

7. Appendix

7.1. Derivation of the Lipschitz Constants

To ensure stability of the gradient-based method, we scale the learning rate of Γ and \mathbf{A} by (the inverse of) approximate upper-bounds of the Lipschitz constants of the gradients, $\frac{1}{L_\Gamma}$ and $\frac{1}{L_A}$, respectively. To simplify calculations, we only bound the Lipschitz constant of the gradients with respect to the matrix approximation term M , as we note this term typically dominates the Lipschitz constant of the gradient. Additionally, for notational simplicity we show the derivation for balanced classes which implies $\frac{1}{m\gamma_j} = \frac{k}{m}$.

Now, consider the function $f(\Gamma, \mathbf{A}) = \frac{\mu k}{2m} \sum_{j=1}^k \left\| \mathbf{Z} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}^\top - \Gamma \text{Diag}(\mathbf{A}_j) \Gamma^\top \right\|_F^2$, and note that with some simple algebra one can show the following equivalence:

$$\begin{aligned} f(\Gamma, \mathbf{A}) &= \frac{\mu k}{2m} \sum_{j=1}^k \left\| \mathbf{Z} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}^\top - \Gamma \text{Diag}(\mathbf{A}_j) \Gamma^\top \right\|_F^2 \\ &= \frac{\mu k}{2m} \langle \mathbf{\Pi} \mathbf{\Pi}^\top, (\mathbf{Z}^\top \mathbf{Z})^{\odot 2} \rangle - \frac{\mu k}{m} \langle \mathbf{A} \mathbf{\Pi}^\top, (\Gamma^\top \mathbf{Z})^{\odot 2} \rangle + \frac{\mu k}{2m} \langle \mathbf{A} \mathbf{A}^\top, (\Gamma^\top \Gamma)^{\odot 2} \rangle \end{aligned} \quad (9)$$

where $(\cdot)^{\odot 2}$ denotes raising to the second power entry-wise.

From this, the relevant gradients are

$$\nabla_\Gamma f = -\frac{2\mu k}{m} \sum_{j=1}^k \left(\underbrace{\mathbf{Z} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}^\top \Gamma \text{Diag}(\mathbf{A}_j)}_{:=g_j(\Gamma)} - \underbrace{\Gamma \text{Diag}(\mathbf{A}_j) \Gamma^\top \Gamma \text{Diag}(\mathbf{A}_j)}_{:=h_j(\Gamma)} \right) \quad (10)$$

$$\nabla_{\mathbf{A}} f = \frac{\mu k}{m} [(\Gamma^\top \Gamma)^{\odot 2} \mathbf{A} - (\Gamma^\top \mathbf{Z})^{\odot 2} \mathbf{\Pi}] \quad (11)$$

We first bound the Lipschitz constant of $\nabla_\Gamma f$. For g_j , we have

$$\frac{\|g_j(\Gamma + \Delta\Gamma) - g_j(\Gamma)\|_F}{\|\Delta\Gamma\|_F} = \frac{\|\mathbf{Z} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}^\top \Delta\Gamma \text{Diag}(\mathbf{A}_j)\|_F}{\|\Delta\Gamma\|_F} \quad (12)$$

$$\leq \|\mathbf{Z} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}^\top\|_F \|\mathbf{A}_j\|_\infty, \quad (13)$$

where the inequality is due to the sub-multiplicative property of the Frobenius norm and the operator norm inequality. Now we consider h_j . Here we make the simplifying assumption that in addition to each column of Γ being a unit vector, Γ also approximately satisfies $\Gamma^\top \Gamma = \mathbf{I}$, which is known to be true near the globally optimal solution from the analysis of [21]. Now, we have the following approximation

$$\bar{h}_j(\Gamma) \approx \Gamma \text{Diag}(\mathbf{A}_j)^2. \quad (14)$$

which implies,

$$\frac{\|\bar{h}_j(\Gamma + \Delta\Gamma) - \bar{h}_j(\Gamma)\|_F}{\|\Delta\Gamma\|_F} \approx \frac{\|\Delta\Gamma \text{Diag}(\mathbf{A}_j)^2\|_F}{\|\Delta\Gamma\|_F} \quad (15)$$

$$\leq \|\mathbf{A}_j\|_\infty^2 \quad (16)$$

Since $f(\Gamma, \mathbf{A})$ has a summation over class j , we get the following approximate upper bound for the Lipschitz constant

$$L_\Gamma = \frac{2\mu k}{m} \sum_{j=1}^k \|\mathbf{Z} \text{Diag}(\mathbf{\Pi}_j) \mathbf{Z}^\top\|_F \|\mathbf{A}_j\|_\infty + \|\mathbf{A}_j\|_\infty^2 \quad (17)$$

For the upper bound of the Lipschitz constant of $\nabla_{\mathbf{A}} f$ we simply have

$$\frac{\|\nabla_{\mathbf{A}} f|_{\mathbf{A}+\Delta\mathbf{A}} - \nabla_{\mathbf{A}} f|_{\mathbf{A}}\|_F}{\|\Delta\mathbf{A}\|_F} = \frac{\|\frac{\mu k}{m} (\Gamma^\top \Gamma)^{\odot 2} \Delta\mathbf{A}\|_F}{\|\Delta\mathbf{A}\|_F} \quad (18)$$

$$\leq \frac{\mu k}{m} \|(\Gamma^\top \Gamma)^{\odot 2}\|_F = L_A \quad (19)$$

7.2. Architecture

We utilize the following architectures for the experiments in Section 5. We use a fairly simple architecture for MNIST, and for the other datasets, we use a slightly modified version of ResNet18. Note d is the feature dimension.

Architecture 1 Neural Network Architecture for MNIST

```
1: Conv2d(in_channel=1, out_channel=32, kernel=3, stride=1)
2: ReLU()
3: Conv2d(in_channel=32, out_channel=64, kernel=3, stride=1)
4: ReLU()
5: MaxPool2d(kernel=2, stride=None)
6: Dropout(p=0.25)
7: Flatten()
8: Linear(12544, d)
9: ReLU()
10: Dropout(p=0.5)
11: Linear(d, d)
12: Normalize()
```

For CIFAR-10/100 and Tiny Imagenet, we use the Torchvision ResNet18 model as the featurizer, but we remove the final layer of the ResNet18 and replace it with the following to reshape the output into the desired feature dimension d . We also normalize the output at the end to fulfill the constraint that the features lie on the unit sphere in the MCR² objective.

Architecture 2 Reshaping Layers for ResNet18

```
1: Linear(512, 512, bias=False)
2: BatchNorm1d()
3: ReLU()
4: Linear(512, d, bias=True)
5: Normalize()
```

For cross-entropy experiments, we add another linear layer on top to map the output of the featurizer to logits.

7.3. Data Augmentation

We utilize the following data augmentations for the experiments in Section 5.

Transformations 1 Transformations for MNIST, CIFAR-10, CIFAR-100

```
1: import torchvision.transforms as transforms
2: TRANSFORM = transforms.Compose([
3:     transforms.RandomCrop(32, padding=8),
4:     transforms.RandomHorizontalFlip(),
5:     transforms.ToTensor()])
```

Transformations 2 Transformations for Tiny ImageNet

```
1: import torchvision.transforms as transforms
2: TRANSFORM = transforms.Compose([
3:     transforms.Resize(32)
4:     transforms.RandomHorizontalFlip()
5:     transforms.ToTensor()])
```

7.4. Additional Experiments

The experiments conducted and presented in the main body, Table 2, were only meant to compare the computational efficiency for different methods under the same conditions. The settings however were not chosen to optimize the classification performance since we did not conduct data augmentation or other training recipes normally adopted, e.g. see [22].

Here we report experimental results on CIFAR-100 and Tiny ImageNet by training with a similar training recipe adopted in [22]. Specifically, as in [22], we use PreAct ResNet-18 as the backbone and train with both cross-entropy loss and V-MCR². All networks are trained by SGD with momentum 0.9, and weight decay of 10^{-4} . Similar to [22], we set the learning rate to be 10^{-1} for the first half of training epochs, divide it by 10 for the next quarter of epochs, and finally divide it again by 10 for the remaining iterations. Additionally, we utilize the same transformations to augment the data. Note that the only difference we make between our strategy and the one in [22] is the choice of batch size and total number of training epochs (both specified in Section 4.3) to ensure fair comparison between CE and V-MCR². In addition to the training strategy, we set $\mu = 10^{-1}$ for V-MCR² as we find a smaller value for μ helps stabilize training in the early stage. We keep all other hyperparameters identical to ones specified in Section 4.3.

Dataset	Objective	Training ΔR	Test Accuracy
CIFAR-100	V-MCR ²	130.2177	0.6951
	CE	-	0.7146
Tiny ImageNet 200	V-MCR ²	134.9504	0.4189
	CE	-	0.4843

Table 3. **Comparison of classification performance.** We evaluate the training ΔR and test accuracy of V-MCR² after 2000 training epochs for CIFAR-100 and Tiny ImageNet. For CE, we evaluate the test accuracy at each epoch and report the highest test accuracy achieved across the 2000 epochs. [22] similarly achieved an accuracy 0.7529 for CIFAR-100 and 0.5647 for Tiny ImageNet for CE training. With V-MCR², we achieve a lower test accuracy. Note that the hyperparameters for V-MCR² such as latch frequency, dictionary size, and μ were not rigorously tuned. Additionally, due to the computational cost of the nearest subspace classification algorithm, we only evaluate the test performance at the 2000 epoch mark.