the Terminal Phase of Deep Learning Training”. In: *PNAS* (2020).
- [16] Sylvestre-Alvise Rebuffi et al. “iCaRL: Incremental Classifier and Representation Learning”. In: *IEEE Conf. Comput. Vis. Pattern Recog.* (2017).
- [17] Andrei A Rusu et al. “Progressive neural networks”. In: *arXiv preprint arXiv:1606.04671* (2016).
- [18] Kate Saenko et al. “Adapting visual category models to new domains”. In: *European conference on computer vision*. Springer. 2010, pp. 213–226.
- [19] Jonathan Schwarz et al. “Progress & Compress: A scalable framework for continual learning”. In: *ICML* (2018).
- [20] Gido van de Ven and Andreas Toliás. “Three scenarios for continual learning”. In: *Adv. Neural Inform. Process. Syst.* (2018).
- [21] Chenshen Wu et al. “Memory replay gans: Learning to generate new categories without forgetting”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 5962–5972.
- [22] Yue Wu et al. “Large scale incremental learning”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019, pp. 374–382.
- [23] Ziyang Wu et al. “Incremental learning via rate reduction”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2021).
- [24] Lu Yu et al. “Semantic Drift Compensation for Class-Incremental Learning”. In: *IEEE Conf. Comput. Vis. Pattern Recog.* (2020).
- [25] Yaodong Yu et al. “Learning Diverse and Discriminative Representations via the Principle of Maximal Coding Rate Reduction”. In: *Adv. Neural Inform. Process. Syst.* (2020).
- [26] Friedemann Zenke, Ben Poole, and Surya Ganguli. “Continual Learning Through Synaptic Intelligence”. In: *ICML* (2017).
- [27] Junting Zhang et al. “Class-incremental learning via deep model consolidation”. In: *The IEEE Winter Conference on Applications of Computer Vision*. 2020, pp. 1131–1140.