

ON THE DYNAMICS UNDER THE AVERAGED SAMPLE MARGIN LOSS AND BEYOND

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent works have studied implicit biases in deep learning, especially the behavior of last-layer features and classifier weights. However, they usually need to simplify the dynamics under gradient descent due to the intractability of loss functions and neural architectures. In this paper, we introduce a concise loss function as a surrogate, namely the Averaged Sample Margin (ASM) loss, which offers more mathematical opportunities to analyze the closed-form dynamics while requiring few simplifications or assumptions, and allows for more practical considerations. Based on the layer-peeled model that views last-layer features as free optimization variables, we build a complete analysis for the unconstrained, regularized, and spherical constrained cases. We show that these dynamics mainly *converge exponentially fast* to a solution depending on the initialization of features and classifier weights, which can help explain why the training of deep neural networks usually takes only a few hundred epochs. Our theoretical results can also aid in providing insights for improvements in practical training with the ASM loss or other losses, such as explicit feature regularization and rescaled learning rate for spherical cases. Finally, we empirically demonstrate these theoretical results and insights with extensive experiments.

1 INTRODUCTION

Deep learning with neural networks has achieved great success in a variety of tasks (Goodfellow et al., 2016), which, however, is not entirely understood in the interpolation and generalization of the learned models (Zhang et al., 2017; Neyshabur et al., 2017; Nakkiran et al., 2019; Bubeck and Sellke, 2021; Mei and Montanari, 2021). Many modules, including loss functions (Lin et al., 2017; Hui and Belkin, 2021) and optimization algorithms (Auer et al., 2002; Duchi et al., 2011; Zeiler, 2012; Kingma and Ba, 2015), play a crucial role in the training of deep neural networks (DNNs), but lack convincing explanations. Due to the complexity of multilayered architectures, recent works are devoted to simplifying modeling to better understand the behavior of neural networks, and then to gain insights for new algorithms, theoretical, and experimental investigations.

To understand the implicit regularization that improves generalization of the trained models in deep learning, many studies have investigated the implicit bias of gradient descent (Hardt et al., 2016; Sekhari et al., 2021), with an emphasis on the behavior of linear predictors (or called classifiers) over linearly separable data (Soudry et al., 2018; Gunasekar et al., 2018; Nacson et al., 2019; Ji and Telgarsky, 2019; Ji et al., 2020; Shamir, 2021). Soudry et al. (2018) proved that gradient descent iterates under exponentially-tailed loss minimization on separable data are biased toward ℓ_2 -maximum-margin solutions and further showed that continuing to optimize can still lead to performance improvements even if the validation loss increases. Shamir (2021) formally showed that standard gradient methods never overfit on separable data. These works impressively expose the implicit regularization induced by optimization algorithms and help to understand the generalization of the learned models, but they mainly focus on the behavior of linear classifiers that is only the last layer of neural networks, while the classifier actually interacts strongly with the features produced by many nonlinear layers and parameterized layers. Thus, the relevant conclusions do not always apply to deep learning. For example, in (Soudry et al., 2018), the convergence rate of gradient descent is rather slow, wherein for almost all datasets, the distance to the maximum-margin solution decreases only as $O(1/\log t)$. However, the training of DNNs typically takes only a few hundred epochs. In this paper, we show an exponential convergence that is more realistic.

Another research line is after the empirical discovery of Neural Collapse by Pappas et al. (2020), which precisely characterizes a pervasive inductive bias of both features and linear classifiers at the terminal phase of training, and has opened a rich area of exploring this phenomenon in simplified mathematical frameworks (Mixon et al., 2022; Fang et al., 2021; Galanti et al., 2021; Zhu et al., 2021; Hui et al., 2022; Tirer and Bruna, 2022; Lu and Steinerberger, 2022; Kothapalli et al., 2022). Neural Collapse provides a clear view of how the last-layer features and linear classifiers behave after interpolation and enables us to understand the benefit in generalization and robustness after achieving zero training error. Although we have a clear picture of the final phase, the intermediate dynamics is hard to analyze as a result of the intractability of CE. Therefore, some work (Mixon et al., 2022; Han et al., 2022; Zhou et al., 2022a; Tirer and Bruna, 2022; Kothapalli et al., 2022) revolves around the more tractable MSE loss that performs comparably to those trained with CE (Demirkaya et al., 2020; Hui and Belkin, 2021), but they still need to make some simplifications or assumptions about the learning process. Mixon et al. (2022) formulate the gradient flow of the unconstrained feature model as a nonlinear ordinary differential equation and then linearize the equation by claiming that nonlinear terms are negligible for models initialized near the origin. Furthermore, to derive exact dynamics, Han et al. (2022) assume that the gradient flow is along the central path which requires the linear classifier to stay MSE-optimal for features throughout the dynamics. Therefore, MSE is still not simple enough to derive exact dynamics in certain mathematical frameworks, making it difficult to grasp the gap between the modeling and practical optimization.

In this paper, within the layer-peeled model (Fang et al., 2021) (or the unconstrained features model (Mixon et al., 2022)), we attempt to analyze the closed-form dynamics under gradient descent with as few simplifications as possible. More specifically, we introduce a new loss function, namely the Averaged Sample Margin (ASM) loss, which has a concise form that intuitively expresses the objective of classification. Compared to CE and MSE, the ASM loss offers more mathematical opportunities to let us glimpse into deep learning with the closed-form dynamics while requiring few simplifications or assumptions, and will prepare us for more practical considerations and more reasonable designs in a later section. The main contributions of our work are highlighted as follows:

- We derive exact dynamics of last-layer features and prototypes in unconstrained and regularized cases. For spherical constrained cases that do not exhibit convexity, Lipschitzness, and β -smoothness, we also prove that gradient descent biases the normalized features towards a global minimizer.
- We provide the corresponding convergence analysis, which shows that the features and classifier weights converge to a solution depending on the initialization rather than induce the neural collapse solution that forms a simplex equiangular tight frame, suggesting that not all losses under gradient descent would lead to neural collapse (as verified in Section 3).
- We prove that the rate of convergence is exponential as a function of $\zeta(t) = \int_0^t \eta(\tau) d\tau$, where $\eta(\tau)$ denotes the learning rate over time. This exponential convergence rate can help explain why the training of deep neural networks usually takes only a few hundred epochs.
- Moreover, we provide some insights for improvements in practical training with the ASM loss or other losses (cf. Section 4).

2 THE AVERAGED SAMPLE MARGIN LOSS

In this paper, we mainly focus on the behavior of last-layer features and classifier weights in classification DNNs. As described in prior works (Fang et al., 2021; Han et al., 2022), we also consider datasets containing inputs from C different classes with N examples in each class. The last-layer features $\mathbf{h}_{i,c} = \mathbf{f}_{\Theta}(\mathbf{x}_{i,c}) \in \mathbb{R}^p$ extracted from the i -th example $\mathbf{x}_{i,c}$ by a number of parameterized layers $\mathbf{f}_{\Theta}: \mathcal{X} \rightarrow \mathbb{R}^p$ are usually simplified as free optimization variables (Mixon et al., 2022; Fang et al., 2021; Han et al., 2022; Ji et al., 2022). The last layer of the network, *i.e.*, the linear classifier, possesses a class prototype $\mathbf{w}_c \in \mathbb{R}^p$ and bias $b_c \in \mathbb{R}$ for each class $c \in [C]$, which predicts a label using the rule $\arg \max_{c'} (\langle \mathbf{w}_{c'}, \mathbf{h}_{i,c} \rangle + b_{c'})$.

¹To obtain the closed-form solution of neural collapse, we note that Zhou et al. (2022c) assume that $p \geq C - 1$, while Han et al. (2022) assume that $p > C$ since the last-layer features are usually of higher dimension than the number of classes. In this work, we will directly emphasize the relationship between p and C in some scenarios and cover all choices of the feature dimension for others.

To better understand the dynamics of features and prototypes based on gradient descent, we consider a surrogate loss that offers more mathematical opportunities than the hard-to-analyze CE loss and the MSE loss. Specifically, we introduce the Averaged Sample Margin (ASM) loss as follows:

$$L_{ASM}(\mathbf{W}\mathbf{h} + \mathbf{b}, y) = -\mathbf{w}_y^\top \mathbf{h} - b_y + \gamma \sum_{j \neq y} (\mathbf{w}_j^\top \mathbf{h} + b_j), \quad (2.1)$$

where $\gamma > 0$ is the trade-off parameter and y denotes the class label of the feature \mathbf{h} . The sample margin $m(\mathbf{h}, y) = \mathbf{w}_y^\top \mathbf{h} + b_y - \max_{j \neq y} (\mathbf{w}_j^\top \mathbf{h} + b_j)$ (Koltchinskii and Panchenko, 2002; Cao et al., 2019) is defined to measure the discriminativeness for a sample, which satisfies $m(\mathbf{h}, y) \leq \frac{1}{k-1} \sum_{j \neq y} [\mathbf{w}_y^\top \mathbf{h} + b_y - (\mathbf{w}_j^\top \mathbf{h} + b_j)]$, i.e., $L_{ASM}(\mathbf{W}\mathbf{h} + \mathbf{b}, y)$ with $\gamma = \frac{1}{k-1}$ averaging the margins over all non-target classes is the lower bound of $-m(\mathbf{h}, y)$. Moreover, $\frac{1}{k-1} \sum_{j \neq y} [\mathbf{w}_y^\top \mathbf{h} + b_y - (\mathbf{w}_j^\top \mathbf{h} + b_j)]$ can also be regarded as the uninged version that removes the max operator and margin term in the hinge loss². Here, we replace $\frac{1}{k-1}$ with an additional parameter γ that balances positive and negative logits to draw general conclusions. Furthermore, the ASM loss can also be regarded as a variant of CE and multi-binary CE loss:

$$L_{ASM}(\mathbf{W}\mathbf{h} + \mathbf{b}, y) \leq \min \left\{ \log(1 + \exp(-\mathbf{w}_y^\top \mathbf{h} - b_y)) + \gamma \sum_{j \neq y} \log(1 + \exp(\mathbf{w}_j^\top \mathbf{h} + b_j)), \right. \\ \left. (\gamma(C-1) - 1)(\mathbf{w}_y^\top + b_y)\mathbf{h} - \gamma(C-1) \log \frac{\exp(\mathbf{w}_y^\top \mathbf{h} + b_y)}{\sum_{i=1}^C \exp(\mathbf{w}_i^\top \mathbf{h} + b_i)} \right\},$$

Intuitively, the ASM loss promotes the learned feature \mathbf{h} to increase the logit with respect to the target class while decreasing the logits of the other classes. If we follow up the layer-peeled model (Fang et al., 2021) to restrict the norms of both features and prototypes, the global minimizer of $\frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C L_{ASM}(\mathbf{W}\mathbf{h}_{i,c}, y_{i,c})$ (the bias term \mathbf{b} is omitted) will lead to *Neural Collapse* (Papayan et al., 2020; Han et al., 2022):

Lemma 2.1 (Neural Collapse under Averaged Sample Margin Loss). *For norm-bounded prototypes and features, i.e., $\|\mathbf{w}_c\|_2 \leq E_1$ and $\|\mathbf{h}_{i,c}\|_2 \leq E_2$, $\forall i \in [N], \forall c \in [C]$, the global minimizer of $\frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C L_{ASM}(\mathbf{W}\mathbf{h}_{i,c}, y_{i,c})$ implies neural collapse when $p \geq C - 1$. More specifically, the global minimizer is uniquely obtained at $\frac{\mathbf{w}_i^\top \mathbf{w}_j}{\|\mathbf{w}_i\|_2 \|\mathbf{w}_j\|_2} = -\frac{1}{C-1}$, $\forall i \neq j$, $\frac{\mathbf{w}_{y_{i,c}}^\top \mathbf{h}_{i,c}}{\|\mathbf{w}_{y_{i,c}}\|_2 \|\mathbf{h}_{i,c}\|_2} = 1$, $\|\mathbf{w}_c\|_2 = E_1$, and $\|\mathbf{h}_{i,c}\|_2 = E_2$, $\forall i \in [N], \forall c \in [C]$.*

This lemma shows that the neural collapse solution is the only global optimal solution to minimize $\frac{1}{CN} \sum_{i,c} L_{ASM}(\mathbf{W}\mathbf{h}_{i,c}, y_{i,c})$ in the norm-bounded case. However, there exists an undesired direction to minimize the ASM loss in unconstrained cases, since the norm of features and prototypes tends to grow to infinity. For example, we can scale up \mathbf{W} and \mathbf{b} to obtain a smaller loss if $L_{ASM}(\mathbf{W}\mathbf{h} + \mathbf{b}, y) < 0$, which will happen analogously to CE (Liu et al., 2016; Wang et al., 2017; Zhou et al., 2022c). In this paper, we will analytically characterize the direction in which features \mathbf{H} and prototypes \mathbf{W} diverge. Specifically, we show that the continual gradient flow exhibits an implicit bias associated with the initialization of features and prototypes.

3 MAIN THEORETICAL RESULTS

In this section, we build a complete analysis of the closed-form dynamics of the last-layer features and prototypes under the ASM loss in unconstrained, regularized, and spherical constrained cases. We then derive the convergence analysis of these dynamics, which mainly shows exponential convergence. **All proof can be found in the Appendix A.**

3.1 UNCONSTRAINED CASE

We first consider the unconstrained case (Mixon et al., 2022; Ji et al., 2022) in which there is no constraint or regularization on features and prototypes, i.e., learning with the following objective

$$\mathcal{L}_{ASM} = \frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C \left[-(\mathbf{w}_{y_{i,c}}^\top \mathbf{h}_{i,c} + b_{y_{i,c}}) + \gamma \sum_{j \neq y_{i,c}} (\mathbf{w}_j^\top \mathbf{h}_{i,c} + b_j) \right], \quad (3.1)$$

²The multi-class hinge loss is $L_{\text{hinge}}(\mathbf{W}\mathbf{h} + \mathbf{b}, y) = \sum_{i \neq y} \max\{0, (\mathbf{w}_i^\top \mathbf{h} + b_i) - (\mathbf{w}_y^\top \mathbf{h} + b_y) + m\}$, where m is the margin term.

which can be reformulated as

$$\mathcal{L}_{ASM} = \frac{1}{CN} \text{Tr} \left((\gamma \mathbf{1}_C \mathbf{1}_{CN}^\top - (1 + \gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top))^\top \mathbf{W}^\top \mathbf{H} \right) + \frac{\gamma^{C-\gamma-1}}{C} \mathbf{1}_C^\top \mathbf{b}, \quad (3.2)$$

where $\mathbf{H} = [\mathbf{h}_{1,1}, \dots, \mathbf{h}_{1,C}, \mathbf{h}_{2,1}, \dots, \mathbf{h}_{N,C}] \in \mathbb{R}^{p \times CN}$ is the matrix resulting from stacking together the feature vectors as columns, $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C] \in \mathbb{R}^{p \times C}$ is the matrix resulting from stacking stacking class prototypes as columns, \otimes denotes the Kronecker product, \mathbf{I}_C is the identity matrix, $\mathbf{1}_C$, $\mathbf{1}_N$, and $\mathbf{1}_{CN}$ are the length- C , $-N$, and $-CN$ vectors of ones, respectively. Without loss of generality, we represent the label set $\{y_{i,c}\}_{1 \leq i \leq N, 1 \leq c \leq C}$ as the columns of the matrix $\mathbf{I}_C \otimes \mathbf{1}_N^\top$.

We follow the unconstrained features and prototypes modeling perspective (Mixon et al., 2022; Ji et al., 2022) that treats \mathbf{H} as a free optimization variable. Within this model, we analyze the continuous dynamics of features \mathbf{H} , prototypes \mathbf{W} and biases \mathbf{b} with gradient flow where *time-of-training* is denoted by the variable t ³.

Let $\mathbf{Y} = \frac{1}{CN} (\gamma \mathbf{1}_C \mathbf{1}_{CN}^\top - (1 + \gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top))$, then $\mathcal{L}_{ASM} = \text{Tr}(\mathbf{Y}^\top \mathbf{W}^\top \mathbf{H}) + \frac{\gamma^{C-\gamma-1}}{C} \mathbf{1}_C^\top \mathbf{b}$. Taking the partial derivative with respect to \mathbf{H} , \mathbf{W} , and \mathbf{b} , respectively, we have:

$$\nabla_{\mathbf{H}} \mathcal{L}_{ASM} = \mathbf{W} \mathbf{Y}, \quad \nabla_{\mathbf{W}} \mathcal{L}_{ASM} = \mathbf{H} \mathbf{Y}^\top, \quad \nabla_{\mathbf{b}} \mathcal{L}_{ASM} = \frac{\gamma^{C-\gamma-1}}{C} \mathbf{1}_C, \quad (3.3)$$

and the corresponding learning dynamics following Gradient Descent (GD) is

$$\mathbf{H}'(t) = \eta_1(t) \mathbf{W}(t) \mathbf{M}, \quad \mathbf{W}'(t) = \eta_2(t) \mathbf{H}(t) \mathbf{M}^\top, \quad \mathbf{b}'(t) = -\eta_2(t) \frac{\gamma^{C-\gamma-1}}{C} \mathbf{1}_C, \quad (3.4)$$

where $\mathbf{M} = -\mathbf{Y} = \frac{1}{CN} ((1 + \gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma \mathbf{1}_C \mathbf{1}_{CN}^\top)$, η_1 and η_2 are the learning rate of the features and prototypes, respectively. The reason to introduce different learning rates is that the representation \mathbf{H} is actually the result of the interaction between a number of nonlinear layers and parameterized layers. This means that even if we use the same learning rate to optimize all parameters of the network, the feature \mathbf{H} assumed to be a free optimization variable is almost impossible to be optimized at this learning rate⁴. Moreover, we consider dynamic learning rates $\eta_1 = \eta_1(t)$ and $\eta_2 = \eta_2(t)$ that are usually adopted in practical implementations, such as the cosine annealing decay (Loshchilov and Hutter, 2017). As can be seen, the dynamics of features \mathbf{H} and prototypes \mathbf{W} are independent of the bias term \mathbf{b} , thus we can analyze the dynamics of \mathbf{H} and \mathbf{W} jointly, and analyze \mathbf{b} independently:

Theorem 3.1 (Dynamics of Features, Prototypes and Biases without Constraints). *Consider the continual gradient flow (Equation 3.4) in which the dynamics follow the gradient descent direction of the Averaged Sample Margin loss in Eq. (3.2). Let $\mathbf{Z}(t) = (\mathbf{H}(t), \mathbf{W}(t))$, if $\eta_1(t_1)\eta_2(t_2) = \eta_1(t_2)\eta_2(t_1)$ for any $t_1, t_2 \geq 0$, we have the following closed-form dynamics*

$$\begin{aligned} \mathbf{Z}(t) = & \Pi_1^+ \mathbf{Z}_0 (\alpha_1^+(t) \mathbf{C}(t) + \beta_1^+(t) \mathbf{I}_{C(N+1)}) + \Pi_1^- \mathbf{Z}_0 (\alpha_1^-(t) \mathbf{C}(t) + \beta_1^-(t) \mathbf{I}_{C(N+1)}) + \Pi_3 \mathbf{Z}_0 \\ & + \Pi_2^+ \mathbf{Z}_0 (\alpha_2^+(t) \mathbf{C}(t) + \beta_2^+(t) \mathbf{I}_{C(N+1)}) + \Pi_2^- \mathbf{Z}_0 (\alpha_2^-(t) \mathbf{C}(t) + \beta_2^-(t) \mathbf{I}_{C(N+1)}), \end{aligned} \quad (3.5)$$

and

$$\mathbf{b}(t) = \mathbf{b}_0 + \frac{(1 + \gamma - \gamma^C) \zeta_2(t)}{C} \mathbf{1}_C, \quad (3.6)$$

where α_ϵ^ϵ , α_2^ϵ , β_1^ϵ and β_2^ϵ for $\epsilon \in \{\pm\}$ are the scalars that only depend on C , N , γ , η_1 and η_2 (where the detailed forms can be seen in the appendix), $\mathbf{Z}_0 = (\mathbf{H}_0, \mathbf{W}_0)$, $\mathbf{C}(t) = \begin{pmatrix} \zeta_1(t) \mathbf{I}_{CN} & 0 \\ 0 & \zeta_2(t) \mathbf{I}_C \end{pmatrix}$,

$\zeta_1(t) = \int_0^t \eta_1(\tau) d\tau$, $\zeta_2(t) = \int_0^t \eta_2(\tau) d\tau$, Π_1^ϵ , Π_2^ϵ and Π_3 for $\epsilon \in \{\pm\}$ are orthogonal projection operators onto the following respective eigenspaces:

$$\begin{aligned} \mathcal{E}_1^\epsilon & := \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}} (\mathbf{W} \otimes \mathbf{1}_N^\top), \mathbf{W} \mathbf{1}_C = 0\}, \\ \mathcal{E}_2^\epsilon & := \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}} \mathbf{h} \mathbf{1}_{CN}^\top, \mathbf{W} = \mathbf{h} \mathbf{1}_C^\top, \mathbf{h} \in \mathbb{R}^p\}, \\ \mathcal{E}_3 & := \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} (\mathbf{I}_C \otimes \mathbf{1}_N) = 0, \mathbf{W} = 0\}. \end{aligned} \quad (3.7)$$

³Intuitively, we interpret $t = 0$ as the initial state, that is $\mathbf{H}(0) = \mathbf{H}_0$, $\mathbf{W}(0) = \mathbf{W}_0$, and $\mathbf{b}(0) = \mathbf{b}_0$.

⁴In this paper, we mainly assume that $\eta_1(t_1)\eta_2(t_2) = \eta_1(t_2)\eta_2(t_1)$ for any pair of values of t , t_1 and t_2 . This condition will be satisfied if and only if $\eta_1(t)$ is a scaled version of $\eta_2(t)$, i.e., $\eta_1(t) = s \cdot \eta_2(t)$

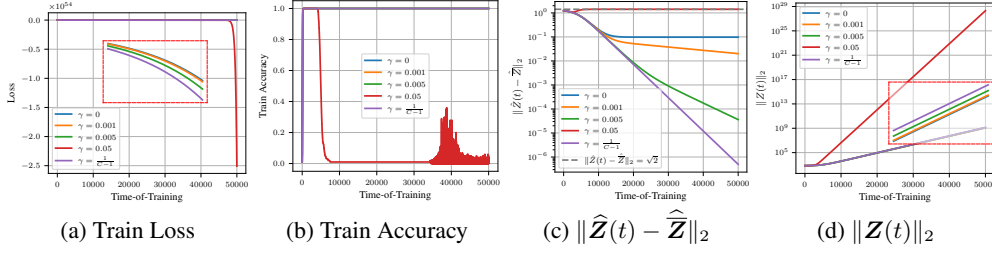


Figure 1: Verification of the behavior of gradient descent iterates in Equation (3.4) with $\gamma \in \{0, 0.05, 0.005, 0.001, \frac{1}{C-1}\}$, where we set $p = 512$, $C = 100$, $N = 10$, $\eta_1(t) = \eta_2(t) = 0.1$ (i.e., $s = \frac{\eta_1(0)}{\eta_2(0)} = 1$, thus $\bar{\mathbf{Z}} = \Pi_1^+ \mathbf{Z}_0$ according to Corollary 3.2), and then randomly initialize \mathbf{H}_0 and \mathbf{W}_0 . **The curve in the red box represents the zoomed-in curve of the last 2,000 epochs.** (a) The training loss. (b) The training accuracy with the prediction rule $\arg \max_c \mathbf{w}_c^\top \mathbf{h}$. As expected, the features align to their corresponding prototypes when $\gamma < \frac{2}{C-2}$. (c) denote the distance between $\hat{\mathbf{Z}}(t) = \mathbf{Z}(t)/\|\mathbf{Z}(t)\|_2$ and $\hat{\mathbf{Z}} = \bar{\mathbf{Z}}/\|\bar{\mathbf{Z}}\|_2$. As expected in Eq. (3.9), the convergence rate is exponential when $0 < \gamma < \frac{2}{C-2}$, and will be fastest if $\gamma = \frac{1}{C-1}$. (d) denotes the norm of $\mathbf{Z}(t)$ that increases exponentially. As can be noticed, $\hat{\mathbf{Z}}(t)$ does not converge to $\hat{\mathbf{Z}}$ but tend to be orthogonal to $\hat{\mathbf{Z}}$ when $\gamma = 0.05 > \frac{2}{C-2}$, that is, $\lim_{t \rightarrow \infty} \|\hat{\mathbf{Z}}(t) - \hat{\mathbf{Z}}\|_2 = \sqrt{2}$.

Remark. Note that \mathcal{E}_1^+ , \mathcal{E}_2^+ and \mathcal{E}_3 are orthogonal to each other, \mathcal{E}_1^+ (or \mathcal{E}_1^-) denotes the subspace where all features are in the same (or opposite) direction of their corresponding prototypes while the mean of prototypes is zero, \mathcal{E}_2^+ denotes the subspace where all features and prototypes collapse respectively into two scaled versions of the same unit vector, and \mathcal{E}_3 denotes the subspace where the mean of all features from the same class is zero with all prototypes being zero. For classification tasks, we expect the features align to their corresponding prototypes with a cosine similarity of 1, i.e., the solution in $\mathcal{E} = \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = k\mathbf{W} \otimes \mathbf{1}_N^T, \mathbf{W} \in \mathbb{R}^{p \times C}, k \in \mathbb{R}^+\}$ that implies two manifestations of Neural Collapse: *variability collapse* and *convergence to self-duality* (Papayan et al., 2020). In the following, we show that $\frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|}$ will converge to a solution in \mathcal{E} :

Corollary 3.2 (Convergence in the Unconstrained Case). *Under the conditions and notation of Theorem 3.1, let $s = \frac{\eta_1(0)}{\eta_2(0)}$, if $0 < \gamma < \frac{2}{C-2}$ (where $C > 2$) or $C = 2$, and $\lim_{t \rightarrow \infty} \zeta_1(t) = \infty$, the gradient flow (as in Eq. (3.4)) will behave as:*

$$e^{-\frac{(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} \mathbf{Z}(t) = \bar{\mathbf{Z}} + \Delta(t), \quad (3.8)$$

where $\bar{\mathbf{Z}} = \left(\frac{1+\sqrt{s}}{2}\mathbf{H}_1^+ + \frac{1-\sqrt{s}}{2}\mathbf{H}_1^-, \frac{1+\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^-\right)$, $(\mathbf{H}_1^+, \mathbf{W}_1^+) = \Pi_1^+ \mathbf{Z}_0$, $(\mathbf{H}_1^-, \mathbf{W}_1^-) = \Pi_1^- \mathbf{Z}_0$, and the residual term $\Delta(t)$ decreases as least as $\|\Delta(t)\| = O\left(e^{\frac{\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}} \cdot \max\{-\gamma C, (C-2)\gamma-2\}}\right)$, and so the normalized $\mathbf{Z}(t)$ converges to $\frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|}$ in

$$\left\| \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} - \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \right\| = O\left(e^{\frac{\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}} \cdot \max\{-\gamma C, (C-2)\gamma-2\}}\right), \quad (3.9)$$

which further indicates $\lim_{t \rightarrow \infty} \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} \in \mathcal{E}$. Moreover, if $\gamma \neq \frac{1}{C-1}$, then $\lim_{t \rightarrow \infty} \frac{\max_i b_i(t)}{\min_i b_i(t)} = 1$.

This corollary shows that even without any explicit regularization, when minimizing the ASM loss using gradient descent, we prove that $\frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|}$ converges to a special point $\frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \in \mathcal{E}$ determined by the initialization of \mathbf{Z}_0 and the ratio s , but does not necessarily perform as a neural collapse solution that forms a simplex equiangular tight frame as CE or MSE do (Papayan et al., 2020; Han et al., 2022). In addition, the rate of convergence is exponential as a function of the integral of the learning rate, i.e., $O\left(e^{\frac{\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}} \cdot \max\{-\gamma C, (C-2)\gamma-2\}}\right)$, which indicates that the convergence of updating both features and prototypes by gradient descent is much faster than $O(1/\log t)$ that only updates prototypes (linear predictors) on linearly separable data (Soudry et al., 2018). In a

sense, this convergence rate explains why training deep neural networks usually takes only several hundred or thousand epochs. Moreover, if $\gamma = \frac{1}{C-1}$, we can obtain the fastest convergence of Eq. (3.9), that is, $\left\| \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} - \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \right\| = O\left(e^{-\frac{\sqrt{\zeta_1(t)\zeta_2(t)}}{(C-1)\sqrt{N}}}\right)$. For example, when $\eta_1(t) = \eta$ is a constant learning rate, then $\zeta_1(t) = \eta t \rightarrow \infty$ as $t \rightarrow \infty$, and the gradient flow in Eq. (3.4) shows an exponential convergence rate of $\frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|}$ to $\frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|}$, but it also leads to an exponential increase with the rate $e^{\frac{(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}}$ in the norm of features and prototypes, which is almost unbearable for the practical training of the models. Therefore, we need to limit the excessive growth of these norms.

3.2 REGULARIZED CASE

We then consider the following regularized optimization problem that explicitly introduces ℓ_2 regularization (often called ‘‘weight decay’’) on features, prototypes, and biases to avoid large values:

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \text{Tr}(\mathbf{Y}^\top \mathbf{W}^\top \mathbf{H}) + \frac{\gamma C - \gamma - 1}{C} \mathbf{1}_C^\top \mathbf{b} + \frac{\lambda}{2} (\|\mathbf{H}\|_F^2 + \|\mathbf{W}\|_F^2 + \|\mathbf{b}\|_2^2), \quad (3.10)$$

where $\|\cdot\|_F$ denotes the Frobenius norm and $\lambda > 0$ is the regularization parameter.

Taking the partial derivative with respect to \mathbf{H} , \mathbf{W} , and \mathbf{b} , we have

$\nabla_{\mathbf{H}} \mathcal{L}_{ASM} = \mathbf{W}\mathbf{Y} + \lambda \mathbf{H}(t)$, $\nabla_{\mathbf{W}} \mathcal{L}_{ASM} = \mathbf{H}\mathbf{Y}^\top + \lambda \mathbf{W}(t)$, $\nabla_{\mathbf{b}} \mathcal{L}_{ASM} = \frac{\gamma C - \gamma - 1}{C} \mathbf{1}_C + \lambda \mathbf{b}(t)$, and the corresponding learning dynamics following gradient descent with different learning rates and ℓ_2 -norm regularization can be formulated as

$$\begin{aligned} \mathbf{H}'(t) &= \eta_1(t) \mathbf{W}(t) \mathbf{M} - \lambda \eta_1(t) \mathbf{H}(t), \\ \mathbf{W}'(t) &= \eta_2(t) \mathbf{H}(t) \mathbf{M}^\top - \lambda \eta_2(t) \mathbf{W}(t), \\ \mathbf{b}'(t) &= -\eta_2(t) \frac{\gamma C - \gamma - 1}{C} \mathbf{1}_C - \lambda \eta_2(t) \mathbf{b}(t). \end{aligned} \quad (3.11)$$

The dynamics of this regularized gradient flow can be proved as follows

Theorem 3.3 (Dynamics of Features, Prototypes, and Biases under Weight Decay). *Consider the continual gradient flow (Equation 3.11) in which the dynamics follows the gradient descent direction of the Averaged Sample Margin loss in Eq. (3.10). Let $\mathbf{Z}(t) = (\mathbf{H}(t), \mathbf{W}(t))$, if $\eta_1(t_1)\eta_2(t_2) = \eta_1(t_2)\eta_2(t_1)$ for any $t_1, t_2 \geq 0$, we have the following closed-form dynamics.*

$$\begin{aligned} \mathbf{Z}(t) &= \Pi_1^+ \mathbf{Z}_0 \begin{pmatrix} a_1^+(t) \mathbf{I}_{CN} & 0 \\ 0 & b_1^+(t) \mathbf{I}_C \end{pmatrix} + \Pi_1^- \mathbf{Z}_0 \begin{pmatrix} a_1^-(t) \mathbf{I}_{CN} & 0 \\ 0 & b_1^-(t) \mathbf{I}_C \end{pmatrix} + \\ &\quad \Pi_2^+ \mathbf{Z}_0 \begin{pmatrix} a_2^+(t) \mathbf{I}_{CN} & 0 \\ 0 & b_2^+(t) \mathbf{I}_C \end{pmatrix} + \Pi_2^- \mathbf{Z}_0 \begin{pmatrix} a_2^-(t) \mathbf{I}_{CN} & 0 \\ 0 & b_2^-(t) \mathbf{I}_C \end{pmatrix} + \\ &\quad \Pi_3 \mathbf{Z}_0 \begin{pmatrix} a_3(t) \mathbf{I}_{CN} & 0 \\ 0 & b_3(t) \mathbf{I}_C \end{pmatrix}, \end{aligned} \quad (3.12)$$

and

$$\mathbf{b}(t) = \phi(t) \left(\mathbf{b}_0 + \frac{1+\gamma-\gamma C}{C} \psi(t) \mathbf{1}_C \right), \quad (3.13)$$

where $\Pi_1^+ \mathbf{Z}_0$, $\Pi_1^- \mathbf{Z}_0$, $\Pi_2^+ \mathbf{Z}_0$, $\Pi_2^- \mathbf{Z}_0$, and $\Pi_3 \mathbf{Z}_0$ follow the definition in Theorem 3.1, a_1^ϵ , a_2^ϵ , b_1^ϵ , b_2^ϵ , a_3 , and b_3 for $\epsilon \in \{\pm\}$ are the scalars that depend only on C , N , γ , λ , η_1 , and η_2 (where the detailed forms can be seen in Appendix A), $\phi(t) = \exp(-\lambda \int_0^t \eta_2(\tau) d\tau)$, and $\psi(t) = \int_0^t \zeta_2(\tau) \exp(\lambda \int_0^\tau \eta_2(s) ds) d\tau$.

The convergence under the regularized case can also be derived as:

Corollary 3.4 (Convergence Under ℓ_2 Regularization). *Under the conditions and notation of Theorem 3.3, let $s = \frac{\eta_1(0)}{\eta_2(0)}$, if $0 < \gamma < \frac{2}{C-2}$ (where $C > 2$) or $C = 2$, and $\lim_{t \rightarrow \infty} \zeta_1(t) = \infty$, then there exist constants π_h^+ , π_h^- , π_w^+ , π_w^- , and ω only depending on λ , γ , s , C , and N , such that the gradient flow (as in Eq. (3.11)) behaves as:*

$$\left\| \frac{\mathbf{H}(t)}{\|\mathbf{H}(t)\|} - \frac{\pi_h^+ \mathbf{H}_1^+ + \pi_h^- \mathbf{H}_1^-}{\|\pi_h^+ \mathbf{H}_1^+ + \pi_h^- \mathbf{H}_1^-\|} \right\| + \left\| \frac{\mathbf{W}(t)}{\|\mathbf{W}(t)\|} - \frac{\pi_w^+ \mathbf{W}_1^+ + \pi_w^- \mathbf{W}_1^-}{\|\pi_w^+ \mathbf{H}_1^+ + \pi_w^- \mathbf{H}_1^-\|} \right\| = O(e^{-\omega \zeta_2(t)}), \quad (3.14)$$

where $(\mathbf{H}_1^+, \mathbf{W}_1^+) = \Pi_1^+ \mathbf{Z}_0$, $(\mathbf{H}_1^-, \mathbf{W}_1^-) = \Pi_1^- \mathbf{Z}_0$.

Furthermore, we have the following convergence results for $\mathbf{Z}(t)$:

- If $\lambda > \frac{1+\gamma}{C\sqrt{N}}$, then $\lim_{t \rightarrow \infty} \|\mathbf{Z}(t)\| = 0$;
- If $\lambda = \frac{1+\gamma}{C\sqrt{N}}$, then $\lim_{t \rightarrow \infty} \mathbf{Z}(t) = \left(\mathbf{H}_1^+ + \frac{1-s}{1+s}\mathbf{H}_1^-, \mathbf{W}_1^+ - \frac{1-s}{1+s}\mathbf{W}_1^-\right)$;
- If $\lambda < \frac{1+\gamma}{C\sqrt{N}}$, then $\lim_{t \rightarrow \infty} \|\mathbf{Z}(t)\| = \infty$.

Remark. The results in Corollary 3.4 suggest that adding an appropriate weight decay on both features and prototypes can avoid impractical effects, since the norm of $\mathbf{Z}(t)$ shrinking to 0 or diverging toward infinity will significantly affect the training of DNNs. Several recent works (Zhu et al., 2021; Zhou et al., 2022a) described that the features are implicitly penalized, but this implicit penalization is fragile, as depicted in Fig. 19. As a consequence, we emphasize *adding explicit regularization to features, not just implicit penalization attached by other components* in Sec. 4.2.

3.3 SPHERICAL CONSTRAINED CASE

We finally consider another constrained case in which features are restricted on the p -sphere $\mathbb{S}^{p-1} = \{\mathbf{x} : \|\mathbf{x}\|_2 = 1, \mathbf{x} \in \mathbb{R}^p\}$ by explicitly performing ℓ_2 normalization to prevent arithmetic overflow or underflow happening in training DNNs, and we fix the prototypes⁵ to satisfy $\mathbf{W}\mathbf{1}_C = \mathbf{0}$, then the optimization problem in Eq. (3.1) can be reformulated as

$$\min_{\mathbf{H}} -\frac{1+\gamma}{CN} \text{Tr}((\mathbf{I}_C \otimes \mathbf{1}_N) \mathbf{W}^\top \widehat{\mathbf{H}}), \quad (3.15)$$

where $\widehat{\mathbf{H}} = (\widehat{\mathbf{h}}_{1,1}, \dots, \widehat{\mathbf{h}}_{c,N})$, and $\widehat{\mathbf{h}}_{i,c} = \frac{\mathbf{h}_{i,c}}{\|\mathbf{h}_{i,c}\|_2}$. We then take the partial derivative with respect to \mathbf{H} , and the discrete dynamical system based on gradient descent is formulated as

$$\mathbf{H}(t+1) = \mathbf{H}(t) + \frac{(1+\gamma)\eta(t)}{CN} \left(\frac{\partial \widehat{\mathbf{H}}}{\partial \mathbf{H}} \Big|_{\mathbf{H}=\mathbf{H}(t)} \right)^\top \mathbf{W}(\mathbf{I}_C \otimes \mathbf{1}_N^\top), \quad (3.16)$$

where $\left(\frac{\partial \widehat{\mathbf{H}}}{\partial \mathbf{H}}\right)^\top$ is an element-wise operator, that is $\left(\frac{\partial \widehat{\mathbf{H}}}{\partial \mathbf{H}}\right)^\top \mathbf{H}' = \left(\frac{1}{\|\mathbf{h}_{i,c}\|_2} (\mathbf{I}_p - \widehat{\mathbf{h}}_{i,c} \widehat{\mathbf{h}}_{i,c}^\top) \mathbf{h}'_{i,c}\right)$ for $\mathbf{H}' \in \mathbb{R}^{p \times CN}$.

Despite the fact that the optimization objective in Eq. (3.15) does not show convexity, Lipschitzness, and β -smoothness on \mathbf{H} due to the ℓ_2 normalization operator, the normalized features that obey the gradient descent iterates in Eq. (3.16) can still converge to their corresponding normalized prototypes, *i.e.*, achieve the global minimum of Eq. (3.15):

Theorem 3.5. *Considering the discrete dynamics in Eq. (3.16), if $\forall i \in [N], c \in [C]$, $\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) > -1$, the learning rate $\eta(t)$ satisfies that $\frac{\eta(t)}{\|\mathbf{h}_{i,c}(t)\|_2}$ is non-increasing, $\frac{\eta(0)(1+\gamma)}{CN\|\mathbf{h}_{i,c}(0)\|_2} \leq \frac{1}{\|\mathbf{w}_c\|_2}$, $\lim_{t \rightarrow \infty} \frac{\eta(t+1)}{\eta(t)} = 1$, and there exists a constant $\varepsilon > 0$, s.t., $\eta(t) > \varepsilon$, then we have*

$$\lim_{t \rightarrow \infty} \left\| \widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N^\top) \right\| = 0, \quad (3.17)$$

and further if $\lim_{t \rightarrow \infty} \|\mathbf{H}(t)\| < \infty$, then there exists a constant $\mu > 0$, such that the error above shows exponential decrease:

$$\left\| \widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N^\top) \right\| = O(e^{-\mu t}). \quad (3.18)$$

Moreover, if $\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) = -1$, then $\mathbf{h}_{i,c}(t) = \mathbf{h}_{i,c}(0)$.

4 INSIGHTS AND EXPERIMENTS

In this section, we provide some insights to better train DNNs according to the conclusions in Sec. 3. We then corroborate our theoretical results and insights with extensive experiments. **More details and results can be found in Appendix B.**

⁵The relevant studies are still few and often require some strict assumptions since the learning dynamics is very complicated when \mathbf{w} participates the optimization process with both feature and prototypes normalization. In this paper, we are going to try a more concise theoretical analysis with fixed prototypes.

⁶This aims to simplify Eq. (2.1) as the objective of feature alignment, that is, $L_{ASM}(\mathbf{W}\widehat{\mathbf{h}}, y) = -\mathbf{w}_y^\top \widehat{\mathbf{h}} + \gamma \sum_{j \neq y} \mathbf{w}_j^\top \widehat{\mathbf{h}} = \frac{(1+\gamma)\|\mathbf{w}\|_2}{2} (\|\widehat{\mathbf{h}} - \widehat{\mathbf{w}}_y\|_2^2 - 2)$, and the global minimum will be obtained at $\widehat{\mathbf{h}} = \widehat{\mathbf{w}}_y$.

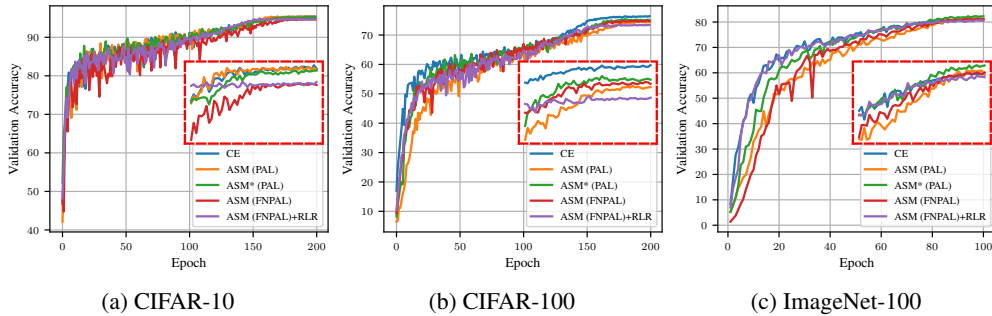


Figure 2: Validation accuracy of different loss functions on CIFAR-10, CIFAR-100, and ImageNet, where * denotes training with explicit feature regularization. PAL and FNPAL denote the model trained with prototype-anchored learning (PAL) and feature-normalized and prototype-anchored learning (FNPAL) (Zhou et al., 2022b). The curve in the red box represents the zoomed-in curve of the last 50 epochs. As can be seen, DNNs trained with the ASM loss can achieve comparative or even better performance compared to those of the CE loss.

4.1 THE ASM LOSS WITH PROTOTYPE-ANCHORED LEARNING

Since directly using the ASM loss will lead to volatile effects, which is mainly reflected in the rapid increases of feature norms and the imbalance between class prototypes when training DNNs with the stochastic gradient method, as shown in Fig. 19. Inspired by recent works (Zhou et al., 2022b; Kasarla et al., 2022) that use the neural collapse solution as an inductive bias, we can anchor prototypes \mathbf{W} during training, and then the dynamics of \mathbf{H} with ℓ_2 -norm-based regularization in Eq. (3.10) can be formulated as the first-order non-homogeneous linear difference equation:

$$\mathbf{H}'(t) = \eta(t)\mathbf{W}\mathbf{M} - \lambda\eta(t)\mathbf{H}(t), \quad (4.1)$$

and the solution to the non-homogeneous linear DE can be easily derived:

Theorem 4.1. Consider the continual gradient flow (Equation (4.1)) in which the prototypes \mathbf{W} are fixed, we have the closed-form dynamics:

$$\mathbf{H}(t) = e^{-\lambda \int_0^t \eta(\tau) d\tau} \mathbf{H}(0) + \frac{1 - e^{-\lambda \int_0^t \eta(\tau) d\tau}}{\lambda} \mathbf{W}\mathbf{M}, \quad (4.2)$$

which further indicates that $\|\mathbf{H}(t) - \frac{1}{\lambda} \mathbf{W}\mathbf{M}\| = O\left(e^{-\lambda \int_0^t \eta(\tau) d\tau}\right)$.

Prototype-anchored learning (PAL) is very effective in alleviating the instability by transforming the classification problem as a feature alignment problem (Zhou et al., 2022b). As shown in Fig. 2, the ASM loss with PAL and FNPAL can achieve comparable or even better results compared to CE.

Table 1: Validation accuracies on long-tailed CIFAR-10/-100 with CE and different explicit feature regularization ($\lambda \in \{0, 5e-6, 1e-5, 5e-5\}$). Imbalance ratio $\rho = \frac{\max_i n_i}{\min_i n_i}$ is the ratio between sample sizes of the most frequent and least frequent classes, and $\rho = 1$ denotes the original CIFAR-10/-100. $\lambda = 0$ denotes the model training with CE. All values are percentages. **Bold** numbers indicate the results that are better than baseline. The best results are underlined. As can be seen, explicit feature regularization effectively improve the performance on long-tailed classification in most cases, even for normal classification.

Dataset	Long-tailed CIFAR-10					Long-tailed CIFAR-100				
	Imbalance Ratio	100	50	20	10	1	100	50	20	10
$\lambda = 0$	67.81	72.93	83.97	88.37	95.28	33.37	39.40	42.96	56.38	75.42
$\lambda = 5e-6$	67.84	72.85	83.17	89.06	95.27	36.00	41.92	50.75	60.13	76.48
$\lambda = 1e-5$	67.74	76.14	84.17	89.19	95.23	36.61	42.36	49.21	58.91	77.34
$\lambda = 5e-5$	69.74	77.29	84.92	88.64	95.39	34.88	42.74	54.72	60.84	76.19

4.2 EXPLICIT FEATURE REGULARIZATION

In this paper, we directly consider explicit feature regularization to avoid excessive growth of feature norms, rather than fully adopt the implicit penalization described in prior works (Fang et al., 2021;

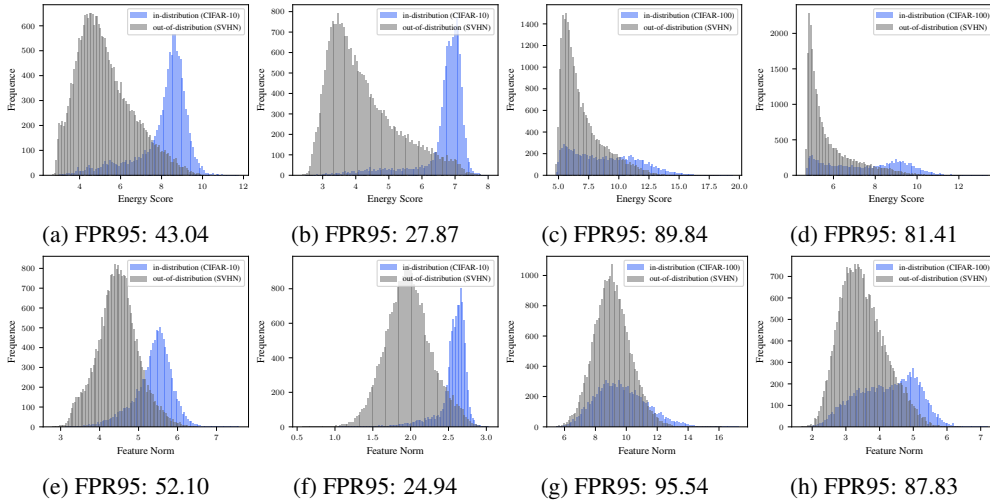


Figure 3: Distribution of energy scores (a-d) (Liu et al., 2020) and feature norms (e-h) from classification models trained without (a & c & e & g) or with (b & d & f & h) explicit feature regularization (EFR) ($\lambda = 1e - 5$). (a & b & e & f) and (c & d & g & h) are from ResNet-18 (He et al., 2016) trained on CIFAR-10 and from ResNet-34 trained on CIFAR-100, respectively. As can be seen, EFR can improve the performance of OOD detection by alleviating the over-confidence of OOD samples and making the energy scores of ID samples more concentrated. More intuitively, comparing (f) to (e) and (h) to (g), EFR effectively limits the growth of feature norms and improves the distinction between ID samples and OOD samples in the feature norm.

Zhou et al., 2022a). Explicit regularization on features can significantly remedy over-confidence and even improve generalization, which can be demonstrated in Tab. 1 and Fig. 3, where experiments on long-tailed classification (Zhong et al., 2021) and out-of-distribution (OOD) detection (Liu et al., 2020) are conducted, respectively. Moreover, adding explicit feature regularization can speed up the convergence of $\widehat{H}(t)$ to \widehat{WM} according to Theorem 4.1, as verified in Fig. 15.

4.3 RESCALED LEARNING RATE FOR THE SPHERICAL CONSTRAINED CASE

As shown in Sec. 3.3, the gradient of the objective in Eq. (3.15) with respect to $\mathbf{h}_{i,c}$ is $-\frac{1+\gamma}{CN\|\mathbf{h}_{i,c}\|_2}(\mathbf{I}_p - \widehat{\mathbf{h}}_{i,c}\widehat{\mathbf{h}}_{i,c}^\top)\mathbf{w}_c$, which shows that the gradient is significantly influenced by $\frac{1}{\|\mathbf{h}_{i,c}\|_2}$. For example, the gradient will disappear when $\|\mathbf{h}_{i,c}\|_2$ is too large, and the gradient will be too large when $\|\mathbf{h}_{i,c}\|_2$ is too small. A natural solution would be scaling up the learning rate $\eta(t)$ with the feature norm $\|\mathbf{h}_{i,c}\|_2$ for each sample⁷, and we can still guarantee the convergence in Theorem 3.5 if $\eta(t)$ is non-increasing and satisfies $\frac{\eta(0)(1+\gamma)}{CN} \leq \frac{1}{\|\mathbf{w}_c\|_2}$. As shown in Fig. 2c, the ASM loss with the rescaled learning rate (RLR) obviously converges faster than without it.

5 CONCLUSION

In this paper, we introduce the Averaged Sample Margin loss (ASM) as a surrogate to analyze the behavior of last-layer features and prototypes. Thanks to the conciseness of the ASM loss, we derive the exact dynamics under gradient descent in unconstrained, regularized, and spherical constrained cases. We then prove that these dynamics will converge exponentially to a particular solution relying on the initialization. Inspired by these results, we further provide some insights for improvements, such as the ASM loss with prototype-anchored learning, explicit feature regularization, and rescaled learning rate for spherical cases. We finally verify these theoretical results and insights with extensive experiments, including numerical analysis, visual classification, imbalanced learning, and out-of-distribution detection.

We expect that the ASM loss can be an excellent tool to help the community further understand the behavior of deep neural networks and beyond the scope of the paper.

⁷We can also implement this strategy by rescaling the loss, but the multiplicative feature norm stops gradients.

REFERENCES

- Peter Auer, Nicolo Cesa-Bianchi, and Claudio Gentile. Adaptive and self-confident on-line learning algorithms. *Journal of Computer and System Sciences*, 64(1):48–75, 2002.
- Sébastien Bubeck and Mark Sellke. A universal law of robustness via isoperimetry. *Advances in Neural Information Processing Systems*, 34:28811–28822, 2021.
- Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018.
- Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *Advances in neural information processing systems*, 32, 2019.
- Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9268–9277, 2019.
- Ahmet Demirkaya, Jiasi Chen, and Samet Oymak. Exploring the role of loss functions in multiclass classification. In *2020 54th annual conference on information sciences and systems (ciss)*, pages 1–5. IEEE, 2020.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training. *Proceedings of the National Academy of Sciences*, 118(43): e2103091118, 2021.
- Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learning. In *International Conference on Learning Representations*, 2021.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- XY Han, Vardan Papyan, and David L Donoho. Neural collapse under mse loss: Proximity to and dynamics on the central path. In *International Conference on Learning Representations*, 2022.
- Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- Like Hui and Mikhail Belkin. Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks. In *International Conference on Learning Representations*, 2021.
- Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding generalization in deep learning. *arXiv preprint arXiv:2202.08384*, 2022.
- Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J Su. An unconstrained layer-peeled perspective on neural collapse. In *International Conference on Learning Representations*, 2022.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Conference on Learning Theory*, pages 1772–1798. PMLR, 2019.
- Ziwei Ji, Miroslav Dudík, Robert E Schapire, and Matus Telgarsky. Gradient descent follows the regularization path for general losses. In *Conference on Learning Theory*, pages 2109–2136. PMLR, 2020.
- Tejaswi Kasarla, Gertjan J Burghouts, Max van Spengler, Elise van der Pol, Rita Cucchiara, and Pascal Mettes. Maximum class separation as inductive bias in one matrix. *arXiv preprint arXiv:2206.08704*, 2022.

- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Vladimir Koltchinskii and Dmitry Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. *The Annals of Statistics*, 30(1):1–50, 2002.
- Vignesh Kothapalli, Ebrahim Rasromani, and Vasudev Awatramani. Neural collapse: A review on modelling principles and generalization. *arXiv preprint arXiv:2206.04041*, 2022.
- A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Computer Science Department, University of Toronto, Tech. Rep.*, 1, 01 2009.
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 33:21464–21475, 2020.
- Weyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-margin softmax loss for convolutional neural networks. In *International Conference on Machine Learning*, pages 507–516. PMLR, 2016.
- Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017.
- Jianfeng Lu and Stefan Steinerberger. Neural collapse under cross-entropy loss. *Applied and Computational Harmonic Analysis*, 59:224–241, 2022.
- Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 75, 2021.
- Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *Sampling Theory, Signal Processing, and Data Analysis*, 20(2):1–13, 2022.
- Mor Shpigel Nacson, J. Lee, Suriya Gunasekar, Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *AISTATS*, 2019.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. In *International Conference on Learning Representations*, 2019.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. *Advances in neural information processing systems*, 30, 2017.
- Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.
- Ayush Sekhari, Karthik Sridharan, and Satyen Kale. Sgd: The role of implicit regularization, batch-size and multiple-epochs. *Advances In Neural Information Processing Systems*, 34:27422–27433, 2021.
- Ohad Shamir. Gradient methods never overfit on separable data. *J. Mach. Learn. Res.*, 22:85:1–85:20, 2021.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Tom Tirer and Joan Bruna. Extended unconstrained features model for exploring deep neural collapse. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 21478–21505. PMLR, 2022.
- Feng Wang, Xiang Xiang, Jian Cheng, and Alan Loddon Yuille. Normface: L2 hypersphere embedding for face verification. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1041–1049, 2017.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representation*, 2017.

- Zhisheng Zhong, Jiequan Cui, Shu Liu, and Jiaya Jia. Improving calibration for long-tailed recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16489–16498, 2021.
- Jinxin Zhou, Xiao Li, Tianyu Ding, Chong You, Qing Qu, and Zhihui Zhu. On the optimization landscape of neural collapse under MSE loss: Global optimality with unconstrained features. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 27179–27202. PMLR, 17–23 Jul 2022a.
- Xiong Zhou, Xianming Liu, Deming Zhai, Junjun Jiang, Xin Gao, and Xiangyang Ji. Prototype-anchored learning for learning with imperfect annotations. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 27245–27267. PMLR, 17–23 Jul 2022b.
- Xiong Zhou, Xianming Liu, Deming Zhai, Junjun Jiang, Xin Gao, and Xiangyang Ji. Learning towards the largest margins. In *International Conference on Learning Representations*, 2022c.
- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. *Advances in Neural Information Processing Systems*, 34:29820–29834, 2021.

Appendix for ‘‘Dynamics of The ASM Loss and Beyond’’

A PROOFS FOR LEMMAS, THEOREMS, PROPOSITIONS AND COROLLARIES

A.1 PROOF OF LEMMA 2.1

Lemma 2.1 (The Neural Collapse of The Averaged Sample Margin loss). *For norm-bounded prototypes and features, i.e., $\|\mathbf{w}_c\|_2 \leq E_1$ and $\|\mathbf{h}_{i,c}\|_2 \leq E_2$, $\forall i \in [N], \forall c \in [C]$, the global minimizer of $\frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C L_{ASM}(\mathbf{W}\mathbf{h}_{i,c}, y_{i,c})$ implies neural collapse when $p \geq C - 1$. More specifically, the global minimizer is uniquely obtained at $\frac{\mathbf{w}_i^\top \mathbf{w}_j}{\|\mathbf{w}_i\|_2 \|\mathbf{w}_j\|_2} = -\frac{1}{C-1}$, $\forall i \neq j$, $\frac{\mathbf{w}_{y_{i,c}}^\top \mathbf{h}_{i,c}}{\|\mathbf{w}_{y_{i,c}}\|_2 \|\mathbf{h}_{i,c}\|_2} = 1$, $\|\mathbf{w}_c\|_2 = E_1$, and $\|\mathbf{h}_{i,c}\|_2 = E_2$, $\forall i \in [N], \forall c \in [C]$.*

Proof. The proof is based on lower bounding the objective $\frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C L_{ASM}(\mathbf{W}\mathbf{h}_{i,c}, y_{i,c})$ by a sequence of inequalities that holds if and only if the solution forms Neural Collapse (Pappan et al., 2020). Let $\hat{\mathbf{w}} = \frac{1}{C} \sum_{c=1}^C \mathbf{w}_c$, according to the definition of L_{SM} , we have

$$\begin{aligned}
& \frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C L_{SM}(\mathbf{W}\mathbf{h}_{i,c}, y_{i,c}) \\
&= \frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C (-\mathbf{w}_{y_{i,c}}^\top \mathbf{h}_{i,c} + \gamma \sum_{j \neq y_{i,c}} \mathbf{w}_j^\top \mathbf{h}_{i,c}) \\
&\geq \frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C (-E_1 E_2 + \gamma (C\hat{\mathbf{w}} - \mathbf{w}_{y_{i,c}})^\top \mathbf{h}_{i,c}) \\
&\geq -\frac{\gamma E_2}{CN} \sum_{i=1}^N \sum_{c=1}^C \|C\hat{\mathbf{w}} - \mathbf{w}_{y_{i,c}}\|_2 - E_1 E_2, \\
&\geq -\gamma E_2 \sqrt{\frac{1}{CN} \sum_{i=1}^N \sum_{c=1}^C \|C\hat{\mathbf{w}} - \mathbf{w}_{y_{i,c}}\|_2^2 - E_1 E_2} \\
&= -\gamma E_2 \sqrt{\frac{1}{C} \sum_{c=1}^C \|\mathbf{w}_c\|_2^2 - C^2 \|\hat{\mathbf{w}}\|_2^2 - E_1 E_2} \\
&\geq -(1 + \gamma) E_1 E_2
\end{aligned}$$

where the first and second inequalities are based on the facts that $\mathbf{w}_{y_{i,c}}^\top \mathbf{h}_{i,c} \leq E_1 E_2$ and $(C\hat{\mathbf{w}} - \mathbf{w}_{y_{i,c}})^\top \mathbf{h}_{i,c} \geq -E_2 \|C\hat{\mathbf{w}} - \mathbf{w}_{y_{i,c}}\|_2$, respectively. In the third equality, we used the Cauchy-Schwarz inequality, and the last inequality we use the facts that $\|\mathbf{w}_c\|_2 \leq E_1$ and $\|\hat{\mathbf{w}}\|_2 \geq 0$.

According the above derivation, the equality holds if and only if $\forall i \in [N], \forall c \in [C]$, $\mathbf{w}_{y_{i,c}}^\top \mathbf{h}_{i,c} = E_1 E_2$, $(C\hat{\mathbf{w}} - \mathbf{w}_{y_{i,c}})^\top \mathbf{h}_{i,c} = -E_2 \|C\hat{\mathbf{w}} - \mathbf{w}_{y_{i,c}}\|_2$, $\|C\hat{\mathbf{w}} - \mathbf{w}_c\|_2 = \|C\hat{\mathbf{w}} - \mathbf{w}_C\|_2$, $\|\mathbf{w}_c\|_2 = E_1$, and $\|\hat{\mathbf{w}}\|_2 = 0$. These equations can be simplified as $\frac{\mathbf{w}_i^\top \mathbf{w}_j}{\|\mathbf{w}_i\|_2 \|\mathbf{w}_j\|_2} = -\frac{1}{C-1}$, $\forall i \neq j$, $\frac{\mathbf{w}_{y_{i,c}}^\top \mathbf{h}_{i,c}}{\|\mathbf{w}_{y_{i,c}}\|_2 \|\mathbf{h}_{i,c}\|_2} = 1$, $\|\mathbf{w}_c\|_2 = E_1$, and $\|\mathbf{h}_{i,c}\|_2 = E_2$, $\forall i \in [N], \forall c \in [C]$, which also implies neural collapse. \square

A.2 PROOF OF THEOREM 3.1

In this section, we will provide the proof of Theorem 3.1. Our analysis will actually rely on the eigenvalues and eigenspaces of five subspaces \mathcal{E}_1^+ , \mathcal{E}_1^- , \mathcal{E}_2^+ , \mathcal{E}_2^- and \mathcal{E}_3 in Theorem 3.1. Their concrete projection operator can be found in Appendix A.8. In the following, we show that these five subspaces are orthogonal:

Lemma A.1. *The following five subspaces are orthogonal to each other and satisfy $\mathbb{R}^{p \times C(N+1)} = \mathcal{E}_1^+ \oplus \mathcal{E}_1^- \oplus \mathcal{E}_2^+ \oplus \mathcal{E}_2^- \oplus \mathcal{E}_3$:*

$$\begin{aligned}\mathcal{E}_1^\epsilon &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}}(\mathbf{W} \otimes \mathbf{1}_N^\top), \mathbf{W}\mathbf{1}_C = 0, \mathbf{W} \in \mathbb{R}^{p \times C}\}, \\ \mathcal{E}_2^\epsilon &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}}\mathbf{h}\mathbf{1}_{CN}^\top, \mathbf{W} = \mathbf{h}\mathbf{1}_C^\top, \mathbf{h} \in \mathbb{R}^p\}, \\ \mathcal{E}_3 &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) = 0, \mathbf{W} = 0, \mathbf{H} \in \mathbb{R}^{p \times CN}\}.\end{aligned}\tag{A.1}$$

where $\epsilon \in \{\pm 1\}$, and $k \neq 0$.

Proof. For $(\mathbf{H}_1, \mathbf{W}_1) = (\frac{1}{\sqrt{N}}(\mathbf{W}_1 \otimes \mathbf{1}_N^\top), \mathbf{W}_1) \in \mathcal{E}_1^+$ and $(\mathbf{H}_2, \mathbf{W}_2) = (-\frac{1}{\sqrt{N}}(\mathbf{W}_2 \otimes \mathbf{1}_N^\top), \mathbf{W}_2) \in \mathcal{E}_1^-$, we have

$$\mathbf{H}_1\mathbf{H}_2^\top + \mathbf{W}_1\mathbf{W}_2^\top = -\frac{1}{N}(\mathbf{W}_1 \otimes \mathbf{1}_N^\top)(\mathbf{W}_2^\top \otimes \mathbf{1}_N) + \mathbf{W}_1\mathbf{W}_2^\top = 0.$$

For $(\mathbf{H}_1, \mathbf{W}_1) = (\frac{1}{\sqrt{N}}\mathbf{h}_1\mathbf{1}_{CN}^\top, \mathbf{h}_1\mathbf{1}_C^\top) \in \mathcal{E}_2^+$ and $(\mathbf{H}_2, \mathbf{W}_2) = (-\frac{1}{\sqrt{N}}\mathbf{h}_2\mathbf{1}_{CN}^\top, \mathbf{h}_2\mathbf{1}_C^\top) \in \mathcal{E}_2^-$, we have

$$\mathbf{H}_1\mathbf{H}_2^\top + \mathbf{W}_1\mathbf{W}_2^\top = -\frac{1}{N}\mathbf{h}_1\mathbf{1}_{CN}^\top\mathbf{1}_{CN}\mathbf{h}_2^\top + \mathbf{h}_1\mathbf{1}_C^\top\mathbf{1}_C\mathbf{h}_2^\top = 0.$$

For $(\mathbf{H}_1, \mathbf{W}_1) = (\epsilon_1\frac{1}{\sqrt{N}}(\mathbf{W}_1 \otimes \mathbf{1}_N^\top), \mathbf{W}_1) \in \mathcal{E}_1^\epsilon$ and $(\mathbf{H}_2, \mathbf{W}_2) = (\epsilon_2\frac{1}{\sqrt{N}}\mathbf{h}_2\mathbf{1}_{CN}^\top, \mathbf{h}_2\mathbf{1}_C^\top) \in \mathcal{E}_2^\epsilon$, since $\mathbf{W}_1\mathbf{1}_C = 0$ and $(\mathbf{W}_1 \otimes \mathbf{1}_N^\top)\mathbf{1}_{CN} = N\mathbf{W}_1\mathbf{1}_C = 0$, we have

$$\mathbf{H}_1\mathbf{H}_2^\top + \mathbf{W}_1\mathbf{W}_2^\top = \frac{\epsilon_1\epsilon_2}{N}(\mathbf{W}_1 \otimes \mathbf{1}_N^\top)\mathbf{1}_{CN}\mathbf{h}_2^\top + \mathbf{W}_1\mathbf{1}_C\mathbf{h}_2^\top = 0.$$

For $(\mathbf{H}_1, \mathbf{W}_1) = (\epsilon\frac{1}{\sqrt{N}}(\mathbf{W}_1 \otimes \mathbf{1}_N^\top), \mathbf{W}_1) \in \mathcal{E}_1^\epsilon$ and $(\mathbf{H}_2, \mathbf{W}_2) = (\mathbf{H}_2, 0) \in \mathcal{E}_3$, we have

$$\mathbf{H}_1\mathbf{H}_2^\top + \mathbf{W}_1\mathbf{W}_2^\top = \frac{\epsilon}{\sqrt{N}}(\mathbf{W}_1 \otimes \mathbf{1}_N^\top)\mathbf{H}_2^\top = \frac{\epsilon}{\sqrt{N}}\mathbf{W}_1(\mathbf{H}_2(\mathbf{I}_C \otimes \mathbf{1}_N))^\top = 0.$$

For $(\mathbf{H}_1, \mathbf{W}_1) = (\epsilon\frac{1}{\sqrt{N}}\mathbf{h}_1\mathbf{1}_{CN}^\top, \mathbf{h}_1\mathbf{1}_C^\top) \in \mathcal{E}_2^\epsilon$ and $(\mathbf{H}_2, \mathbf{W}_2) = (\mathbf{H}_2, 0) \in \mathcal{E}_3$, since $\mathbf{H}_2(\mathbf{I}_C \otimes \mathbf{1}_N) = 0$, we have

$$\mathbf{H}_1\mathbf{H}_2^\top + \mathbf{W}_1\mathbf{W}_2^\top = \frac{\epsilon}{\sqrt{N}}\mathbf{h}_1\mathbf{1}_{CN}^\top\mathbf{H}_2^\top = \frac{\epsilon}{\sqrt{N}}\mathbf{h}_1(\mathbf{H}_2(\mathbf{I}_C \otimes \mathbf{1}_N)\mathbf{1}_C^\top)^\top = 0$$

To sum up, we prove that the five subspaces $\mathcal{E}_1^+, \mathcal{E}_1^-, \mathcal{E}_2^+, \mathcal{E}_2^-, \mathcal{E}_3$ are orthogonal to each other. Moreover, we have

$$\dim \mathcal{E}_1^+ = \dim \mathcal{E}_1^- = p(C-1), \quad \dim \mathcal{E}_2^+ = \dim \mathcal{E}_2^- = p, \quad \dim \mathcal{E}_3 = pC(N-1).$$

Since these dimensions sum to $pC(N+1) = \dim(\mathbb{R}^{p \times CN} \oplus \mathbb{R}^{p \times C})$, then $\mathbb{R}^{p \times C(N+1)} = \mathcal{E}_1^+ \oplus \mathcal{E}_1^- \oplus \mathcal{E}_2^+ \oplus \mathcal{E}_2^- \oplus \mathcal{E}_3$. \square

Theorem 3.1 (Dynamics of Features, Prototypes and Biases without Constraints). *Consider the continual gradient flow (Equation 3.4) in which the dynamics follow the gradient descent direction of Sample Margin loss in Eq. (3.2). Let $\mathbf{Z}(t) = (\mathbf{H}(t), \mathbf{W}(t))$, if $\eta_1(t_1)\eta_2(t_2) = \eta_1(t_2)\eta_2(t_1)$ for any $t_1, t_2 \geq 0$, we have the following closed-form dynamics*

$$\begin{aligned}\mathbf{Z}(t) = & \Pi_1^+ \mathbf{Z}_0 (\alpha_1^+(t)\mathbf{C}(t) + \beta_1^+(t)\mathbf{I}_{C(N+1)}) + \Pi_1^- \mathbf{Z}_0 (\alpha_1^-(t)\mathbf{C}(t) + \beta_1^-(t)\mathbf{I}_{C(N+1)}) + \Pi_3 \mathbf{Z}_0 \\ & + \Pi_2^+ \mathbf{Z}_0 (\alpha_2^+(t)\mathbf{C}(t) + \beta_2^+(t)\mathbf{I}_{C(N+1)}) + \Pi_2^- \mathbf{Z}_0 (\alpha_2^-(t)\mathbf{C}(t) + \beta_2^-(t)\mathbf{I}_{C(N+1)}),\end{aligned}\tag{A.2}$$

and

$$\mathbf{b}(t) = \mathbf{b}_0 + \frac{1 + \gamma - \gamma C}{C} \int_0^t \eta_2(\tau) d\tau,\tag{A.3}$$

where $\alpha_1^\epsilon, \alpha_2^\epsilon, \beta_1^\epsilon$ and β_2^ϵ for $\epsilon \in \{\pm\}$ are the scalars that only depend on C, N, γ, η_1 and η_2 (where the detailed forms can be seen in the appendix), $\mathbf{Z}_0 = (\mathbf{H}_0, \mathbf{W}_0)$, $\mathbf{C}(t) =$

$\begin{pmatrix} \int_0^t \eta_1(\tau) d\tau \mathbf{I}_{CN} & 0 \\ 0 & \int_0^t \eta_2(\tau) d\tau \mathbf{I}_C \end{pmatrix}$, Π_1^ϵ , Π_2^ϵ and Π_3 for $\epsilon \in \{\pm\}$ are orthogonal projection operators onto the following respective eigenspaces:

$$\begin{aligned} \mathcal{E}_1^\epsilon &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}}(\mathbf{W} \otimes \mathbf{1}_N^\top), \mathbf{W}\mathbf{1}_C = 0\}, \\ \mathcal{E}_2^\epsilon &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}}\mathbf{h}\mathbf{1}_{CN}^\top, \mathbf{W} = \mathbf{h}\mathbf{1}_C^\top, \mathbf{h} \in \mathbb{R}^p\}, \\ \mathcal{E}_3 &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) = 0, \mathbf{W} = 0\}. \end{aligned} \quad (\text{A.4})$$

Proof. Writing $\mathbf{Z}(t) = (\mathbf{H}(t), \mathbf{W}(t))$, then the unsolved portion of the system is given by

$$\mathbf{Z}'(t) = \mathbf{Z}(t) \begin{pmatrix} 0 & \mathbf{M}^\top \\ \mathbf{M} & 0 \end{pmatrix} \begin{pmatrix} \eta_1(t)\mathbf{I}_{CN} & 0 \\ 0 & \eta_2(t)\mathbf{I}_C \end{pmatrix}, \quad (\text{A.5})$$

where $\mathbf{M} = \frac{1}{CN}((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{1}_C\mathbf{1}_{CN}^\top)$.

Let $\mathbf{A}(t) = \begin{pmatrix} 0 & \mathbf{M}^\top \\ \mathbf{M} & 0 \end{pmatrix} \begin{pmatrix} \eta_1(t)\mathbf{I}_{CN} & 0 \\ 0 & \eta_2(t)\mathbf{I}_C \end{pmatrix} = \begin{pmatrix} 0 & \eta_2(t)\mathbf{M}^\top \\ \eta_1(t)\mathbf{M} & 0 \end{pmatrix}$, then the equation above can be reformulated as the initial-value problem associated with the linear ordinary differential equation:

$$\mathbf{Z}'(t) = \mathbf{Z}(t)\mathbf{A}(t), \quad \mathbf{Z}(0) = \mathbf{Z}_0. \quad (\text{A.6})$$

For any t_1, t_2 , we have the matrix commutator of $\mathbf{A}(t_1)$ and $\mathbf{A}(t_2)$

$$\begin{aligned} &[\mathbf{A}(t_1), \mathbf{A}(t_2)] \\ &= \mathbf{A}(t_1)\mathbf{A}(t_2) - \mathbf{A}(t_2)\mathbf{A}(t_1) \\ &= \begin{pmatrix} 0 & \eta_2(t_1)\mathbf{M}^\top \\ \eta_1(t_1)\mathbf{M} & 0 \end{pmatrix} \begin{pmatrix} 0 & \eta_2(t_2)\mathbf{M}^\top \\ \eta_1(t_2)\mathbf{M} & 0 \end{pmatrix} - \mathbf{A}(t_2)\mathbf{A}(t_1) \\ &= \begin{pmatrix} \eta_2(t_1)\eta_1(t_2)\mathbf{M}^\top\mathbf{M} & 0 \\ 0 & \eta_1(t_1)\eta_2(t_2)\mathbf{M}\mathbf{M}^\top \end{pmatrix} - \mathbf{A}(t_2)\mathbf{A}(t_1) \\ &= \begin{pmatrix} (\eta_2(t_1)\eta_1(t_2) - \eta_2(t_2)\eta_1(t_1))\mathbf{M}^\top\mathbf{M} & 0 \\ 0 & (\eta_1(t_1)\eta_2(t_2) - \eta_2(t_1)\eta_1(t_2))\mathbf{M}\mathbf{M}^\top \end{pmatrix} \\ &= 0 \end{aligned}$$

where the last equality is based on the fact that $\eta_2(t_1)\eta_1(t_2) = \eta_2(t_2)\eta_1(t_1)$. Therefore, according to Magnus approach, we have

$$\mathbf{Z}(t) = \mathbf{Z}_0 \exp\left(\int_0^t \mathbf{A}(\tau) d\tau\right) = \mathbf{Z}_0 \exp\begin{pmatrix} 0 & \int_0^t \eta_2(\tau) d\tau \mathbf{M}^\top \\ \int_0^t \eta_1(\tau) d\tau \mathbf{M} & 0 \end{pmatrix}. \quad (\text{A.7})$$

Let $\zeta_1(t) = \int_0^t \eta_1(\tau) d\tau$, $\zeta_2(t) = \int_0^t \eta_2(\tau) d\tau$, $\mathbf{B} = \begin{pmatrix} 0 & \mathbf{M}^\top \\ \mathbf{M} & 0 \end{pmatrix}$, $\mathbf{C}(t) = \begin{pmatrix} \zeta_1(t)\mathbf{I}_{CN} & 0 \\ 0 & \zeta_2(t)\mathbf{I}_C \end{pmatrix}$, and $\mathbf{L}(t) = \mathbf{B}\mathbf{C}(t)$, we have

$$\mathbf{Z}(t) = \mathbf{Z}_0 \exp(\mathbf{L}(t)) = \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\mathbf{L}(t))^k}{k!}. \quad (\text{A.8})$$

Moreover, we have

$$\begin{aligned} (\mathbf{L}(t))^2 &= \begin{pmatrix} 0 & \zeta_2(t)\mathbf{M}^\top \\ \zeta_1(t)\mathbf{M} & 0 \end{pmatrix} \begin{pmatrix} 0 & \zeta_2(t)\mathbf{M}^\top \\ \zeta_1(t)\mathbf{M} & 0 \end{pmatrix} \\ &= \begin{pmatrix} \zeta_1(t)\zeta_2(t)\mathbf{M}^\top\mathbf{M} & 0 \\ 0 & \zeta_1(t)\zeta_2(t)\mathbf{M}\mathbf{M}^\top \end{pmatrix} \\ &= \zeta_1(t)\zeta_2(t)\mathbf{B}^2, \end{aligned} \quad (\text{A.9})$$

thus we obtain

$$\begin{aligned}
\mathbf{Z}(t) &= \mathbf{Z}_0 \left(\sum_{k=0}^{\infty} \frac{(\mathbf{L}(t))^{2k+1}}{(2k+1)!} + \sum_{k=0}^{\infty} \frac{(\mathbf{L}(t))^{2k}}{(2k)!} \right) \\
&= \mathbf{Z}_0 \left(\sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k+1} \mathbf{C}(t)}{(2k+1)!} + \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k}}{(2k)!} \right) \\
&= \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k+1} \mathbf{C}(t)}{(2k+1)!} + \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k}}{(2k)!}
\end{aligned} \tag{A.10}$$

Looking at the above equation, we just need to analyze the eigenspaces and eigenvalues of \mathbf{B} .

Considering the following five subspaces:

$$\begin{aligned}
\mathcal{E}_1^\epsilon &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}}(\mathbf{W} \otimes \mathbf{1}_N^\top), \mathbf{W}\mathbf{1}_C = 0\}, \\
\mathcal{E}_2^\epsilon &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H} = \epsilon \cdot \frac{1}{\sqrt{N}}\mathbf{h}\mathbf{1}_{CN}^\top, \mathbf{W} = \mathbf{h}\mathbf{1}_C^\top, \mathbf{h} \in \mathbb{R}^P\}, \\
\mathcal{E}_3 &:= \{(\mathbf{H}, \mathbf{W}) : \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) = 0, \mathbf{W} = 0\}.
\end{aligned} \tag{A.11}$$

where $\epsilon \in \{\pm\}$. According to Lemma A.1, these five subspaces are orthogonal to each other and satisfy $\mathbb{R}^{p \times C(N+1)} = \mathcal{E}_1^+ \oplus \mathcal{E}_1^- \oplus \mathcal{E}_2^+ \oplus \mathcal{E}_2^- \oplus \mathcal{E}_3$.

In the following, we will prove that \mathcal{E}_1^ϵ , \mathcal{E}_2^ϵ and \mathcal{E}_3 are five eigenspaces of \mathbf{B} . More specifically, each nonzero member of each claimed eigenspace is an eigenvector, and the claimed eigenspaces have distinct eigenvalues.

Note that for $(\mathbf{H}, \mathbf{W}) \in \mathbb{R}^{p \times CN} \oplus \mathbb{R}^{p \times C}$, we have $(\mathbf{H}, \mathbf{W})\mathbf{B} = (\mathbf{W}\mathbf{M}^\top, \mathbf{H}\mathbf{W})$.

For $(\mathbf{H}, \mathbf{W}) \in \mathcal{E}_1^\epsilon$, we have $\mathbf{H} = \frac{\epsilon}{\sqrt{N}}\mathbf{W} \otimes \mathbf{1}_N^\top$ and $\mathbf{W}\mathbf{1}_C = 0$, thus

$$\begin{aligned}
\mathbf{W}\mathbf{M} &= \frac{1}{CN}\mathbf{W}((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{1}_C\mathbf{1}_{CN}^\top) \\
&= \frac{(1+\gamma)}{CN}\mathbf{W} \otimes \mathbf{1}_N^\top \\
&= \frac{\epsilon(1+\gamma)}{C\sqrt{N}}\mathbf{H}, \\
\mathbf{H}\mathbf{M}^\top &= \frac{1}{CN}[\epsilon \cdot \frac{1}{\sqrt{N}}(\mathbf{W} \otimes \mathbf{1}_N^\top)][(1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{1}_C\mathbf{1}_{CN}^\top]^\top \\
&= \frac{\epsilon}{CN\sqrt{N}}[(\mathbf{W} \otimes \mathbf{1}_N^\top)((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N) - \gamma(\mathbf{1}_C \otimes \mathbf{1}_N)\mathbf{1}_C^\top)] \\
&= \frac{\epsilon(1+\gamma)}{C\sqrt{N}}\mathbf{W},
\end{aligned}$$

i.e., (\mathbf{H}, \mathbf{W}) is an eigenvector of \mathbf{B} with eigenvalue $\frac{\epsilon(1+\gamma)}{C\sqrt{N}}$.

For $(\mathbf{H}, \mathbf{W}) \in \mathcal{E}_2^\epsilon$, we have $\mathbf{H} = \frac{\epsilon}{\sqrt{N}}\mathbf{h}\mathbf{1}_{CN}^\top$ and $\mathbf{W} = \mathbf{h}\mathbf{1}_C^\top$, thus

$$\begin{aligned}
\mathbf{W}\mathbf{M} &= \frac{1}{CN}\mathbf{h}\mathbf{1}_C^\top((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{1}_C\mathbf{1}_{CN}^\top) \\
&= \frac{1}{CN}\mathbf{h}((1+\gamma)\mathbf{1}_C(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{C}\mathbf{1}_{CN}^\top) \\
&= \frac{(1+\gamma-\gamma C)}{CN}\mathbf{h}\mathbf{1}_{CN}^\top \\
&= \frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}}\mathbf{H}, \\
\mathbf{H}\mathbf{M}^\top &= \frac{1}{CN\sqrt{N}}(\epsilon \cdot \mathbf{h}\mathbf{1}_{CN}^\top)((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{1}_C\mathbf{1}_{CN}^\top)^\top \\
&= \frac{\epsilon}{CN}\mathbf{h}((1+\gamma)\mathbf{1}_{CN}^\top(\mathbf{I}_C \otimes \mathbf{1}_N) - \gamma\mathbf{1}_{CN\sqrt{N}}^\top\mathbf{1}_{CN}\mathbf{1}_C^\top) \\
&= \frac{\epsilon}{CN\sqrt{N}}\mathbf{h}((1+\gamma)N\mathbf{1}_C^\top - \gamma CN\mathbf{1}_C^\top) \\
&= \frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}}\mathbf{h}\mathbf{1}_C^\top \\
&= \frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}}\mathbf{W},
\end{aligned}$$

i.e., (\mathbf{H}, \mathbf{W}) is an eigenvector of \mathbf{B} with eigenvalue $\frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}}$.

For $(\mathbf{H}, \mathbf{W}) \in \mathcal{E}_3$, we have $\mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) = 0$ and $\mathbf{W} = 0$, thus

$$\begin{aligned} \mathbf{W}\mathbf{M} &= \frac{1}{C\sqrt{N}} \cdot 0((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{1}_C\mathbf{1}_{CN}^\top) = 0, \\ \mathbf{H}\mathbf{M}^\top &= \frac{1}{C\sqrt{N}}\mathbf{H}((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N^\top) - \gamma\mathbf{1}_C\mathbf{1}_{CN}^\top)^\top \\ &= \frac{1}{C\sqrt{N}}\mathbf{H}((1+\gamma)(\mathbf{I}_C \otimes \mathbf{1}_N) - \gamma\mathbf{1}_{CN}\mathbf{1}_C^\top) \\ &= -\frac{\gamma}{C\sqrt{N}}\mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N)\mathbf{1}_C\mathbf{1}_C^\top \\ &= 0 \end{aligned}$$

i.e., (\mathbf{H}, \mathbf{W}) is an eigenvector of \mathbf{B} with eigenvalue 0.

Overall, letting Π_i^ϵ denote orthogonal projection onto \mathcal{E}_i^ϵ , we have the spectral decomposition

$$\mathbf{B} = \frac{1}{C\sqrt{N}} [(1+\gamma)(\Pi_1^+ - \Pi_1^-) + (1+\gamma - \gamma C)(\Pi_2^+ - \Pi_2^-)]. \quad (\text{A.12})$$

We then provide the concrete formulation of $\mathbf{Z}(t) = \mathbf{Z}_0 \exp(\mathbf{L}(t))$ by the orthogonal projection of \mathbf{Z}_0 onto each eigenspace of \mathbf{B} , *i.e.*,

$$\mathbf{Z}_0 = \Pi_1^+ \mathbf{Z}_0 + \Pi_1^- \mathbf{Z}_0 + \Pi_2^+ \mathbf{Z}_0 + \Pi_2^- \mathbf{Z}_0 + \Pi_3 \mathbf{Z}_0.$$

Decomposition along $\Pi_1^\epsilon \mathbf{Z}_0$. First, $\Pi_1^\epsilon \mathbf{Z}_0 \mathbf{B} = \frac{\epsilon(1+\gamma)}{C\sqrt{N}} \Pi_1^\epsilon \mathbf{Z}_0$, so $\Pi_1^\epsilon \mathbf{Z}_0 \mathbf{B}^k = \left(\frac{\epsilon(1+\gamma)}{C\sqrt{N}}\right)^k \Pi_1^\epsilon \mathbf{Z}_0$ for $k \geq 0$, then

$$\begin{aligned} &\Pi_1^\epsilon \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k+1} \mathbf{C}(t)}{(2k+1)!} \\ &= \Pi_1^\epsilon \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \left(\frac{\epsilon(1+\gamma)}{C\sqrt{N}}\right)^{2k+1} \mathbf{C}(t)}{(2k+1)!} \\ &= \frac{\Pi_1^\epsilon \mathbf{Z}_0 \mathbf{C}(t)}{\sqrt{\zeta_1(t)\zeta_2(t)}} \sum_{k=0}^{\infty} \frac{\left(\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right)^{2k+1}}{(2k+1)!} \\ &= \frac{\Pi_1^\epsilon \mathbf{Z}_0 \mathbf{C}(t)}{2\sqrt{\zeta_1(t)\zeta_2(t)}} \left(e^{\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} - e^{-\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} \right), \end{aligned} \quad (\text{A.13})$$

and

$$\begin{aligned} &\Pi_1^\epsilon \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k}}{(2k)!} \\ &= \Pi_1^\epsilon \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \left(\frac{\epsilon(1+\gamma)}{C\sqrt{N}}\right)^{2k}}{(2k)!} \\ &= \frac{\Pi_1^\epsilon \mathbf{Z}_0}{2} \left(e^{\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} + e^{-\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} \right), \end{aligned} \quad (\text{A.14})$$

which is based on the facts that $\frac{e^x - e^{-x}}{2} = \sum_{k=0}^{\infty} \frac{x^{2k+1}}{(2k+1)!}$ and $\frac{e^x + e^{-x}}{2} = \sum_{k=0}^{\infty} \frac{x^{2k}}{(2k)!}$. Thus we have

$$\Pi_1^\epsilon \mathbf{Z}_0 \exp(\mathbf{L}(t)) = \Pi_1^\epsilon \mathbf{Z}_0 (\alpha_1^\epsilon \mathbf{C}(t) + \beta_1^\epsilon \mathbf{I}_{C(N+1)}), \quad (\text{A.15})$$

with

$$\begin{aligned} \alpha_1^\epsilon(t) &= \frac{\exp\left(\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right) - \exp\left(-\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right)}{2\sqrt{\zeta_1(t)\zeta_2(t)}}, \\ \beta_1^\epsilon(t) &= \frac{\exp\left(\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right) + \exp\left(-\frac{\epsilon(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right)}{2}. \end{aligned} \quad (\text{A.16})$$

Decomposition along $\Pi_2^\epsilon \mathbf{Z}_0$. Similarly, for $\Pi_2^\epsilon \mathbf{Z}_0 \mathbf{B} = \frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}} \Pi_2^\epsilon$, we have

$$\Pi_2^\epsilon \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k+1} \mathbf{C}(t)}{(2k+1)!} = \frac{\Pi_2^\epsilon \mathbf{Z}_0 \mathbf{C}(t)}{2\sqrt{\zeta_1(t)\zeta_2(t)}} \left(e^{\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} - e^{-\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} \right),$$

and

$$\Pi_2^\epsilon \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k}}{(2k)!} = \frac{\Pi_2^\epsilon \mathbf{Z}_0}{2} \left(e^{\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} + e^{-\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} \right).$$

Thus we have

$$\Pi_2^\epsilon \mathbf{Z}_0 \exp(\mathbf{L}(t)) = \Pi_2^\epsilon \mathbf{Z}_0 (\alpha_2^\epsilon \mathbf{C}(t) + \beta_2^\epsilon \mathbf{I}_{C(N+1)}), \quad (\text{A.17})$$

with

$$\begin{aligned} \alpha_2^\epsilon(t) &= \frac{\exp\left(\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right) - \exp\left(-\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right)}{2\sqrt{\zeta_1(t)\zeta_2(t)}}, \\ \beta_2^\epsilon(t) &= \frac{\exp\left(\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right) + \exp\left(-\frac{\epsilon(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}\right)}{2}. \end{aligned} \quad (\text{A.18})$$

Decomposition along $\Pi_3 \mathbf{Z}_0$. Since each vector in \mathcal{E}_3 is a eigenvector of \mathbf{B} with eigenvalue 0, then we have

$$\Pi_3 \mathbf{Z}_0 \left(\sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k+1} \mathbf{C}(t)}{(2k+1)!} + \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k}}{(2k)!} \right) = \Pi_3 \mathbf{Z}_0 \quad (\text{A.19})$$

Note that \mathcal{E}_1^+ , \mathcal{E}_1^- , \mathcal{E}_2^+ , \mathcal{E}_2^- and \mathcal{E}_3 are orthogonal subspace of $\mathbb{R}^{p \times CN} \oplus \mathbb{R}^{C \times p}$, thus

$$\begin{aligned} \mathbf{Z}(t) &= \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k+1} \mathbf{C}(t)}{(2k+1)!} + \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\zeta_1(t)\zeta_2(t))^k \mathbf{B}^{2k}}{(2k)!} \\ &= (\Pi_1^+ \mathbf{Z}_0 + \Pi_1^- \mathbf{Z}_0 + \Pi_2^+ \mathbf{Z}_0 + \Pi_2^- \mathbf{Z}_0 + \Pi_3 \mathbf{Z}_0) \exp(\mathbf{L}(t)) \\ &= \Pi_1^+ \mathbf{Z}_0 (\alpha_1^+(t) \mathbf{C}(t) + \beta_1^+(t) \mathbf{I}_{C(N+1)}) + \Pi_1^- \mathbf{Z}_0 (\alpha_1^-(t) \mathbf{C}(t) + \beta_1^-(t) \mathbf{I}_{C(N+1)}) + \Pi_3 \mathbf{Z}_0 \\ &\quad + \Pi_2^+ \mathbf{Z}_0 (\alpha_2^+(t) \mathbf{C}(t) + \beta_2^+(t) \mathbf{I}_{C(N+1)}) + \Pi_2^- \mathbf{Z}_0 (\alpha_2^-(t) \mathbf{C}(t) + \beta_2^-(t) \mathbf{I}_{C(N+1)}) \end{aligned} \quad (\text{A.20})$$

Moreover, since $\mathbf{b}'(t) = -\eta_2(t) \frac{\gamma C - \gamma - 1}{C} \mathbf{1}_C$, we obtain

$$\mathbf{b}(t) = \mathbf{b}(0) + \int_0^t -\eta_2(\tau) \frac{\gamma C - \gamma - 1}{C} d\tau \mathbf{1}_C = \mathbf{b}_0 + \frac{(1 + \gamma - \gamma C)\zeta_2(t)}{C} \mathbf{1}_C, \quad (\text{A.21})$$

with $\zeta_2(t) = \int_0^t \eta_2(\tau) d\tau$. \square

A.3 PROOF OF COROLLARY 3.2

Corollary 3.2 Under the conditions and notation of Theorem 3.1, let $s = \frac{\eta_1(0)}{\eta_2(0)}$, if $0 < \gamma < \frac{2}{C-2}$ (where $C > 2$) or $C = 2$, and $\lim_{t \rightarrow \infty} \zeta_1(t) = \infty$, then the gradient flow (as in Eq. (3.4)) will behave as:

$$e^{-\frac{(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}} \mathbf{Z}(t) = \left(\frac{1+\sqrt{s}}{2} \mathbf{H}_1^+ + \frac{1-\sqrt{s}}{2} \mathbf{H}_1^-, \frac{1+\sqrt{s}}{2\sqrt{s}} \mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}} \mathbf{W}_1^- \right) + \mathbf{\Delta}(t), \quad (\text{A.22})$$

where $(\mathbf{H}_1^+, \mathbf{W}_1^+) = \Pi_1^+ \mathbf{Z}_0$, $(\mathbf{H}_1^-, \mathbf{W}_1^-) = \Pi_1^- \mathbf{Z}_0$, and the residual term $\mathbf{\Delta}(t)$ decreases at least as $\|\mathbf{\Delta}(t)\| = O\left(e^{\frac{\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}} \cdot \max\{-\gamma C, (C-2)\gamma-2\}}\right)$, and so let $\bar{\mathbf{Z}} = \left(\frac{1+\sqrt{s}}{2} \mathbf{H}_1^+ + \frac{1-\sqrt{s}}{2} \mathbf{H}_1^-, \frac{1+\sqrt{s}}{2\sqrt{s}} \mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}} \mathbf{W}_1^-\right)$, the normalized $\mathbf{Z}(t)$ converges to the normalized $\bar{\mathbf{Z}}$ in

$$\left\| \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} - \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \right\| = O\left(e^{\frac{\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}} \cdot \max\{-\gamma C, (C-2)\gamma-2\}}\right), \quad (\text{A.23})$$

which further indicates $\lim_{t \rightarrow \infty} \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} \in \mathcal{E}$. Moreover, if $\gamma \neq \frac{1}{C-1}$, then $\lim_{t \rightarrow \infty} \frac{\max_i b_i(t)}{\min_i b_i(t)} = 1$.

Proof. Let $(\mathbf{H}_1^\epsilon, \mathbf{W}_1^\epsilon) = \Pi_1^\epsilon \mathbf{Z}_0$, $(\mathbf{H}_2^\epsilon, \mathbf{W}_2^\epsilon) = \Pi_2^\epsilon \mathbf{Z}_0$ and $(\mathbf{H}_3, \mathbf{W}_3) = \Pi_3 \mathbf{Z}_0$ for $\epsilon \in \{\pm 1\}$, according to Theorem 3.1, we have

$$\begin{aligned}\mathbf{H}(t) &= \sum_{\substack{i \in \{1,2\} \\ \epsilon \in \{\pm\}}} (\alpha_i^\epsilon(t) \zeta_1(t) + \beta_i^\epsilon(t)) \mathbf{H}_i^\epsilon + \mathbf{H}_3, \\ \mathbf{W}(t) &= \sum_{\substack{i \in \{1,2\} \\ \epsilon \in \{\pm\}}} (\alpha_i^\epsilon(t) \zeta_2(t) + \beta_i^\epsilon(t)) \mathbf{W}_i^\epsilon + \mathbf{W}_3.\end{aligned}\tag{A.24}$$

Since $\eta_1(t_1)\eta_2(t_2) = \eta_1(t_2)\eta_2(t_1)$, then $\eta_1(t) = \frac{\eta_1(0)}{\eta_2(0)}\eta_2(t)$. Let $s = \frac{\eta_1(0)}{\eta_2(0)}$, $p(t) = \frac{(1+\gamma)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}$, and $q(t) = \frac{(1+\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}}$, we have $\zeta_1(t) = \int_0^t \eta_1(\tau) d\tau = s \int_0^t \eta_2(\tau) d\tau = s\zeta_2(t)$, and then

$$\begin{aligned}\alpha_1^\epsilon(t)\zeta_1(t) + \beta_1^\epsilon(t) &= \frac{1+\sqrt{s}}{2}e^{\epsilon p(t)} + \frac{1-\sqrt{s}}{2}e^{-\epsilon p(t)} = \frac{1+\epsilon\sqrt{s}}{2}e^{\epsilon p(t)} + O(e^{-p(t)}), \\ \alpha_1^\epsilon(t)\zeta_2(t) + \beta_1^\epsilon(t) &= \frac{1+\sqrt{s}}{2\sqrt{s}}e^{\epsilon p(t)} - \frac{1-\sqrt{s}}{2\sqrt{s}}e^{-\epsilon p(t)} = \frac{\epsilon+\sqrt{s}}{2\sqrt{s}}e^{\epsilon p(t)} + O(e^{-p(t)}), \\ \alpha_2^\epsilon(t)\zeta_1(t) + \beta_2^\epsilon(t) &= \frac{1+\sqrt{s}}{2}e^{\epsilon q(t)} + \frac{1-\sqrt{s}}{2}e^{-\epsilon q(t)} = \frac{1+\epsilon\sqrt{s}}{2}e^{\epsilon q(t)} + O(e^{-q(t)}), \\ \alpha_2^\epsilon(t)\zeta_2(t) + \beta_2^\epsilon(t) &= \frac{1+\sqrt{s}}{2\sqrt{s}}e^{\epsilon q(t)} - \frac{1-\sqrt{s}}{2\sqrt{s}}e^{-\epsilon q(t)} = \frac{\epsilon+\sqrt{s}}{2\sqrt{s}}e^{\epsilon q(t)} + O(e^{-q(t)}).\end{aligned}$$

Since $0 < \gamma < \frac{2}{C-2}$ (where $C > 2$) or $C = 2$, then we have $p(t) - q(t) = \frac{\gamma C \sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}} > 0$, $p(t) + q(t) = \frac{(2+2\gamma-\gamma C)\sqrt{\zeta_1(t)\zeta_2(t)}}{C\sqrt{N}} > 0$, and substitute these results into Equation (A.24) to obtain

$$\begin{aligned}e^{-p(t)}\mathbf{H}(t) &= \frac{1+\sqrt{s}}{2}\mathbf{H}_1^+ + \frac{1-\sqrt{s}}{2}\mathbf{H}_1^- + \mathbf{\Delta}_1(t), \\ e^{-p(t)}\mathbf{W}(t) &= \frac{1+\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^- + \mathbf{\Delta}_2(t),\end{aligned}\tag{A.25}$$

where $\|\mathbf{\Delta}_1(t)\| = O(e^{\max\{q(t)-p(t), -q(t)-p(t)\}})$ and $\|\mathbf{\Delta}_2(t)\| = O(e^{\max\{q(t)-p(t), -q(t)-p(t)\}})$. Therefore, we have

$$e^{-p(t)}\mathbf{Z}(t) = \left(\frac{1+\sqrt{s}}{2}\mathbf{H}_1^+ + \frac{1-\sqrt{s}}{2}\mathbf{H}_1^-, \frac{1+\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^- \right) + \mathbf{\Delta}(t),\tag{A.26}$$

where $\mathbf{\Delta}(t) = (\mathbf{\Delta}_1(t), \mathbf{\Delta}_2(t))$ and $\|\mathbf{\Delta}(t)\| \leq \|\mathbf{\Delta}_1(t)\| + \|\mathbf{\Delta}_2(t)\| = O(e^{\max\{q(t)-p(t), -q(t)-p(t)\}})$.

Let $\bar{\mathbf{Z}} = \left(\frac{1+\sqrt{s}}{2}\mathbf{H}_1^+ + \frac{1-\sqrt{s}}{2}\mathbf{H}_1^-, \frac{1+\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^- \right)$, we have

$$\left\| \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} - \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \right\| = \left\| \frac{\bar{\mathbf{Z}} + \mathbf{\Delta}(t)}{\|\bar{\mathbf{Z}} + \mathbf{\Delta}(t)\|} - \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \right\| \leq \frac{2\|\bar{\mathbf{Z}}\|\|\mathbf{\Delta}(t)\|}{\|\bar{\mathbf{Z}} + \mathbf{\Delta}(t)\|\|\bar{\mathbf{Z}}\|} = \frac{2\|\mathbf{\Delta}(t)\|}{\|\bar{\mathbf{Z}} + \mathbf{\Delta}(t)\|},\tag{A.27}$$

thus $\left\| \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} - \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \right\| = O(e^{\max\{q(t)-p(t), -q(t)-p(t)\}})$, and further $\lim_{t \rightarrow \infty} \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} = \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|}$ when $\lim_{t \rightarrow \infty} \zeta_1(t) = \infty$.

According to the definition of \mathcal{E}_1 and \mathcal{E}_2 , we have

$$\bar{\mathbf{Z}} = \left(\frac{\sqrt{s}}{\sqrt{N}} \left(\frac{1+\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^- \right) \otimes \mathbf{1}_N^\top, \frac{1+\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^+ - \frac{1-\sqrt{s}}{2\sqrt{s}}\mathbf{W}_1^- \right),\tag{A.28}$$

thus we have $\bar{\mathbf{Z}} \in \mathcal{E}$, then $\lim_{t \rightarrow \infty} \frac{\mathbf{Z}(t)}{\|\mathbf{Z}(t)\|} = \frac{\bar{\mathbf{Z}}}{\|\bar{\mathbf{Z}}\|} \in \mathcal{E}$.

Moreover, we have $\mathbf{b}(t) = \mathbf{b}_0 + \frac{(1+\gamma-\gamma C)\zeta_2(t)}{C}\mathbf{1}_C$, then $\forall i, j$,

$$\lim_{t \rightarrow \infty} \frac{b_i(t)}{b_j(t)} = \lim_{t \rightarrow \infty} \frac{b_i(0) + \frac{1+\gamma-\gamma C}{C}\zeta_2(t)}{b_j(0) + \frac{1+\gamma-\gamma C}{C}\zeta_2(t)} = 1,$$

thus $\lim_{t \rightarrow \infty} \frac{\max_i b_i(t)}{\max_i b_i(t)} = 1$. \square

A.4 PROOF OF THEOREM 3.3

Theorem 3.3 (Dynamics of Features and Prototypes with ℓ_2 Regularization) *Under the conditions and notations of Theorem 3.1, consider the continual gradient flow (Equation 3.11) in which the dynamics follow the gradient descent direction of Sample Margin loss in Eq. (3.10). Let $\mathbf{Z}(t) = (\mathbf{H}(t), \mathbf{W}(t))$, if $\eta_1(t_1)\eta_2(t_2) = \eta_1(t_2)\eta_2(t_1)$ for any $t_1, t_2 \geq 0$, we have the following closed-form dynamics*

$$\begin{aligned} \mathbf{Z}(t) = & \Pi_1^+ \mathbf{Z}_0 \begin{pmatrix} a_1^+(t) \mathbf{I}_{CN} & 0 \\ 0 & b_1^+(t) \mathbf{I}_C \end{pmatrix} + \Pi_1^- \mathbf{Z}_0 \begin{pmatrix} a_1^-(t) \mathbf{I}_{CN} & 0 \\ 0 & b_1^-(t) \mathbf{I}_C \end{pmatrix} + \\ & \Pi_2^+ \mathbf{Z}_0 \begin{pmatrix} a_2^+(t) \mathbf{I}_{CN} & 0 \\ 0 & b_2^+(t) \mathbf{I}_C \end{pmatrix} + \Pi_2^- \mathbf{Z}_0 \begin{pmatrix} a_2^-(t) \mathbf{I}_{CN} & 0 \\ 0 & b_2^-(t) \mathbf{I}_C \end{pmatrix} + \\ & \Pi_3 \mathbf{Z}_0 \begin{pmatrix} a_3(t) \mathbf{I}_{CN} & 0 \\ 0 & b_3(t) \mathbf{I}_C \end{pmatrix} \end{aligned} \quad (\text{A.29})$$

$$\mathbf{Z}(t) = \mathbf{Z}(t) \begin{pmatrix} \phi_1(t) \mathbf{I}_{CN} & 0 \\ 0 & \phi_2(t) \mathbf{I}_C \end{pmatrix}, \text{ and } \mathbf{b}(t) = \phi_3(t) \left(\mathbf{b}(0) - \frac{\gamma C - 1 - \gamma}{C} \psi(t) \mathbf{1}_C \right), \quad (\text{A.30})$$

where $\mathbf{Z}(t)$ is defined in Eq. (3.5), $\phi_i(t) = e^{-\int_0^t \lambda_i(r) dr}$ for $i \in \{1, 2, 3\}$, and $\psi(t) = \int_0^t \eta(s) e^{\int_0^s \lambda_3(r) dr} ds$.

Proof. According to the gradient flow in Eq. (3.11), and the notations in the proof of Theorem 3.1, we have:

$$\mathbf{Z}'(t) = \mathbf{Z}(t) \mathbf{A}(t) - \lambda \mathbf{Z}(t) \begin{pmatrix} \eta_1(t) \mathbf{I}_{CN} & 0 \\ 0 & \eta_2(t) \mathbf{I}_C \end{pmatrix}, \quad (\text{A.31})$$

i.e., $\mathbf{Z}'(t) = \mathbf{Z}(t) \mathbf{A}(t)$, where $\mathbf{A}(t) = \mathbf{A}(t) - \lambda \mathbf{\Lambda}(t)$, and $\mathbf{\Lambda}(t) = \begin{pmatrix} \eta_1(t) \mathbf{I}_{CN} & 0 \\ 0 & \eta_2(t) \mathbf{I}_C \end{pmatrix}$.

For any t_1, t_2 , we have the matrix commutator of $\mathbf{A}(t_1)$ and $\mathbf{A}(t_2)$

$$\begin{aligned} & [\mathbf{A}(t_1), \mathbf{A}(t_2)] \\ &= \mathbf{A}(t_1) \mathbf{A}(t_2) - \mathbf{A}(t_2) \mathbf{A}(t_1) \\ &= [\mathbf{A}(t_1) - \lambda \mathbf{\Lambda}(t_1)] [\mathbf{A}(t_2) - \lambda \mathbf{\Lambda}(t_2)] - [\mathbf{A}(t_2) - \lambda \mathbf{\Lambda}(t_2)] [\mathbf{A}(t_1) - \lambda \mathbf{\Lambda}(t_1)] \\ &= -\lambda [\mathbf{\Lambda}(t_1) \mathbf{A}(t_2) + \mathbf{A}(t_1) \mathbf{\Lambda}(t_2) - \mathbf{\Lambda}(t_2) \mathbf{A}(t_1) - \mathbf{A}(t_2) \mathbf{\Lambda}(t_1)] \\ &= -\lambda \begin{pmatrix} 0 & [\eta_1(t_1)\eta_2(t_2) - \eta_1(t_2)\eta_2(t_1)] \mathbf{M}^\top \\ [\eta_2(t_1)\eta_1(t_2) - \eta_2(t_2)\eta_1(t_2)] \mathbf{M} & 0 \end{pmatrix} \\ &= 0 \end{aligned} \quad (\text{A.32})$$

where the last equality is based on the fact that $\eta_2(t_1)\eta_1(t_2) = \eta_2(t_2)\eta_1(t_1)$. Therefore, according to Magnus approach, we have

$$\mathbf{Z}(t) = \mathbf{Z}_0 \exp \left(\int_0^t \mathbf{A}(\tau) d\tau \right) = \mathbf{Z}_0 \exp \begin{pmatrix} -\lambda \zeta_1(t) \mathbf{I}_{CN} & \zeta_2(t) \mathbf{M}^\top \\ \zeta_1(t) \mathbf{M} & -\lambda \zeta_2(t) \mathbf{I}_C \end{pmatrix}, \quad (\text{A.33})$$

where $\zeta_1(t) = \int_0^t \eta(\tau) d\tau$ and $\zeta_2(t) = \int_0^t \eta(\tau) d\tau$.

Let $\mathbf{B} = \mathbf{B} - \lambda \mathbf{I}_{C(N+1)}$, we have

$$\mathbf{Z}(t) = \mathbf{Z}_0 \exp(\mathbf{B}\mathbf{C}(t)) = \mathbf{Z}_0 \sum_{k=0}^{\infty} \frac{(\mathbf{B}\mathbf{C}(t))^k}{k!}. \quad (\text{A.34})$$

We again consider the orthogonal decomposition of \mathbf{Z}_0 , i.e., $\mathbf{Z}_0 = (\Pi_1^+ + \Pi_1^- + \Pi_2^+ + \Pi_2^- + \Pi_3) \mathbf{Z}_0$. As mentioned in the proof of Theorem 3.1, we have

$$\begin{aligned} \Pi_1^\epsilon \mathbf{Z}_0 \mathbf{B} &= \frac{\epsilon(1+\gamma)}{C\sqrt{N}} \Pi_1^\epsilon \mathbf{Z}_0, & \Pi_1^\epsilon \mathbf{Z}_0 \mathbf{B} &= \frac{\epsilon(1+\gamma) - \lambda C \sqrt{N}}{C\sqrt{N}} \Pi_1^\epsilon \mathbf{Z}_0, \\ \Pi_2^\epsilon \mathbf{Z}_0 \mathbf{B} &= \frac{\epsilon(1+\gamma - \gamma C)}{C\sqrt{N}} \Pi_2^\epsilon \mathbf{Z}_0, & \Rightarrow \Pi_2^\epsilon \mathbf{Z}_0 \mathbf{B} &= \frac{\epsilon(1+\gamma - \gamma C) - \lambda C \sqrt{N}}{C\sqrt{N}} \Pi_2^\epsilon \mathbf{Z}_0, \\ \Pi_3 \mathbf{Z}_0 \mathbf{B} &= 0, & \Pi_3 \mathbf{Z}_0 \mathbf{B} &= -\lambda \Pi_3 \mathbf{Z}_0. \end{aligned} \quad (\text{A.35})$$

Therefore, for any $D = (\mathbf{H}, \mathbf{W}) \in \{\Pi_1^\epsilon \mathbf{Z}_0, \Pi_2^\epsilon \mathbf{Z}_0, \Pi_3 \mathbf{Z}_0\}$ (where $\mathbf{H} \in \mathbb{R}^{p \times CN}$ and $\mathbf{W} \in \mathbb{R}^{p \times C}$) and the corresponding eigenvalue $\sigma \in \left\{ \frac{\epsilon(1+\gamma)}{C\sqrt{N}}, \frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}}, 0 \right\}$, we have

$$D\mathbf{B} = \sigma D, \mathbf{H}\mathbf{M}^\top = \sigma \mathbf{W}, \text{ and } \mathbf{W}\mathbf{M} = \sigma \mathbf{H}. \quad (\text{A.36})$$

In the following, we will prove that there exist two scalars $a(t)$ and $b(t)$, such that $D \exp(\mathbf{B}\mathbf{C}(t)) = D \begin{pmatrix} a(t)\mathbf{I}_{CN} & 0 \\ 0 & b(t)\mathbf{I}_C \end{pmatrix}$.

First, we prove that $D(\mathbf{B}\mathbf{C}(t))^k$ can be represented as $D(\mathbf{B}\mathbf{C}(t))^k = (a_k(t)\mathbf{H}, b_k(t)\mathbf{W})$ by induction, where $a_k(t), b_k(t) \in \mathbb{R}$.

For $k = 0$, we have $D(\mathbf{B}\mathbf{C}(t))^0 = (\mathbf{H}, \mathbf{W})$, i.e., $a_0 = b_0 = 1$. Assume that $D(\mathbf{B}\mathbf{C}(t))^n = (a_n(t)\mathbf{H}, b_n(t)\mathbf{W})$ for $k = n$. Then for $k = n + 1$, we have

$$\begin{aligned} & D(\mathbf{B}\mathbf{C}(t))^{n+1} \\ &= D(\mathbf{B}\mathbf{C}(t))^n(\mathbf{B}\mathbf{C}(t)) \\ &= (a_n(t)\mathbf{H}, b_n(t)\mathbf{W})(\mathbf{B}\mathbf{C}(t)) \\ &= (a_n(t)\mathbf{H}, b_n(t)\mathbf{W}) \begin{pmatrix} -\lambda\mathbf{I}_{CN} & \mathbf{M}^\top \\ \mathbf{M} & -\lambda\mathbf{I}_C \end{pmatrix} \begin{pmatrix} \zeta_1(t)\mathbf{I}_{CN} & 0 \\ 0 & \zeta_2(t)\mathbf{I}_C \end{pmatrix}, \quad (\text{A.37}) \\ &= (b_n(t)\mathbf{W}\mathbf{M} - \lambda a_n(t)\mathbf{H}, a_n(t)\mathbf{H}\mathbf{M}^\top - \lambda b_n(t)\mathbf{W}) \begin{pmatrix} \zeta_1(t)\mathbf{I}_{CN} & 0 \\ 0 & \zeta_2(t)\mathbf{I}_C \end{pmatrix} \\ &= (\zeta_1(t)(\sigma b_n(t) - \lambda a_n(t))\mathbf{H}, \zeta_2(t)(\sigma a_n(t) - \lambda b_n(t))\mathbf{H}) \end{aligned}$$

thus $a_{n+1}(t) = \zeta_1(t)(\sigma b_n(t) - \lambda a_n(t))$ and $b_{n+1}(t) = \zeta_2(t)(\sigma a_n(t) - \lambda b_n(t))$.

To sum up, we have shown by induction that $D(\mathbf{B}\mathbf{C}(t))^k$ can be represented as $D(\mathbf{B}\mathbf{C}(t))^k = (a_k(t)\mathbf{H}, b_k(t)\mathbf{W}) = D \begin{pmatrix} a_k(t)\mathbf{I}_{CN} & 0 \\ 0 & b_k(t)\mathbf{I}_C \end{pmatrix}$, and $\begin{pmatrix} a_k(t) \\ b_k(t) \end{pmatrix}$ satisfies

$$\begin{pmatrix} a_k(t) \\ b_k(t) \end{pmatrix} = \begin{pmatrix} -\lambda\zeta_1(t) & \sigma\zeta_1(t) \\ \sigma\zeta_2(t) & -\lambda\zeta_2(t) \end{pmatrix} \begin{pmatrix} a_{k-1}(t) \\ b_{k-1}(t) \end{pmatrix} \quad \text{with} \quad \begin{pmatrix} a_0 \\ b_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad (\text{A.38})$$

i.e., $\begin{pmatrix} a_k(t) \\ b_k(t) \end{pmatrix} = (\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t)))^k \begin{pmatrix} 1 \\ 1 \end{pmatrix}$, where $\mathbf{S}(\sigma, \lambda, \zeta_1, \zeta_2) = \begin{pmatrix} -\lambda\zeta_1 & \sigma\zeta_1 \\ \sigma\zeta_2 & -\lambda\zeta_2 \end{pmatrix}$.

Therefore, we have

$$D \exp(\mathbf{B}\mathbf{C}(t)) = D \sum_{k=0}^{\infty} \frac{(\mathbf{B}\mathbf{C}(t))^k}{k!} = D \begin{pmatrix} a(t)\mathbf{I}_{CN} & 0 \\ 0 & b(t)\mathbf{I}_C \end{pmatrix}, \quad (\text{A.39})$$

with $a(t) = \sum_{k=0}^{\infty} \frac{a_k(t)}{k!}$ and $b(t) = \sum_{k=0}^{\infty} \frac{b_k(t)}{k!}$, i.e.,

$$\begin{pmatrix} a(t) \\ b(t) \end{pmatrix} = \sum_{k=0}^{\infty} \frac{1}{k!} \begin{pmatrix} a_k(t) \\ b_k(t) \end{pmatrix} = \sum_{k=0}^{\infty} \frac{(\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t)))^k}{k!} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \exp(\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t))) \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

Next, we are going to derive a concrete expression of $\exp(\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t)))$. Let the determinant $|\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t)) - \theta\mathbf{I}| = 0$, we can derive that the eigenvalues of $\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t))$ are $\theta_1 = \frac{-\lambda(\zeta_1(t) + \zeta_2(t)) - \sqrt{\lambda^2(\zeta_1(t) - \zeta_2(t))^2 + 4\sigma^2\zeta_1(t)\zeta_2(t)}}{2} < 0$ and $\theta_2 = \frac{-\lambda(\zeta_1(t) + \zeta_2(t)) + \sqrt{\lambda^2(\zeta_1(t) - \zeta_2(t))^2 + 4\sigma^2\zeta_1(t)\zeta_2(t)}}{2} = \frac{(\sqrt{\lambda^2(s-1)^2 + 4s\sigma^2} - \lambda(s+1))\zeta_2(t)}{2}$, and the corresponding eigenvectors are

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ \frac{\lambda\zeta_1(t) + \theta_1}{\sigma\zeta_1(t)} \end{pmatrix}, \text{ and } \mathbf{v}_2 = \begin{pmatrix} 1 \\ \frac{\lambda\zeta_1(t) + \theta_2}{\sigma\zeta_1(t)} \end{pmatrix}.$$

Let $\mathbf{P} = (\mathbf{v}_1, \mathbf{v}_2)$, we have

$$\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t)) = \mathbf{P} \begin{pmatrix} \theta_1 & 0 \\ 0 & \theta_2 \end{pmatrix} \mathbf{P}^{-1}, \quad \mathbf{P}^{-1} = \frac{\sigma\zeta_1(t)}{\theta_2 - \theta_1} \begin{pmatrix} \frac{\lambda\eta_1(t) + \theta_2}{\sigma\zeta_1(t)} & -1 \\ -\frac{\lambda\zeta_1(t) + \theta_1}{\sigma\zeta_1(t)} & 1 \end{pmatrix}, \quad (\text{A.40})$$

and

$$\begin{aligned}
\exp(\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t))) &= \mathbf{P} \begin{pmatrix} e^{\theta_1} & 0 \\ 0 & e^{\theta_2} \end{pmatrix} \mathbf{P}^{-1} \\
&= \frac{\sigma \zeta_1(t)}{\theta_2 - \theta_1} \begin{pmatrix} 1 & 1 \\ \frac{\lambda \zeta_1(t) + \theta_1}{\sigma \zeta_1(t)} & \frac{\lambda \zeta_1(t) + \theta_2}{\sigma \zeta_1(t)} \end{pmatrix} \begin{pmatrix} e^{\theta_1} & 0 \\ 0 & e^{\theta_2} \end{pmatrix} \begin{pmatrix} \frac{\lambda \eta_1(t) + \theta_2}{\sigma \zeta_1(t)} & -1 \\ -\frac{\lambda \zeta_1(t) + \theta_1}{\sigma \zeta_1(t)} & 1 \end{pmatrix}, \quad (\text{A.41}) \\
&= \frac{1}{\theta_2 - \theta_1} \begin{pmatrix} (\lambda \zeta_1(t) + \theta_2)e^{\theta_1} - (\lambda \zeta_1(t) + \theta_1)e^{\theta_2} & \sigma \zeta_1(t)e^{\theta_2} - \sigma \zeta_1(t)e^{\theta_1} \\ (\lambda \zeta_1(t) + \theta_2)e^{\theta_2} - (\lambda \zeta_1(t) + \theta_1)e^{\theta_1} & \sigma \zeta_2(t)e^{\theta_2} - \sigma \theta_2(t)e^{\theta_1} \end{pmatrix}
\end{aligned}$$

thus

$$\begin{aligned}
a(t) &= \frac{e^{\theta_2}}{\theta_2 - \theta_1} [(\sigma \zeta_1(t) - \lambda \zeta_1(t) - \theta_1) - (\sigma \zeta_1(t) - \lambda \zeta_1(t) - \theta_2)e^{\theta_1 - \theta_2}], \\
b(t) &= \frac{e^{\theta_2}}{\theta_2 - \theta_1} [(\lambda \zeta_1(t) + \theta_2 + \sigma \zeta_2(t)) + (\lambda \zeta_1(t) + \theta_1 - \sigma \zeta_2(t))e^{\theta_1 - \theta_2}]. \quad (\text{A.42})
\end{aligned}$$

Finally, since \mathcal{E}_1^+ , \mathcal{E}_1^- , \mathcal{E}_2^+ , \mathcal{E}_2^- and \mathcal{E}_3 are orthogonal subspace of $\mathbb{R}^{p \times CN} \oplus \mathbb{R}^{C \times p}$, we obtain that

$$\begin{aligned}
\mathbf{Z}(t) &= (\Pi_1^+ \mathbf{Z}_0 + \Pi_1^- \mathbf{Z}_0 + \Pi_2^+ \mathbf{Z}_0 + \Pi_2^- \mathbf{Z}_0 + \Pi_3 \mathbf{Z}_0) \exp(\mathbf{L}(t)) \\
&= \Pi_1^+ \mathbf{Z}_0 \begin{pmatrix} a_1^+(t) \mathbf{I}_{CN} & 0 \\ 0 & b_1^+(t) \mathbf{I}_C \end{pmatrix} + \Pi_1^- \mathbf{Z}_0 \begin{pmatrix} a_1^-(t) \mathbf{I}_{CN} & 0 \\ 0 & b_1^-(t) \mathbf{I}_C \end{pmatrix} + \\
&\quad \Pi_2^+ \mathbf{Z}_0 \begin{pmatrix} a_2^+(t) \mathbf{I}_{CN} & 0 \\ 0 & b_2^+(t) \mathbf{I}_C \end{pmatrix} + \Pi_2^- \mathbf{Z}_0 \begin{pmatrix} a_2^-(t) \mathbf{I}_{CN} & 0 \\ 0 & b_2^-(t) \mathbf{I}_C \end{pmatrix} + \\
&\quad \Pi_3 \mathbf{Z}_0 \begin{pmatrix} a_3(t) \mathbf{I}_{CN} & 0 \\ 0 & b_3(t) \mathbf{I}_C \end{pmatrix}, \quad (\text{A.43})
\end{aligned}$$

where $a_1^\epsilon(t)$, $b_1^\epsilon(t)$, $a_2^\epsilon(t)$, $b_2^\epsilon(t)$, $a_3(t)$ and $b_3(t)$ satisfy

$$\begin{aligned}
\begin{pmatrix} a_1^\epsilon(t) \\ b_1^\epsilon(t) \end{pmatrix} &= \exp\left(\mathbf{S}\left(\frac{\epsilon(1+\gamma)}{C\sqrt{N}}, \lambda, \zeta_1(t), \zeta_2(t)\right)\right) \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \\
\begin{pmatrix} a_2^\epsilon(t) \\ b_2^\epsilon(t) \end{pmatrix} &= \exp\left(\mathbf{S}\left(\frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}}, \lambda, \zeta_1(t), \zeta_2(t)\right)\right) \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad (\text{A.44}) \\
\begin{pmatrix} a_3(t) \\ b_3(t) \end{pmatrix} &= \exp(\mathbf{S}(0, \lambda, \zeta_1(t), \zeta_2(t))) \begin{pmatrix} 1 \\ 1 \end{pmatrix}.
\end{aligned}$$

Moreover, since $\mathbf{b}'(t) = \eta_2(t) \frac{1+\gamma-\gamma C}{C} \mathbf{1}_C - \lambda \eta_2(t) \mathbf{b}(t)$ is a first-order linear differential equation, then we have

$$\mathbf{b}(t) = \phi(t) \left(\mathbf{b}_0 + \frac{1+\gamma-\gamma C}{C} \psi(t) \mathbf{1}_C \right), \quad (\text{A.45})$$

where $\phi(t) = \exp(-\lambda \int_0^t \eta_2(\tau) d\tau)$ and $\psi(t) = \int_0^t \zeta_2(\tau) \exp(\lambda \int_0^\tau \eta_2(s) ds) d\tau$. \square

A.5 PROOF OF COROLLARY 3.4

Corollary 3.4. *Under the conditions and notation of Theorem 3.3, let $s = \frac{\eta_1(0)}{\eta_2(0)}$, if $0 < \gamma < \frac{2}{C-2}$ (where $C > 2$) or $C = 2$, and $\lim_{t \rightarrow \infty} \zeta_1(t) = \infty$, then there exist constants π_h^+ , π_h^- , π_w^+ , π_w^- and ω only depending on λ , γ , s , C and N , such that the gradient flow behaves as:*

$$\left\| \frac{\mathbf{H}(t)}{\|\mathbf{H}(t)\|} - \frac{\pi_h^+ \mathbf{H}_1^+ + \pi_h^- \mathbf{H}_1^-}{\|\pi_h^+ \mathbf{H}_1^+ + \pi_h^- \mathbf{H}_1^-\|} \right\| + \left\| \frac{\mathbf{W}(t)}{\|\mathbf{W}(t)\|} - \frac{\pi_w^+ \mathbf{W}_1^+ + \pi_w^- \mathbf{W}_1^-}{\|\pi_w^+ \mathbf{H}_1^+ + \pi_w^- \mathbf{H}_1^-\|} \right\| = O(e^{-\omega \zeta_2(t)}), \quad (\text{A.46})$$

where $(\mathbf{H}_1^+, \mathbf{W}_1^+) = \Pi_1^+ \mathbf{Z}_0$, $(\mathbf{H}_1^-, \mathbf{W}_1^-) = \Pi_1^- \mathbf{Z}_0$. Furthermore, we have the following results:

- If $\lambda > \frac{1+\gamma}{C\sqrt{N}}$, then $\lim_{t \rightarrow \infty} \|\mathbf{Z}(t)\| = 0$;
- If $\lambda = \frac{1+\gamma}{C\sqrt{N}}$, then $\lim_{t \rightarrow \infty} \mathbf{H}(t) = \mathbf{H}_1^+ + \frac{1-s}{1+s} \mathbf{H}_1^-$, $\lim_{t \rightarrow \infty} \mathbf{W}(t) = \mathbf{W}_1^+ - \frac{1-s}{1+s} \mathbf{W}_1^-$;
- If $\lambda < \frac{1+\gamma}{C\sqrt{N}}$, then $\lim_{t \rightarrow \infty} \|\mathbf{Z}(t)\| = \infty$.

Proof. Since $\zeta_1(t) = s\zeta_2(t)$, then the eigenvalues of $\mathbf{S}(\sigma, \lambda, \zeta_1(t), \zeta_2(t))$ are $\theta_1 = -\frac{\lambda(s+1) + \sqrt{\lambda^2(s-1)^2 + 4s\sigma^2}}{2}\zeta_2(t) \rightarrow -\infty$ as $t \rightarrow \infty$ and $\theta_2 = \frac{(\sqrt{\lambda^2(s-1)^2 + 4s\sigma^2} - \lambda(s+1))}{2}\zeta_2(t)$.

Let $\omega_1(\sigma, \lambda, s) = -\frac{\lambda(s+1) + \sqrt{\lambda^2(s-1)^2 + 4s\sigma^2}}{2} < 0$ and $\omega_2(\sigma, \lambda, s) = \frac{(\sqrt{\lambda^2(s-1)^2 + 4s\sigma^2} - \lambda(s+1))}{2}$. For brevity, let ω_1 and ω_2 denote $\omega_1(\sigma, \lambda, s)$ and $\omega_2(\sigma, \lambda, s)$, respectively, and then we can reformulate $a(t)$ and $b(t)$ as

$$\begin{aligned} a(t) &= \frac{e^{\omega_2\zeta_2(t)}}{\omega_2 - \omega_1} \left[(s\sigma - s\lambda - \omega_1) - (s\sigma - s\lambda - \omega_2)e^{(\omega_1 - \omega_2)\zeta_2(t)} \right] \\ &= \frac{s\sigma - s\lambda - \omega_1}{\omega_2 - \omega_1} e^{\omega_2\zeta_2(t)} + O(e^{(\omega_1)\zeta_2(t)}), \\ b(t) &= \frac{e^{\omega_2\zeta_2(t)}}{\omega_2 - \omega_1} \left[(s\lambda + \omega_2 + \sigma) + (s\lambda + \omega_1 - \sigma)e^{(\omega_1 - \omega_2)\zeta_2(t)} \right] \\ &= \frac{s\lambda + \omega_2 + \sigma}{\omega_2 - \omega_1} e^{\omega_2\zeta_2(t)} + O(e^{\omega_1\zeta_2(t)}). \end{aligned} \quad (\text{A.47})$$

Moreover, according to Theorem 3.3, we have

$$\mathbf{H}(t) = \sum_{\substack{\epsilon \in \{\pm\} \\ i \in \{1,2\}}} a_i^\epsilon(t) \mathbf{H}_i^\epsilon + a_3(t) \mathbf{H}_3, \quad \mathbf{W}(t) = \sum_{\substack{\epsilon \in \{\pm\} \\ i \in \{1,2\}}} b_i^\epsilon(t) \mathbf{W}_i^\epsilon + b_3(t) \mathbf{W}_3, \quad (\text{A.48})$$

with

$$\begin{aligned} \begin{pmatrix} a_1^\epsilon(t) \\ b_1^\epsilon(t) \end{pmatrix} &= \exp\left(\mathbf{S}\left(\frac{\epsilon(1+\gamma)}{C\sqrt{N}}, \lambda, \zeta_1(t), \zeta_2(t)\right)\right) \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \\ \begin{pmatrix} a_2^\epsilon(t) \\ b_2^\epsilon(t) \end{pmatrix} &= \exp\left(\mathbf{S}\left(\frac{\epsilon(1+\gamma-\gamma C)}{C\sqrt{N}}, \lambda, \zeta_1(t), \zeta_2(t)\right)\right) \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \\ \begin{pmatrix} a_3(t) \\ b_3(t) \end{pmatrix} &= \exp\left(\mathbf{S}(0, \lambda, \zeta_1(t), \zeta_2(t))\right) \begin{pmatrix} 1 \\ 1 \end{pmatrix}. \end{aligned} \quad (\text{A.49})$$

Since $0 < \gamma < \frac{2}{C-2}$ (where $C > 2$) or $C = 2$, thus we have $\frac{1+\gamma}{C\sqrt{N}} > \frac{|1+\gamma-\gamma C|}{C\sqrt{N}}$, and then $\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right) > \omega_2\left(\frac{|1+\gamma-\gamma C|}{C\sqrt{N}}, \lambda, s\right)$. When $t \rightarrow \infty$, the dominant terms in $\mathbf{H}(t)$ and $\mathbf{W}(t)$ are the ones whose coefficient contains $\exp\left(\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)\zeta_2(t)\right)$, i.e., $a_1^+(t)$, $a_1^-(t)$, $b_1^+(t)$ and $b_1^-(t)$. Let $(\mathbf{H}_1^\epsilon, \mathbf{W}_1^\epsilon) = \Pi_1^\epsilon \mathbf{Z}_0$, $(\mathbf{H}_2^\epsilon, \mathbf{W}_2^\epsilon) = \Pi_2^\epsilon \mathbf{Z}_0$ and $(\mathbf{H}_3, \mathbf{W}_3) = \Pi_3 \mathbf{Z}_0$ for $\epsilon \in \{\pm 1\}$, thus we have

$$\begin{aligned} e^{-\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)\zeta_2(t)} \mathbf{H}(t) &= \pi_h^+(\lambda, \gamma, s, C, N) \mathbf{H}_1^+ + \pi_h^-(\lambda, \gamma, s, C, N) \mathbf{H}_1^- + \mathbf{\Delta}_1, \\ e^{-\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)\zeta_2(t)} \mathbf{W}(t) &= \pi_w^+(\lambda, \gamma, s, C, N) \mathbf{W}_1^+ + \pi_w^-(\lambda, \gamma, s, C, N) \mathbf{W}_1^- + \mathbf{\Delta}_2, \end{aligned} \quad (\text{A.50})$$

where $\mathbf{\Delta}_1$ and $\mathbf{\Delta}_2$ decrease to zero as least as $O\left(e^{\frac{\left(\omega_2\left(\frac{|1+\gamma-\gamma C|}{C\sqrt{N}}, \lambda, s\right) - \omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)\right)\zeta_2(t)}{C\sqrt{N}}}\right)$, and

$$\begin{aligned} \pi_h^+(\lambda, \gamma, s, C, N) &= \frac{s\frac{1+\gamma}{C\sqrt{N}} - s\lambda - \omega_1\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)}{\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right) - \omega_1\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)}, \\ \pi_h^-(\lambda, \gamma, s, C, N) &= \frac{-s\frac{1+\gamma}{C\sqrt{N}} - s\lambda - \omega_1\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)}{\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right) - \omega_1\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)}, \\ \pi_w^+(\lambda, \gamma, s, C, N) &= \frac{s\lambda + \omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right) + \frac{1+\gamma}{C\sqrt{N}}}{\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right) - \omega_1\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)}, \\ \pi_w^-(\lambda, \gamma, s, C, N) &= \frac{s\lambda + \omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right) - \frac{1+\gamma}{C\sqrt{N}}}{\omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right) - \omega_1\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)}. \end{aligned} \quad (\text{A.51})$$

Therefore, we have

$$\begin{aligned}\lim_{t \rightarrow \infty} \frac{\mathbf{H}(t)}{\|\mathbf{H}(t)\|} &= \frac{\pi_h^+(\lambda, \gamma, s, C, N) \mathbf{H}_1^+ + \pi_h^-(\lambda, \gamma, s, C, N) \mathbf{H}_1^-}{\|\pi_h^+(\lambda, \gamma, s, C, N) \mathbf{H}_1^+ + \pi_h^-(\lambda, \gamma, s, C, N) \mathbf{H}_1^-\|}, \\ \lim_{t \rightarrow \infty} \frac{\mathbf{W}(t)}{\|\mathbf{W}(t)\|} &= \frac{\pi_w^+(\lambda, \gamma, s, C, N) \mathbf{W}_1^+ + \pi_w^-(\lambda, \gamma, s, C, N) \mathbf{W}_1^-}{\|\pi_w^+(\lambda, \gamma, s, C, N) \mathbf{W}_1^+ + \pi_w^-(\lambda, \gamma, s, C, N) \mathbf{W}_1^-\|},\end{aligned}\tag{A.52}$$

and the rate of convergence is $O\left(e^{\left(\frac{\omega_2\left(\frac{1+\gamma-\gamma C}{C\sqrt{N}}, \lambda, s\right) - \omega_2\left(\frac{1+\gamma}{C\sqrt{N}}, \lambda, s\right)\right)\zeta_2(t)}\right)$.

Moreover, we have the following conclusions:

- If $\lambda = \frac{1+\gamma}{C\sqrt{N}}$, we have $\omega_2 = 0$, $\omega_1 = -\lambda(s+1)$, $\pi_h^+ = 1$, $\pi_h^- = \frac{1-s}{1+s}$, $\pi_w^+ = 1$ and $\pi_w^- = -\frac{1-s}{1+s}$. Since $\lim_{t \rightarrow \infty} \zeta_2(t) = \infty$, we have

$$\lim_{t \rightarrow \infty} \mathbf{H}(t) = \mathbf{H}_1^+ + \frac{1-s}{1+s} \mathbf{H}_1^-, \quad \lim_{t \rightarrow \infty} \mathbf{W}(t) = \mathbf{W}_1^+ - \frac{1-s}{1+s} \mathbf{W}_1^-. \tag{A.53}$$

- If $\lambda > \frac{1+\gamma}{C\sqrt{N}}$, we have $\omega_2 > 0$, and then $\lim_{t \rightarrow \infty} \|\mathbf{Z}(t)\| = 0$ since $\lim_{t \rightarrow \infty} \zeta_2(t) = \infty$.
- If $\lambda < \frac{1+\gamma}{C\sqrt{N}}$, we have $\omega_2 < 0$, and then $\lim_{t \rightarrow \infty} \|\mathbf{Z}(t)\| = \infty$ since $\lim_{t \rightarrow \infty} \zeta_2(t) = \infty$.

So far the proof has been completed. \square

A.6 PROOF OF THEOREM 3.5

Lemma A.2. For $\mathbf{h}(t), \mathbf{w} \in \mathbb{R}^p$, $\eta(t) > 0$, let $\hat{\mathbf{v}} = \frac{\hat{\mathbf{v}}}{\|\hat{\mathbf{v}}\|_2}$ denote the ℓ_2 -normalized vector of \mathbf{v} , considering the discrete dynamical system $\mathbf{h}(t+1) = \mathbf{h}(t) + \frac{\eta(t)}{\|\mathbf{h}(t)\|_2} \left(\mathbf{I}_p - \hat{\mathbf{h}}(t)\hat{\mathbf{h}}^\top(t)\right) \mathbf{w}$, if $\hat{\mathbf{w}}^\top \hat{\mathbf{h}}(0) > -1$, the learning rate $\eta(t)$ satisfies that $\lim_{t \rightarrow \infty} \frac{\eta(t+1)}{\eta(t)} = 1$, $\frac{\eta(t)}{\|\mathbf{h}(t)\|_2}$ is non-increasing with $\frac{\eta(0)}{\|\mathbf{h}(0)\|_2} < \frac{1}{\|\mathbf{w}\|_2}$, and there exists a constant $\varepsilon > 0$, s.t., $\eta(t) > \varepsilon$, s.t., $\eta(t) > \varepsilon$, then we have

$$\lim_{t \rightarrow \infty} \left\| \frac{\mathbf{h}(t)}{\|\mathbf{h}(t)\|_2} - \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right\| = 0. \tag{A.54}$$

Proof. For brevity, let $\alpha_t = \frac{\eta(t)\|\mathbf{w}\|_2}{\|\mathbf{h}(t)\|_2^2}$, $\xi_t = \frac{\eta(t+1)}{\eta(t)}$, and $\beta_t = \hat{\mathbf{w}}^\top \hat{\mathbf{h}}(t)$ denote the cosine similarity between $\hat{\mathbf{w}}$ and $\hat{\mathbf{h}}(t)$, then we can easily derive that $\alpha_t > 0$ and $\beta_t > -1$ for all $t \geq 0$.

We will show that α_t is monotonically decreasing and β_t is monotonically increasing. Note that $\mathbf{h}(t)$ is orthogonal with $\left(\mathbf{I}_p - \hat{\mathbf{h}}(t)\hat{\mathbf{h}}^\top(t)\right) \mathbf{w}$, thus we have

$$\|\mathbf{h}(t+1)\|_2^2 = \|\mathbf{h}(t)\|_2^2 + \frac{\eta^2(t)}{\|\mathbf{h}(t)\|_2^2} \left\| \left(\mathbf{I}_p - \hat{\mathbf{h}}(t)\hat{\mathbf{h}}^\top(t)\right) \mathbf{w} \right\|_2^2 \geq \|\mathbf{h}(t)\|_2^2,$$

which indicates $\|\mathbf{h}(t)\|_2$ is monotonically increasing as a function of t , and then α_t is monotonically decreasing since $\frac{\eta(t)\|\mathbf{w}\|_2}{\|\mathbf{h}(t)\|_2^2}$ is non-increasing.

Moreover, we can rearrange the discrete dynamics and formulate $\mathbf{h}(t+1)$ as a positive combination of $\mathbf{h}(t)$ and \mathbf{w} :

$$\mathbf{h}(t+1) = \left(1 - \frac{\eta(t)\mathbf{w}^\top \hat{\mathbf{h}}(t)}{\|\mathbf{h}(t)\|_2}\right) \mathbf{h}(t) + \frac{\eta(t)}{\|\mathbf{h}(t)\|_2} \mathbf{w}, \tag{A.55}$$

so $\hat{\mathbf{w}}^\top \hat{\mathbf{h}}(t+1) \geq \hat{\mathbf{w}}^\top \hat{\mathbf{h}}(t)$, i.e., β_t is monotonically increasing, which is based on the facts that $1 - \frac{\eta(t)\mathbf{w}^\top \hat{\mathbf{h}}(t)}{\|\mathbf{h}(t)\|_2} \geq 1 - \frac{\eta(t)\|\mathbf{w}\|_2}{\|\mathbf{h}(t)\|_2} \geq 1 - \frac{\eta(0)\|\mathbf{w}\|_2}{\|\mathbf{h}(0)\|_2} > 0$, $\frac{\eta(t)}{\|\mathbf{h}(t)\|_2} > 0$, and $\frac{\mathbf{y}^\top (\mathbf{x} + k\mathbf{y})}{\|\mathbf{x} + k\mathbf{y}\|_2} \geq \frac{\mathbf{y}^\top \mathbf{x}}{\|\mathbf{x}\|_2}$ holds for all $k > 0$ and $\mathbf{x}, \mathbf{y} \neq 0$.

We can formulate the discrete iterations of α_t and β_t from the Eq. (A.55) as follows:

$$\begin{aligned}
\beta_{t+1} &= \frac{\mathbf{w}^\top \mathbf{h}(t+1)}{\|\mathbf{w}\|_2 \|\mathbf{h}(t+1)\|_2} \\
&= \frac{\mathbf{w}^\top \left(\mathbf{h}(t) + \frac{\eta(t)}{\|\mathbf{h}(t)\|_2} \left(\mathbf{I}_p - \widehat{\mathbf{h}}(t) \widehat{\mathbf{h}}^\top(t) \right) \mathbf{w} \right)}{\|\mathbf{w}\|_2 \sqrt{\|\mathbf{h}(t)\|_2^2 + \frac{\eta^2(t)}{\|\mathbf{h}(t)\|_2^2} \left\| \left(\mathbf{I}_p - \widehat{\mathbf{h}}(t) \widehat{\mathbf{h}}^\top(t) \right) \mathbf{w} \right\|_2^2}} \\
&= \frac{\beta_t + \frac{\eta(t) \|\mathbf{w}\|_2}{\|\mathbf{h}(t)\|_2^2} (1 - \beta_t^2)}{\sqrt{1 + \frac{\eta^2(t) \|\mathbf{w}\|_2^2}{\|\mathbf{h}(t)\|_2^2} (1 - \beta_t^2)}} \\
&= \frac{\beta_t + \alpha_t (1 - \beta_t^2)}{\sqrt{1 + \alpha_t^2 (1 - \beta_t^2)}}, \\
\alpha_{t+1} &= \frac{\eta(t+1) \|\mathbf{w}\|_2}{\|\mathbf{h}(t+1)\|_2^2} \\
&= \frac{\xi_t \eta(t) \|\mathbf{w}\|_2}{\|\mathbf{h}(t)\|_2^2 + \frac{\eta^2(t)}{\|\mathbf{h}(t)\|_2^2} \left\| \left(\mathbf{I}_p - \widehat{\mathbf{h}}(t) \widehat{\mathbf{h}}^\top(t) \right) \mathbf{w} \right\|_2^2} \\
&= \frac{\xi_t \alpha_t}{1 + \alpha_t^2 (1 - \beta_t^2)},
\end{aligned} \tag{A.56}$$

with $\beta_0 = \widehat{\mathbf{w}}^\top \widehat{\mathbf{h}}(0) > -1$, $\alpha_0 = \frac{\eta(0) \|\mathbf{w}\|_2}{\|\mathbf{h}(0)\|_2^2} > 0$, $\xi_t \leq 1$ and $\lim_{t \rightarrow \infty} \xi_t = 1$.

To prove $\lim_{t \rightarrow \infty} \left\| \frac{\mathbf{h}(t)}{\|\mathbf{h}(t)\|_2} - \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right\| = 0$, we just need prove $\lim_{t \rightarrow \infty} \beta_t = 1$. Note that α_t is monotonic decreasing and lower bounded by 0, then the sequence (α_t) is convergent. Similarly, the sequence (β_t) is convergent. Let $a = \lim_{t \rightarrow \infty} \alpha_t$ and $b = \lim_{t \rightarrow \infty} \beta_t$, we obtain

$$\lim_{t \rightarrow \infty} \alpha_{t+1} = \lim_{t \rightarrow \infty} \frac{\xi_t \alpha_t}{1 + \alpha_t^2 (1 - \beta_t^2)} \Rightarrow a = \frac{a}{1 + \lim_{t \rightarrow \infty} \alpha_t^2 (1 - \beta_t^2)}, \tag{A.57}$$

thus $a = 0$ or $\lim_{t \rightarrow \infty} \alpha_t^2 (1 - \beta_t^2) = 0$, i.e., $a = 0$ or $b = 1$. Therefore, the limits of α_t and β_t exist if and only if $\lim_{t \rightarrow \infty} \alpha_t = 0$ or $\lim_{t \rightarrow \infty} \beta_t = 1$. In the following, we will prove that the limit of β_t must be equal to 1.

Firstly, we prove a simpler result when $\beta_0 > 0$:

Lemma A.3. *For the discrete dynamical system in Eq. (A.56), if $\beta_0 \geq 0$, then $\lim_{t \rightarrow \infty} \beta_t = 1$.*

Proof. As aforementioned, due to the existence of the limit of α_t , we have $\lim_{t \rightarrow \infty} \alpha_t = 0$ or $\lim_{t \rightarrow \infty} \beta_t = 1$. Thus, we just need to prove that $\lim_{t \rightarrow \infty} \beta_t = 1$ as $\lim_{t \rightarrow \infty} \alpha_t = 0$.

When $\lim_{t \rightarrow \infty} \alpha_t = 0$, then there exists τ , such that $\forall t > \tau$, $\alpha_t \leq 1$.

According to the iterations in Eq. (A.56), we can derive that

$$\begin{aligned}
\frac{1 - \beta_{t+1}^2}{\alpha_{t+1}} &= \frac{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2 - (\alpha_t + \beta_t - \alpha_t \beta_t^2)^2}{\xi_t \alpha_t} \\
&= \frac{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2 - \alpha_t^2 - \beta_t^2 - \alpha_t^2 \beta_t^4 - 2\alpha_t \beta_t + 2\alpha_t^2 \beta_t^2 + 2\alpha_t \beta_t^3}{\xi_t \alpha_t} \\
&= \frac{1 - \beta_t^2 + \alpha_t^2 (\beta_t^2 - \beta_t^4) - 2\alpha_t (\beta_t - \beta_t^3)}{\xi_t \alpha_t} \\
&= \frac{1 - \beta_t^2}{\xi_t \alpha_t} \cdot (1 + \alpha_t^2 \beta_t^2 - 2\alpha_t \beta_t) \\
&= \frac{1 - \beta_t^2}{\alpha_t} \cdot \frac{(1 - \alpha_t \beta_t)^2}{\xi_t},
\end{aligned}$$

then $\forall t > \tau$,

$$\begin{aligned}
1 - \beta_{t+1}^2 &= \alpha_{t+1} \cdot \frac{1 - \beta_0^2}{\alpha_0} \prod_{i=0}^t \frac{(1 - \alpha_i \beta_i)^2}{\xi_i} \\
&= \alpha_{t+1} \cdot \frac{\eta(0)(1 - \beta_0^2)}{\alpha_0 \eta(t+1)} \cdot \prod_{i=0}^{\tau} (1 - \alpha_i \beta_i)^2 \cdot \prod_{i=\tau+1}^t (1 - \alpha_i \beta_i)^2 \quad (\text{A.58}) \\
&\leq \alpha_{t+1} \cdot \frac{\eta(0)(1 - \beta_0^2)}{\varepsilon \alpha_0} \cdot \prod_{i=0}^{\tau} (1 - \alpha_i \beta_i)^2,
\end{aligned}$$

where the inequality is based on the fact that $1 - \alpha_i \beta_i \in (0, 1]$ when $0 \leq \beta_0 \leq \beta_i \leq 1$, $\eta(t+1) \geq \varepsilon$, and $\alpha_i \leq 1$ for $i > \tau$. Since $\frac{\eta(0)(1 - \beta_0^2)}{\varepsilon \alpha_0} \cdot \prod_{i=0}^{\tau} (1 - \alpha_i \beta_i)^2$ is a constant, we obtain

$$\lim_{t \rightarrow \infty} 1 - \beta_{t+1}^2 \leq \lim_{t \rightarrow \infty} \alpha_{t+1} \cdot \frac{\eta(0)(1 - \beta_0^2)}{\varepsilon \alpha_0} \cdot \prod_{i=0}^{\tau} (1 - \alpha_i \beta_i)^2 = 0,$$

as $\lim_{t \rightarrow \infty} \alpha_{t+1} = 0$. This reveals $\lim_{t \rightarrow \infty} \beta_t^2 = 1$. Furthermore, since $\beta_t \geq 0$, we then have $\lim_{t \rightarrow \infty} \beta_t = 1$. \square

Next, we are going to prove $\lim_{t \rightarrow \infty} \beta_t = 1$ when $-1 < \beta_0 < 0$. According to Lemma A.3, we just need prove that $\exists \tau > 0$, s.t., $\beta_\tau \geq 0$.

For the sake of contradiction, suppose that $\beta_t < 0$ for all $t > 0$, we then have $\lim_{t \rightarrow \infty} \alpha_t = 0$. As a consequence, we obtain

$$\alpha_t + \beta_t - \alpha_t \beta_t^2 < 0, \quad \forall t \geq 0, \quad (\text{A.59})$$

and we know that $\exists t' > 0$, such that

$$\alpha_t < \frac{\varepsilon}{\eta(0)} \alpha_0 (1 - \beta_0^2), \quad \forall t \geq t'. \quad (\text{A.60})$$

According to the iterations in Eq. (A.56), we can derive that

$$\begin{aligned}
\alpha_{t+1}(1 - \beta_{t+1}^2) &= \frac{\xi_t \alpha_t}{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2} \left(1 - \frac{(\alpha_t + \beta_t - \alpha_t \beta_t^2)^2}{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2} \right) \\
&= \xi_t \alpha_t \left(\frac{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2 - (\alpha_t + \beta_t - \alpha_t \beta_t^2)^2}{(1 + \alpha_t^2 - \alpha_t^2 \beta_t^2)^2} \right) \\
&= \xi_t \alpha_t (1 - \beta_t^2) \left(\frac{1 - \alpha_t \beta_t}{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2} \right)^2 \\
&= \xi_t \alpha_t (1 - \beta_t^2) \left(1 - \frac{\alpha_t(\alpha_t + \beta_t - \alpha_t \beta_t^2)}{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2} \right)^2,
\end{aligned}$$

Since $1 - \frac{\alpha_t(\alpha_t + \beta_t - \alpha_t \beta_t^2)}{1 + \alpha_t^2 - \alpha_t^2 \beta_t^2} \geq 1$, then for $t \geq t'$,

$$\alpha_t \geq \alpha_t (1 - \beta_t^2) \geq \xi_{t-1} \alpha_{t-1} (1 - \beta_{t-1}^2) \geq \dots \geq \alpha_0 (1 - \beta_0^2) \prod_{i=0}^{t-1} \xi_i \geq \frac{\varepsilon}{\eta(0)} \alpha_0 (1 - \beta_0^2), \quad (\text{A.61})$$

which contradicts the fact in Eq. (A.60). Thus, $\exists \tau > 0$, s.t. $\beta_\tau \geq 0$. Consider the dynamical system with an initial time τ , we have $\lim_{t \rightarrow \infty} \beta_t = 1$ according to Lemma A.3.

To sum up, we have proven that $\lim_{t \rightarrow \infty} \beta_t = 1$ when $\beta_0 > -1$ and $\alpha_0 > 0$. \square

Theorem 3.5 *Considering the discrete dynamics in Eq. (3.16), if $\forall i \in [N], c \in [C]$, $\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) > -1$, the learning rate $\eta(t)$ satisfies that $\frac{\eta(t)}{\|\mathbf{h}_{i,c}(t)\|_2}$ is non-increasing, $\frac{\eta(0)(1+\gamma)}{CN\|\mathbf{h}_{i,c}(0)\|_2} \leq \frac{1}{\|\mathbf{w}_c\|_2}$, $\lim_{t \rightarrow \infty} \frac{\eta(t+1)}{\eta(t)} = 1$, and there exists a constant $\varepsilon > 0$, s.t., $\eta(t) > \varepsilon$, then we have*

$$\lim_{t \rightarrow \infty} \left\| \widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N^\top) \right\| = 0, \quad (\text{A.62})$$

and further if $\lim_{t \rightarrow \infty} \|\mathbf{H}(t)\| < \infty$, then there exists a constant $\mu > 0$, such that the error above shows exponential convergence:

$$\left\| \widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N^\top) \right\| \leq O(e^{-\mu t}). \quad (\text{A.63})$$

Moreover, if $\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) = -1$, then $\mathbf{h}_{i,c}(t) = \mathbf{h}_{i,c}(0)$.

Proof. Since $\mathbf{H}(t+1) = \mathbf{H}(t) + \frac{(1+\gamma)\eta(t)}{CN} \left(\frac{\partial \widehat{\mathbf{H}}}{\partial \mathbf{H}} \Big|_{\mathbf{H}=\mathbf{H}(t)} \right)^\top \mathbf{W}(\mathbf{I}_C \otimes \mathbf{1}_N^\top)$, then for $i \in [N]$, $c \in [C]$,

$$\mathbf{h}_{i,c}(t+1) = \mathbf{h}_{i,c}(t) + \frac{(1+\gamma)\eta(t)}{CN \|\mathbf{h}_{i,c}(t)\|_2} \left(\mathbf{I}_p - \widehat{\mathbf{h}}_{i,c}(t) \widehat{\mathbf{h}}_{i,c}^\top(t) \right) \mathbf{w}_c. \quad (\text{A.64})$$

According to Lemma A.2, when $\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) > -1$ and $\frac{\eta(0)(1+\gamma)}{CN} < \frac{\|\mathbf{h}_{i,c}(0)\|_2^2}{\|\mathbf{w}_c\|_2}$, we have $\lim_{t \rightarrow \infty} \|\widehat{\mathbf{h}}_{i,c} - \widehat{\mathbf{w}}_c\| = 0$, then $\lim_{t \rightarrow \infty} \left\| \widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N^\top) \right\| = 0$.

If further $\lim_{t \rightarrow \infty} \|\mathbf{H}(t)\|_2 < \infty$, let $L = \sup_{i,c,t} \|\mathbf{h}_{i,c}(t)\|$. According to the proof of Lemma A.2, $\forall i, c$, we have $\lim_{t \rightarrow \infty} \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(t) = 1$, then for a given constant $\delta > 0$, $\exists \tau > 0$, s.t., $\forall t > \tau$, $\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(t) \geq \delta$. Consider $t > \tau$, we have

$$\begin{aligned} & 1 - \left(\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(t+1) \right)^2 \\ &= \frac{\|\mathbf{h}_{i,c}(t)\|_2^2}{\|\mathbf{h}_{i,c}(t+1)\|_2^2} \left(1 - \left(\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(t) \right)^2 \right) \left(1 - \frac{(1+\gamma)\eta(t)\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(t)\|_2^2} \cdot \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(t) \right)^2 \\ &\leq \left(1 - \left(\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) \right)^2 \right) \prod_{j=0}^t \left(1 - \frac{(1+\gamma)\eta(j)\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(j)\|_2^2} \cdot \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(j) \right)^2 \\ &\leq \left(1 - \left(\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) \right)^2 \right) \prod_{j=0}^{\tau} \left(1 - \frac{(1+\gamma)\eta(j)\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(j)\|_2^2} \cdot \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(j) \right)^2 \prod_{j=\tau+1}^t \left(1 - \frac{(1+\gamma)\varepsilon\delta\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(j)\|_2^2} \right)^2 \\ &\leq \left(1 - \left(\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) \right)^2 \right) \prod_{j=0}^{\tau} \left(1 - \frac{(1+\gamma)\eta(j)\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(j)\|_2^2} \cdot \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(j) \right)^2 \left(1 - \frac{(1+\gamma)\varepsilon\delta\|\mathbf{w}_c\|_2}{CN L^2} \right)^{2(t-\tau)}. \end{aligned}$$

where the first, the second, and the third inequalities are based on the facts that $\frac{\|\mathbf{h}_{i,c}(t)\|_2^2}{\|\mathbf{h}_{i,c}(t+1)\|_2^2} \leq 1$, $1 - \frac{(1+\gamma)\eta(j+1)\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(j)\|_2^2} \cdot \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(j) \leq 1 - \frac{(1+\gamma)\varepsilon\delta\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(j)\|_2^2}$ for $t > \tau$, and $\|\mathbf{h}_{i,c}(j)\|_2 \leq L$, respectively.

Let $c_1 = \max_{i,c} \left(1 - \left(\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) \right)^2 \right) \prod_{j=0}^{\tau} \left(1 - \frac{(1+\gamma)\eta(j)\|\mathbf{w}_c\|_2}{CN \|\mathbf{h}_{i,c}(j)\|_2^2} \cdot \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(j) \right)^2 \left(1 - \frac{(1+\gamma)\varepsilon\delta\|\mathbf{w}_c\|_2}{CN L^2} \right)^{-2\tau}$, and $\mu = \min_c -2 \log \left(1 - \frac{(1+\gamma)\varepsilon\delta\|\mathbf{w}_c\|_2}{CN L^2} \right)$, then $1 - \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(t+1) \leq \frac{c_1 e^{-\mu t}}{1+\delta}$.

Therefore, we have

$$\left\| \widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N^\top) \right\|_2^2 = 2 \sum_{i,c} \left(1 - \widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(t+1) \right) \leq \frac{2c_1 C N e^{-\mu t}}{(1+\delta)}, \quad (\text{A.65})$$

i.e., $\left\| \widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N^\top) \right\| = O(e^{-\mu t})$.

Moreover, if $\widehat{\mathbf{w}}_c^\top \widehat{\mathbf{h}}_{i,c}(0) = -1$, we have

$$\mathbf{h}_{i,c}(t+1) = \mathbf{h}_{i,c}(t) + \frac{(1+\gamma)\eta(t)}{CN \|\mathbf{h}_{i,c}(t)\|_2} \left(\mathbf{I}_p - \widehat{\mathbf{h}}_{i,c}(t) \widehat{\mathbf{h}}_{i,c}^\top(t) \right) \mathbf{w}_c = \mathbf{h}_{i,c}(t), \quad (\text{A.66})$$

thus $\mathbf{h}_{i,c}(t) = \mathbf{h}_{i,c}(0)$. \square

A.7 PROOF OF THEOREM 4.1

Theorem 4.1 Consider the continual gradient flow (Equation (4.1)) in which the prototypes \mathbf{W} is fixed, we have the closed-form dynamics:

$$\mathbf{H}(t) = e^{-\lambda \int_0^t \eta(\tau) d\tau} \mathbf{H}(0) + \frac{1 - e^{-\lambda \int_0^t \eta(\tau) d\tau}}{\lambda} \mathbf{W} \mathbf{M}, \quad (\text{A.67})$$

which further indicates that $\|\mathbf{H}(t) - \frac{1}{\lambda} \mathbf{W} \mathbf{M}\| = O\left(e^{-\lambda \int_0^t \eta(\tau) d\tau}\right)$.

Proof. For the first order non-homogeneous linear difference equation in Eq. (4.1), the solution is

$$\begin{aligned} \mathbf{H}(t) &= e^{-\lambda \int_0^t \eta(\tau) d\tau} \left(\mathbf{H}(0) + \int_0^t \eta(s) e^{\lambda \int_0^s \eta(\tau) d\tau} ds \mathbf{W} \mathbf{M} \right) \\ &= e^{-\lambda \int_0^t \eta(\tau) d\tau} \mathbf{H}(0) + e^{-\lambda \int_0^t \eta(\tau) d\tau} \int_0^t \frac{1}{\lambda} de^{\lambda \int_0^s \eta(\tau) d\tau} \mathbf{W} \mathbf{M} \\ &= e^{-\lambda \int_0^t \eta(\tau) d\tau} \mathbf{H}(0) + e^{-\lambda \int_0^t \eta(\tau) d\tau} \frac{e^{\lambda \int_0^t \eta(\tau) d\tau} - 1}{\lambda} \mathbf{W} \mathbf{M} \\ &= e^{-\lambda \int_0^t \eta(\tau) d\tau} \mathbf{H}(0) + \frac{1 - e^{-\lambda \int_0^t \eta(\tau) d\tau}}{\lambda} \mathbf{W} \mathbf{M}, \end{aligned} \quad (\text{A.68})$$

and then $\|\mathbf{H}(t) - \frac{1}{\lambda} \mathbf{W} \mathbf{M}\| = \|e^{-\lambda \int_0^t \eta(\tau) d\tau} (\mathbf{H}(0) - \frac{\mathbf{W} \mathbf{M}}{\lambda})\| = O(e^{-\lambda \int_0^t \eta(\tau) d\tau})$. \square

A.8 THE PROJECTIONS ONTO \mathcal{E}_1^+ , \mathcal{E}_1^- , \mathcal{E}_2^+ , \mathcal{E}_2^- AND \mathcal{E}_3

Lemma A.4. Let \mathcal{S} denote the subspace $\mathcal{S} = \{\mathbf{W} : \mathbf{W} \mathbf{1}_n = 0, \mathbf{W} \in \mathbb{R}^{m \times n}\}$, then the projection of a point $\mathbf{A} \in \mathbb{R}^{m \times n}$ onto \mathcal{S} can be denoted as $\Pi_{\mathcal{S}} \mathbf{A} = \mathbf{A} (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top)$.

Proof. Let $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n) \in \mathcal{S}$ and $\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n) \in \mathbb{R}^{m \times n}$, we have

$$\|\mathbf{W} - \mathbf{A}\|_F^2 = \sum_{i=1}^n \|\mathbf{w}_i - \mathbf{a}_i\|_2^2 \geq \frac{1}{n} \left\| \sum_{i=1}^n \mathbf{w}_i - \sum_{i=1}^n \mathbf{a}_i \right\|_2^2 = \frac{1}{n} \|\mathbf{W} \mathbf{1}_n - \mathbf{A} \mathbf{1}_n\|_2^2 = \frac{1}{n} \|\mathbf{A} \mathbf{1}_n\|_2^2$$

where we used the Cauchy-Schwarz inequality, and the equality holds if and only if $\mathbf{w}_i - \mathbf{a}_i = \mathbf{w}_n - \mathbf{a}_n, \forall i \in [n]$, and $\sum_{i=1}^n \mathbf{w}_i = 0$, i.e., $\mathbf{W} = \mathbf{A} (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top)$. Therefore, the projection of \mathbf{A} onto \mathcal{S} is $\Pi_{\mathcal{S}} \mathbf{A} = \arg \min_{\mathbf{W} \in \mathcal{S}} \|\mathbf{W} - \mathbf{A}\|_F^2 = \mathbf{A} (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top)$. \square

Lemma A.5. Let \mathcal{S} denote the subspace $\{\mathbf{W} : \mathbf{W} (\mathbf{I}_c \otimes \mathbf{1}_n) = 0, \mathbf{W} \in \mathbb{R}^{m \times cn}\}$, then the projection of a point $\mathbf{A} \in \mathbb{R}^{m \times cn}$ onto \mathcal{S} can be denoted as $\Pi_{\mathcal{S}} \mathbf{A} = \mathbf{A} (\mathbf{I}_{cn} - \frac{1}{n} \mathbf{I}_c \otimes \mathbf{1}_n \mathbf{1}_n^\top)$.

Proof. The proof is similar to Lemma A.4. We can simply let $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_n) \in \mathcal{S}$ and $\mathbf{A} = (\mathbf{A}_1, \dots, \mathbf{A}_n) \in \mathbb{R}^{m \times cn}$, where $\mathbf{W}_i, \mathbf{A}_i \in \mathbb{R}^{m \times c}$, then we have

$$\|\mathbf{W} - \mathbf{A}\|_F^2 = \sum_{i=1}^n \|\mathbf{W}_i - \mathbf{A}_i\|_F^2 \geq \frac{1}{n} \left\| \sum_{i=1}^n \mathbf{W}_i - \sum_{i=1}^n \mathbf{A}_i \right\|_2^2 = \frac{1}{n} \|\mathbf{A} (\mathbf{I}_c \otimes \mathbf{1}_n)\|_2^2$$

where the equality holds if and only if $\mathbf{W}_i - \mathbf{A}_i = \mathbf{W}_n - \mathbf{A}_n, \forall i \in [n]$, and $\sum_{i=1}^n \mathbf{W}_i = 0$, i.e., $\mathbf{W} = \mathbf{A} (\mathbf{I}_{cn} - \frac{1}{n} \mathbf{I}_c \otimes \mathbf{1}_n \mathbf{1}_n^\top)$. Therefore, the projection of \mathbf{A} onto \mathcal{S} is $\Pi_{\mathcal{S}} \mathbf{A} = \mathbf{A} (\mathbf{I}_{cn} - \frac{1}{n} \mathbf{I}_c \otimes \mathbf{1}_n \mathbf{1}_n^\top)$. \square

Lemma A.6. For $\mathbf{H} \in \mathbb{R}^{p \times CN}$, $\mathbf{W} \in \mathbb{R}^{p \times C}$, the projection of (\mathbf{H}, \mathbf{W}) onto \mathcal{E}_1^ϵ is

$$\Pi_1^\epsilon(\mathbf{H}, \mathbf{W}) = \left(\frac{\epsilon}{\sqrt{N}} (\mathbf{P} \otimes \mathbf{1}_N), \mathbf{P} \right), \quad (\text{A.69})$$

where $\mathbf{P} = \frac{1}{2} \left(\frac{\epsilon}{\sqrt{N}} \mathbf{H} (\mathbf{I}_C \otimes \mathbf{1}_N) + \mathbf{W} \right) (\mathbf{I}_C - \frac{1}{C} \mathbf{1}_C \mathbf{1}_C^\top)$.

Proof. Let $\mathcal{S} = \{\mathbf{Z} : \mathbf{Z}\mathbf{1}_C = 0, \mathbf{Z} \in \mathbb{R}^{p \times C}\}$ and $\mathbf{H} = \{\mathbf{H}_1, \dots, \mathbf{H}_N\}$ (where $\mathbf{H}_i \in \mathbb{R}^{p \times C}$), the minimizer of $\mathbf{Z} \in \mathcal{S}$ is

$$\begin{aligned}
& \arg \min_{\mathbf{Z} \in \mathcal{S}} \left\| \frac{\epsilon}{\sqrt{N}} (\mathbf{Z} \otimes \mathbf{1}_N^\top) - \mathbf{H} \right\|_F^2 + \|\mathbf{Z} - \mathbf{W}\|_F^2 \\
&= \arg \min_{\mathbf{Z} \in \mathcal{S}} \sum_{i=1}^N \left\| \frac{1}{\sqrt{N}} \mathbf{Z} - \epsilon \mathbf{H}_i \right\|_F^2 + \|\mathbf{Z} - \mathbf{W}\|_F^2 \\
&= \arg \min_{\mathbf{Z} \in \mathcal{S}} \|\mathbf{Z}\|_F^2 - \frac{2\epsilon}{\sqrt{N}} \sum_{i=1}^N \langle \mathbf{Z}, \mathbf{H}_i \rangle + \|\mathbf{H}\|_F^2 + \|\mathbf{Z} - \mathbf{W}\|_F^2 \\
&= \arg \min_{\mathbf{Z} \in \mathcal{S}} \left\| \mathbf{Z} - \frac{\epsilon}{\sqrt{N}} \sum_{i=1}^N \mathbf{H}_i \right\|_F^2 - \left\| \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{H}_i \right\|_F^2 + \|\mathbf{H}\|_F^2 + \|\mathbf{Z} - \mathbf{W}\|_F^2 \\
&= \arg \min_{\mathbf{Z} \in \mathcal{S}} \left\| \mathbf{Z} - \frac{\epsilon}{\sqrt{N}} \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) \right\|_F^2 + \|\mathbf{Z} - \mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2 - \left\| \frac{\epsilon}{\sqrt{N}} \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) \right\|_F^2 \\
&= \arg \min_{\mathbf{Z} \in \mathcal{S}} \left\| \mathbf{Z} - \frac{1}{2} \left(\frac{\epsilon}{\sqrt{N}} \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) + \mathbf{W} \right) \right\|_F^2 \\
&= \frac{1}{2} \left(\frac{\epsilon}{\sqrt{N}} \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) + \mathbf{W} \right) (\mathbf{I}_C - \frac{1}{C} \mathbf{1}_C \mathbf{1}_C^\top).
\end{aligned} \tag{A.70}$$

Thus, $\Pi_1^\epsilon(\mathbf{H}, \mathbf{W}) = \left(\frac{\epsilon}{\sqrt{N}} (\mathbf{P} \otimes \mathbf{1}_N^\top), \mathbf{P} \right)$ with $\mathbf{P} = \frac{1}{2} \left(\frac{\epsilon}{\sqrt{N}} \mathbf{H}(\mathbf{I}_C \otimes \mathbf{1}_N) + \mathbf{W} \right) (\mathbf{I}_C - \frac{1}{C} \mathbf{1}_C \mathbf{1}_C^\top)$. \square

Lemma A.7. For $\mathbf{H} \in \mathbb{R}^{p \times CN}$, $\mathbf{W} \in \mathbb{R}^{p \times C}$, the projection of (\mathbf{H}, \mathbf{W}) onto \mathcal{E}_2^ϵ is

$$\Pi_2^\epsilon(\mathbf{H}, \mathbf{W}) = \left(\frac{\epsilon}{\sqrt{N}} \mathbf{h} \mathbf{1}_{CN}^\top, \mathbf{h} \mathbf{1}_C^\top \right), \tag{A.71}$$

where $\mathbf{h} = \frac{1}{2C} \left(\frac{\epsilon}{\sqrt{N}} \mathbf{H} \mathbf{1}_{CN} + \mathbf{W} \mathbf{1}_C \right)$.

Proof. We have

$$\begin{aligned}
& \arg \min_{\mathbf{h} \in \mathbb{R}^p} \left\| \frac{\epsilon}{\sqrt{N}} \mathbf{h} \mathbf{1}_{CN}^\top - \mathbf{H} \right\|_F^2 + \|\mathbf{h} \mathbf{1}_C^\top - \mathbf{W}\|_F^2 \\
&= \arg \min_{\mathbf{h} \in \mathbb{R}^p} CN \left\| \frac{1}{\sqrt{N}} \mathbf{h} - \frac{\epsilon}{CN} \mathbf{H} \mathbf{1}_{CN} \right\|_2^2 + C \|\mathbf{h} - \frac{1}{C} \mathbf{W} \mathbf{1}_C\|_2^2 \\
&= \arg \min_{\mathbf{h} \in \mathbb{R}^p} \left\| \mathbf{h} - \frac{\epsilon}{C\sqrt{N}} \mathbf{H} \mathbf{1}_{CN} \right\|_2^2 + \left\| \mathbf{h} - \frac{1}{C} \mathbf{W} \mathbf{1}_C \right\|_2^2 \\
&= \arg \min_{\mathbf{h} \in \mathbb{R}^p} \left\| \mathbf{h} - \frac{1}{2C} \left(\frac{\epsilon}{\sqrt{N}} \mathbf{H} \mathbf{1}_{CN} + \mathbf{W} \mathbf{1}_C \right) \right\|_2^2 \\
&= \frac{1}{2C} \left(\frac{\epsilon}{\sqrt{N}} \mathbf{H} \mathbf{1}_{CN} + \mathbf{W} \mathbf{1}_C \right)
\end{aligned} \tag{A.72}$$

\square

Lemma A.8. For $\mathbf{H} \in \mathbb{R}^{p \times CN}$, $\mathbf{W} \in \mathbb{R}^{p \times C}$, the projection of (\mathbf{H}, \mathbf{W}) onto \mathcal{E}_3 is

$$\Pi_3(\mathbf{H}, \mathbf{W}) = \left(\mathbf{H}(\mathbf{I}_{CN} - \frac{1}{N} \mathbf{I}_C \otimes \mathbf{1}_N \mathbf{1}_N^\top), 0 \right). \tag{A.73}$$

Proof. This can be easily derived by Lemma A.5. \square

B EXPERIMENTS

In this section, we provide experimental details, including datasets, network architectures, optimization methods, hyperparameter settings, and more results.

B.1 NUMERICAL EXPERIMENTS

For numerical experiments in Figures 1, 4, 5, 6, 7, 8, and 9, we set $p = 512$, $C = 100$, $N = 10$, and then randomly initialize \mathbf{H}_0 and \mathbf{W}_0 . We use the SGD optimizer to optimize these free variables.

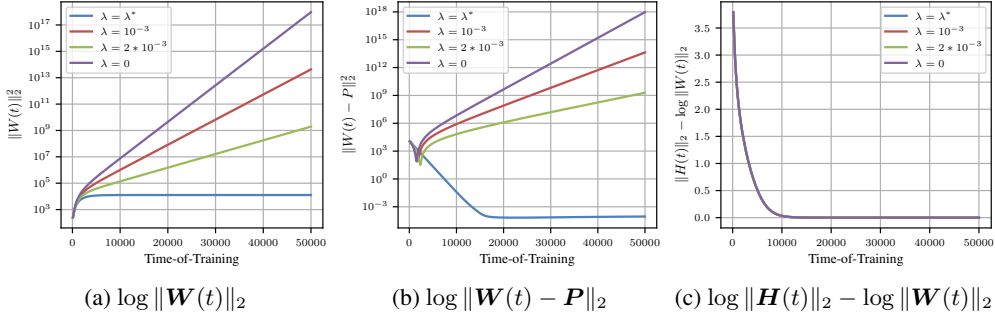


Figure 4: Behavior of gradient descent iterates under the Averaged Sample Margin loss with different weight decay coefficients and $\eta_1(t) = \eta_2(t) = 0.1$ (i.e., $s = 1$), where \mathbf{P} denotes the component \mathbf{W}_1^+ in the projection $\Pi_1^+ \mathbf{Z}_0$ calculated according to Lemma A.6. (a) The logarithm of the norm of $\mathbf{W}(t)$. As expected, the norm increases exponentially when $\lambda < \lambda^*$; (b) The difference between $\mathbf{W}(t)$ and \mathbf{P} . As expected (Corollary 3.4), $\mathbf{W}(t)$ converges to \mathbf{P} when $\lambda = \lambda^*$, while other differences are dominated by $\|\mathbf{W}(t)\|_2$; (c) The difference in ℓ_2 norm between $\mathbf{H}(t)$ and $\mathbf{W}(t)$. The convergence is the same even if the weight decay is different.

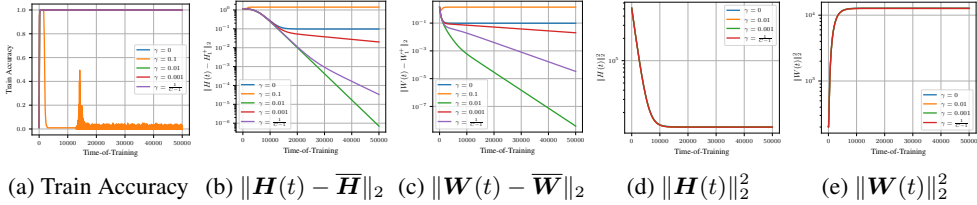


Figure 5: Verification of the behavior of regularized gradient descent iterates in Equation (3.11) with $\gamma \in \{0, 0.1, 0.01, 0.001, \frac{1}{C-1}\}$. We set $p = 512$, $C = 100$, $N = 10$, $\lambda = \frac{(1+\gamma)}{C\sqrt{N}}$, $\eta_1(t) = \eta_2(t) = 0.1$ (i.e., $s = 1$), thus we have $\lim_{t \rightarrow \infty} \mathbf{Z}(t) = \Pi_1^+ \mathbf{Z}_0$, according to Corollary 3.4, and then randomly initialize \mathbf{H}_0 and \mathbf{W}_0 . (a) The training accuracy with the prediction rule $\arg \max_c \mathbf{w}_c^\top \mathbf{h}$. As expected, the features align to their corresponding prototypes when $\gamma < \frac{2}{C-2}$. (b) The ℓ_2 distance between $\mathbf{H}(t)$ and \mathbf{H}_1^+ . As expected Theorem 3.4, the distance will decrease as exponential rate when $0 < \gamma < \frac{2}{C-2}$. (c) The ℓ_2 distance between $\mathbf{W}(t)$ and \mathbf{W}_1^+ . (d) and (e) denote the norm of features and prototypes, respectively. As can be seen, $\|\mathbf{H}\|_2$ and $\|\mathbf{W}\|_2$ do not grow exponentially as in the unconstrained case, which confirms that weight decay can avoid excessive growth of feature norm and prototype norm.

B.2 VISUAL CLASSIFICATION

For classification experiments in Fig. 2, Fig. 12, Fig. 13, Fig. 14, we experiment with ResNet-18, ResNet-34, and ResNet-50 (He et al., 2016) trained on CIFAR-10, CIFAR-100 (Krizhevsky and Hinton, 2009), and ImageNet-100 that takes the first 100 classes of ImageNet (Deng et al., 2009), respectively. The networks are trained for 200 epochs and 100 epochs for CIFAR-10/-100 and ImageNet-100, respectively. For all training, we use SGD optimizer with momentum 0.9 and cosine

Table 2: Test accuracies on imbalanced CIFAR-10 under different explicit feature regularization.

Dataset	Imbalanced CIFAR-10							
	long-tailed				step			
Imbalance Ratio	100	50	20	10	100	50	20	10
baseline	67.81	72.93	83.97	88.37	61.24	68.10	78.73	85.49
$\lambda = 5e-6$	67.84	72.85	83.17	89.06	60.79	68.41	80.20	86.69
$\lambda = 1e-5$	67.74	76.14	84.17	89.19	61.50	67.71	80.97	87.18
$\lambda = 5e-5$	69.74	77.29	84.92	88.64	60.69	70.27	81.27	87.17

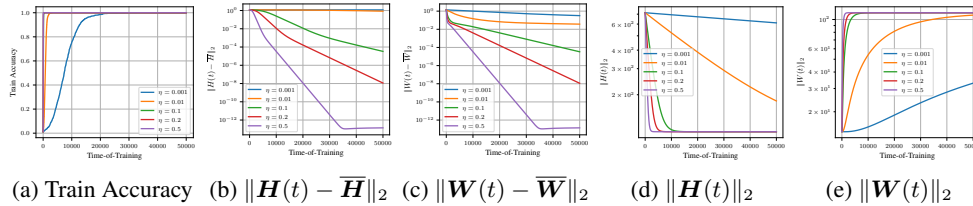


Figure 6: Verification of the behavior of regularized gradient descent iterates in Equation (3.11) with different learning rates ($\eta \in \{0.001, 0.01, 0.1, 0.2, 0.5\}$). We set $p = 512$, $C = 100$, $N = 10$, $\eta_1 = \eta_2 = \eta$ ($s = 1$), $\gamma = \frac{1}{C-1}$, and $\lambda = \frac{1+\gamma}{C\sqrt{N}}$. As can be seen, features and prototypes converge to $(\bar{\mathbf{H}}, \bar{\mathbf{W}})$ exponentially, and larger learning rates can accelerate convergence.

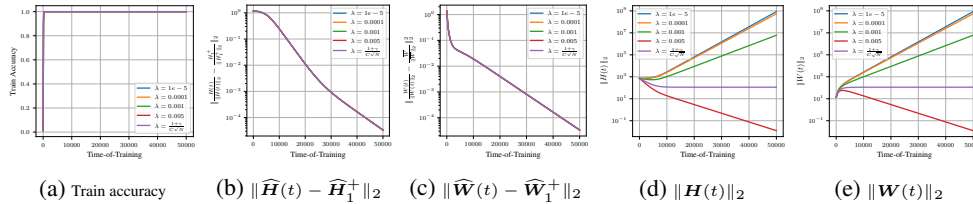


Figure 7: Verification of the behavior of regularized gradient descent iterates in Equation (3.14) with different weight decay coefficients ($\lambda = \{1e - 5, 1e - 4, 1e - 3, 5e - 3, \frac{1+\gamma}{C\sqrt{N}}\}$). We set $p = 512$, $C = 100$, $N = 10$, $\eta_1(t) = \eta_2(t) = 0.1$ (i.e., $s = 1$), $\gamma = \frac{1}{C-1}$, where $\bar{\mathbf{H}} = \pi_h^+ \mathbf{H}_1^+ + \pi_h^- \mathbf{H}_1^-$ and $\bar{\mathbf{W}} = \pi_w^+ \mathbf{W}_1^+ + \pi_w^- \mathbf{W}_1^-$ in Corollary 3.4. (a) The logarithm of the norm of $\mathbf{W}(t)$. As expected, the norm increases exponentially when $\lambda < \lambda^* = \frac{1+\gamma}{C\sqrt{N}}$; (b) The difference between $\mathbf{W}(t)$ and \mathbf{P} . As expected in Corollary 3.4, $\mathbf{W}(t)$ converges to \mathbf{P} when $\lambda = \lambda^*$, while other differences are dominated by $\|\mathbf{W}(t)\|_2$; (c) The difference in ℓ_2 norm between $\mathbf{H}(t)$ and $\mathbf{W}(t)$. The convergence is the same even if the weight decay is different.

learning rate annealing Loshchilov and Hutter (2017) with T_{\max} being the corresponding epochs. The initial learning rate is set to 0.1, weight decay is set to 5×10^{-4} , and batch size is set to 256. Typical data augmentations including random width/height shift and horizontal flip are applied. Moreover, to use the PAL and FNPAL (Zhou et al., 2022b) that anchors prototypes with a neural collapse solution, we remove the ReLU layer before the linear classifier in the last layer.

B.3 IMBALANCED CLASSIFICATION

For the experiments of imbalanced learning in Tab. 1, Tab. 2, Tab. 3, Fig. 10, and Fig. 11, we utilize the same network architectures, and optimization settings as visual classification. We only use the imbalanced versions of CIFAR-10 and CIFAR-100 by following the setting in (Zhou et al., 2022c). The number of training examples is reduced for per class, and the test set keeps unchanged, where we use the imbalance ratio $\rho = \frac{\max_i n_i}{\min_i n_i}$ to denote the ratio between sample sizes of the most frequent and least frequent class. Moreover, long-tailed imbalance (Cui et al., 2019) that utilizes

Table 3: Test accuracies on imbalanced CIFAR-100 under different explicit feature regularization.

Dataset	Imbalanced CIFAR-100							
	long-tailed				step			
Imbalance Ratio	100	50	20	10	100	50	20	10
baseline	33.37	39.40	42.96	56.38	40.89	42.69	51.92	57.52
$\lambda = 5e - 6$	36.00	41.92	50.75	60.13	41.90	43.85	47.80	56.74
$\lambda = 1e - 5$	36.61	42.36	49.21	58.91	41.48	43.77	49.64	56.49
$\lambda = 5e - 5$	34.88	42.74	54.72	60.84	40.97	43.20	48.96	57.97

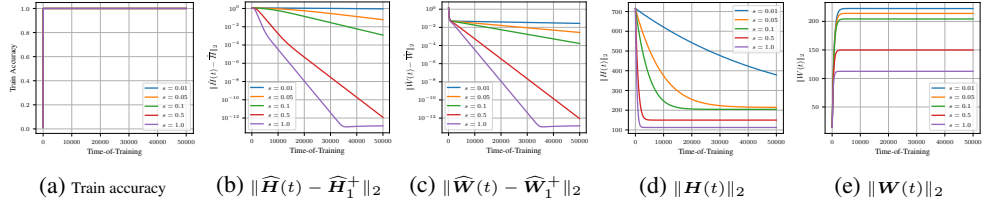


Figure 8: Verification of the behavior of regularized gradient descent iterates in Equation (3.14) with different scale parameters $s = \frac{\eta_1(0)}{\eta_2(0)} \in \{0.01, 0.05, 0.1, 0.5, 1.0\}$. We set $p = 512$, $C = 100$, $N = 10$, $\eta_2 = 0.5$, $\gamma = \frac{1}{C-1}$, where $\widehat{\mathbf{H}} = \pi_h^+ \mathbf{H}_1^+ + \pi_h^- \mathbf{H}_1^-$ and $\widehat{\mathbf{W}} = \pi_w^+ \mathbf{W}_1^+ + \pi_w^- \mathbf{W}_1^-$ in Corollary 3.4. As can be seen, larger scale parameter s can achieve faster convergence speed.

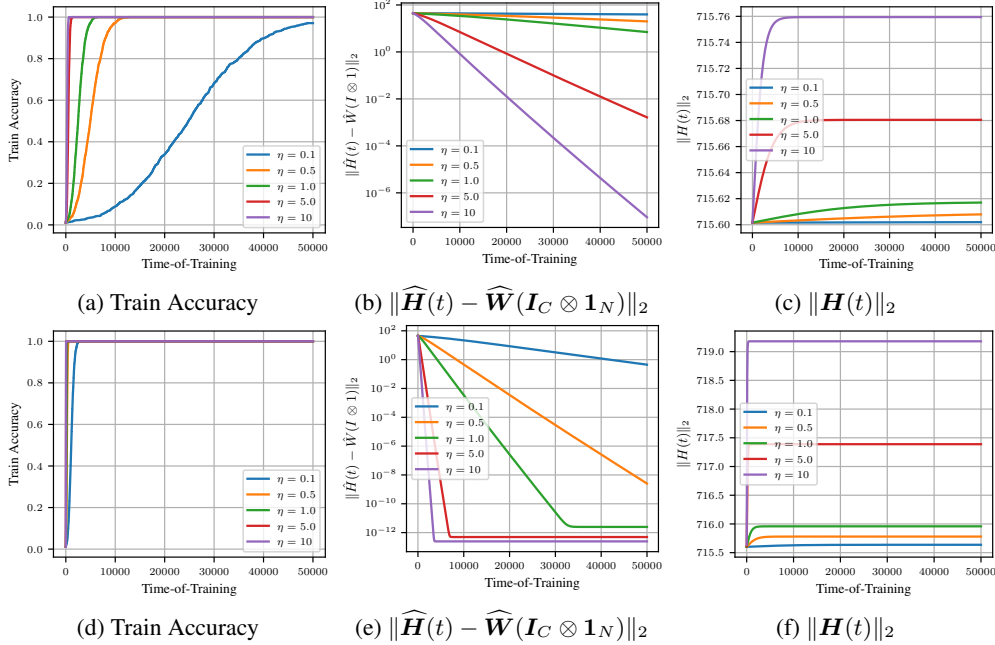


Figure 9: Verification of the behavior of discrete gradient descent iterates in Equation (3.15) under anchored prototypes with different learning rates $\eta \in \{0.1, 0.5, 1.0, 5.0, 10\}$ and without (a-c) or with (d-f) rescaled learning rates. We set $p = 512$, $C = 100$, and $N = 10$. As expected in Theorem 3.5, the feature norm $\|\mathbf{H}(t)\|_2$ is non-decreasing, and the error $\|\widehat{\mathbf{H}}(t) - \widehat{\mathbf{W}}(\mathbf{I}_C \otimes \mathbf{1}_N)\|_2$ shows exponential decrease.

an exponential decay in samples sizes and step imbalance (Buda et al., 2018)(that sets all minority classes to have the same number of samples, as do all majority classes) are considered.

For imbalanced learning, we utilize expected calibration error (ECE) to measure calibration of the models (Zhong et al., 2021), where all predictions are grouped into several interval bins of equal size and then calculate the error between the accuracy and confidence for each interval bin, *i.e.*,

$$\text{ECE} = \sum_{b=1}^B \frac{|\mathcal{S}_b|}{N} |\text{acc}(\mathcal{S}_b) - \text{conf}(\mathcal{S}_b)| \times 100\%, \quad (\text{B.1})$$

where N denotes the number of predictions, B is the number of interval bins, \mathcal{S}_b is the set of samples whose prediction scores fall into Bin- b , $\text{acc}(\cdot)$ and $\text{conf}(\cdot)$ denote the accuracy and predicted confidence of \mathcal{S}_b , respectively.

As shown in Tab. 2 and Tab. 3, explicit feature regularization can improve imbalanced learning on CIFAR-10/100 in most cases.

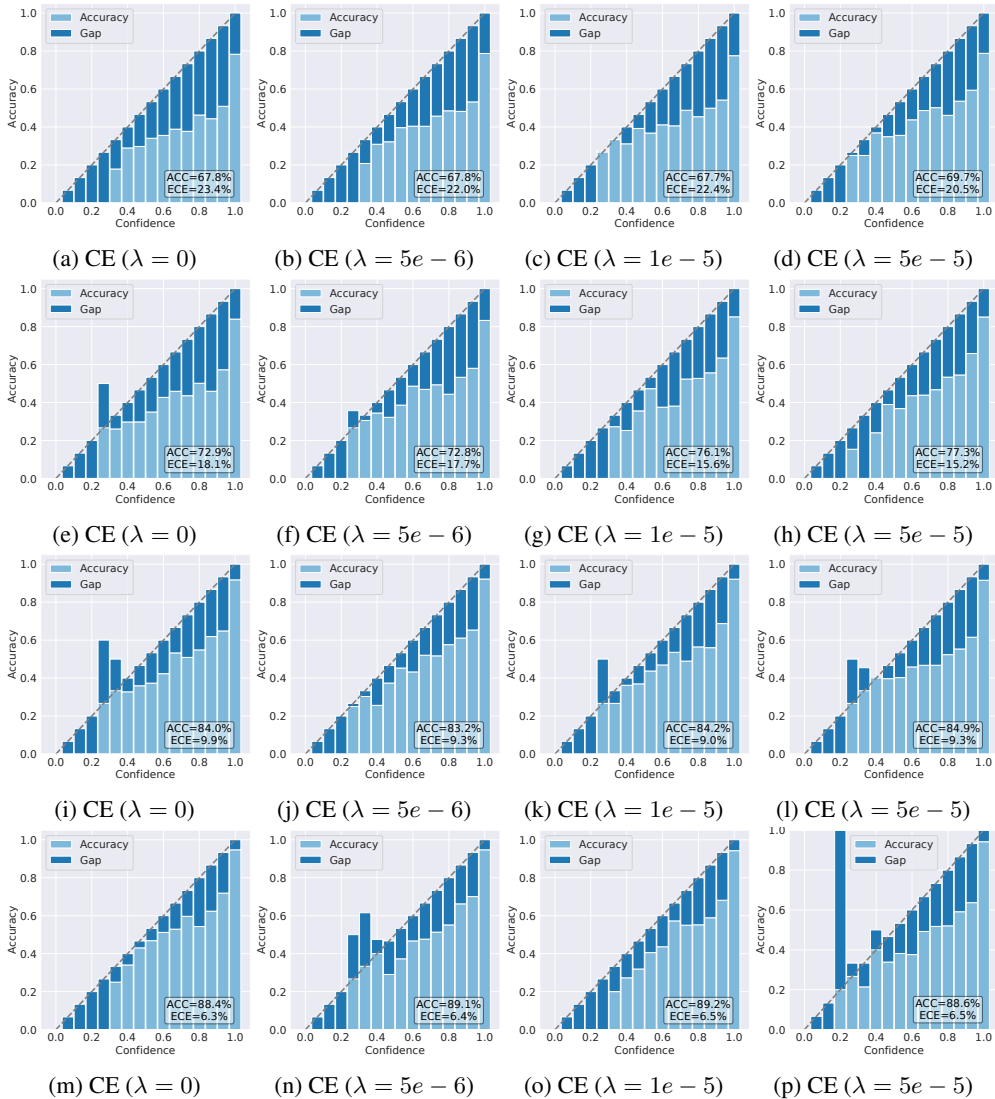


Figure 10: Reliability diagrams of ResNet-34 (He et al., 2016) trained by CE on CIFAR-10-LT with imbalance ratio $\rho \in \{100, 50, 20, 10\}$ under different explicit feature regularization ($\lambda \in \{0, 5e-6, 1e-5, 5e-5\}$). As can be seen, an appropriate larger weight decay can improve both accuracy and confidence

B.4 OUT-OF-DISTRIBUTION DETECTION

For the experiments of OOD detection in Fig. 3, Tab. 4, Fig. 16, Fig. 17, and Fig. 18, we use a ResNet-18 on CIFAR-10 and a ResNet-34 on CIFAR-100 to train the classification models, and use their test dataset as the in-distribution data \mathcal{D}_{in}^{test} . For the OOD test dataset \mathcal{D}_{out}^{test} , we simply use a common benchmark: SVHN (Netzer et al., 2011). We measure the performance with the following metrics: (1) the false positive rate (FPR95) of OOD examples when true positive rate of in-distribution examples is at 95%; (2) the area under the receiver operating characteristic curve (AUROC); and (3) the area under the precision-recall curve (AUPR). We then consider the softmax-based score (Hendrycks and Gimpel, 2016), energy-based score (Liu et al., 2020), and our proposed feature norm-based score to assessing the improvement of explicit feature regularization over the normal training.

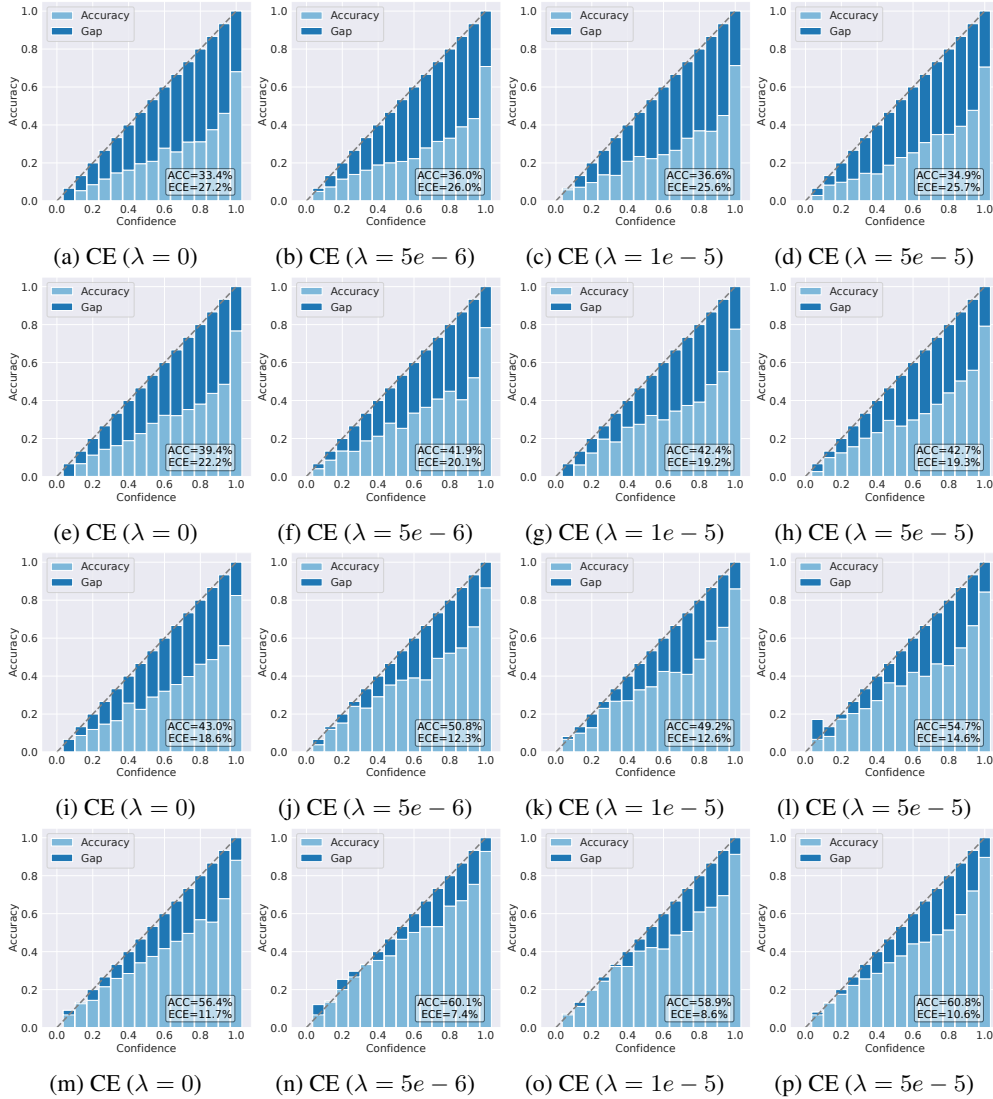


Figure 11: Reliability diagrams of ResNet-34 (He et al., 2016) trained by CE on CIFAR-100-LT with imbalance ratio $\rho \in \{100, 50, 20, 10\}$ under different explicit feature regularization ($\lambda \in \{0.0, 5e-6, 1e-5, 5e-5\}$), where ECE denotes the expected calibration error (Zhong et al., 2021). As can be seen, an appropriate larger weight decay can improve both accuracy and confidence calibration.

C OTHER POTENTIAL INSIGHTS

C.1 A GOOD INITIALIZATION OF PROTOTYPES

As depicted in Sec. 3 and Appendix A.8, the dynamics under the ASM loss is dependent on the initialization of both features and prototypes, such as $\Pi_1^+ \mathbf{Z}_0 = (\frac{1}{\sqrt{N}}(\mathbf{P} \otimes \mathbf{1}_N^T), \mathbf{P})$, where $\mathbf{P} = \frac{1}{2} \left(\frac{1}{\sqrt{N}} \mathbf{H}_0 (\mathbf{I}_C \otimes \mathbf{1}_N) + \mathbf{W}_0 \right) (\mathbf{I}_C - \frac{1}{C} \mathbf{1}_C \mathbf{1}_C^T)$. However, these features \mathbf{H}_0 extracted from a dataset by some nonlinear layers and parameterized layers are practically intractable, but we can elaborately initialize \mathbf{W}_0 and highlight its role in the whole. To do this, we consider two ways: (1) Initializing the structure of \mathbf{W}_0 . Inspired by the neural collapse solution that maximizes class separation, we can initialize \mathbf{W}_0 as this structure, *i.e.*, $\hat{\mathbf{w}}_i^T \hat{\mathbf{w}}_j = \frac{1}{C-1}, \forall i \neq j$; (2) Increasing the importance of \mathbf{W}_0 . A simple strategy is scaling up \mathbf{W}_0 , thereby implicitly weakening the importance of \mathbf{H}_0 . However, it is difficult to handle the initialization of features because they are obtained by a complex processing a large dataset, thus we seek to initialize the prototypes in the last layer of the network.

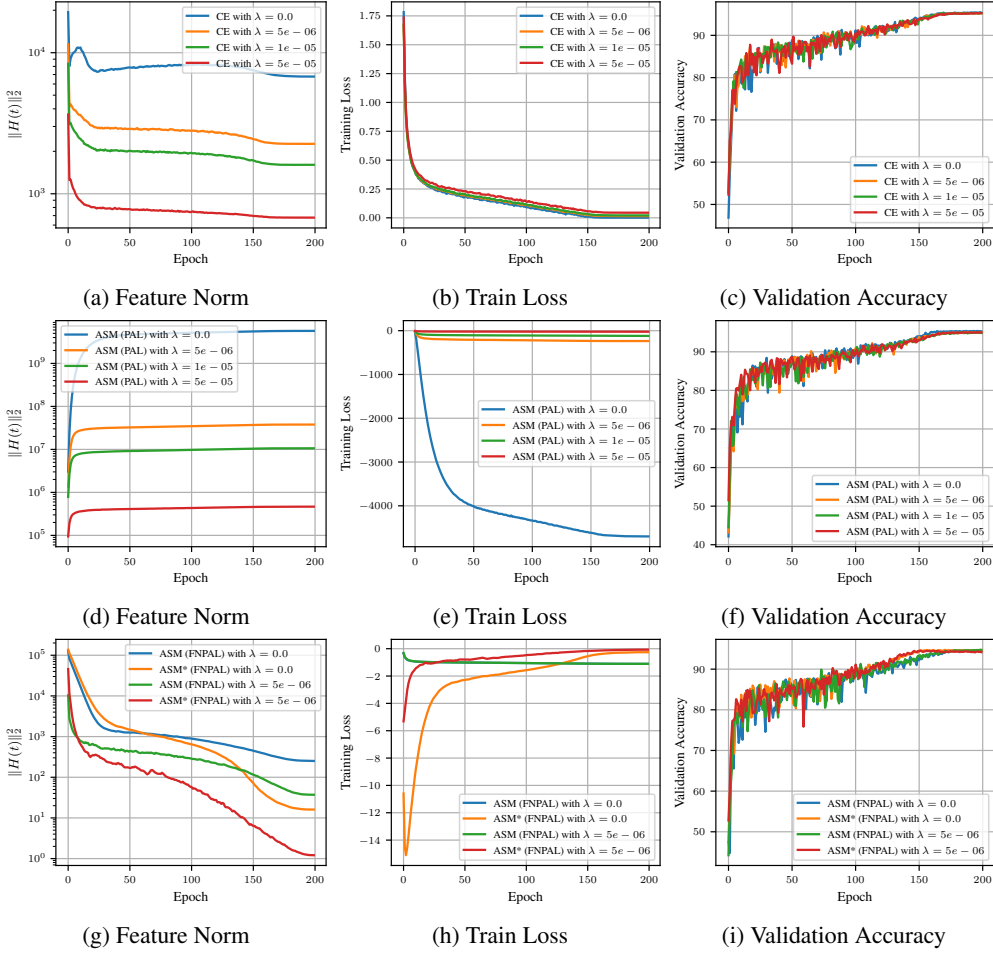


Figure 12: Behavior of visual classification on CIFAR-10 with CE, ASM (PAL) and ASM (FNPAL) under different weight decay coefficients.

C.2 REFINED DECISION-MAKINGS

Recalling the rule— $\arg \max_{c'} \langle \mathbf{w}_{c'}, \mathbf{h} \rangle + b_{c'}$ that makes decision by selecting the class with the largest logit (where the inner product $\langle \mathbf{w}_{c'}, \mathbf{h} \rangle$ is dominant), which may not be good to directly use the learned features and prototypes, since learning with the ASM within limited iterations (that means $\zeta_1(t) < \infty$) will introduce some residual $\Delta(t)$ caused by gradient descent regardless of the unconstrained case or regularized case.

Example C.1. *If we add a perturbation Δ for all features while adding $s\Delta$ for all prototypes, then the perturbed decision-making will be $\arg \max_{c'} \langle \mathbf{w}_{c'} + s\Delta, \mathbf{h} + \Delta \rangle + b_{c'}$, which may not be equivalent to $\arg \max_{c'} \langle \mathbf{w}_{c'}, \mathbf{h} \rangle + b_{c'}$.*

C.3 ADJUSTED SAMPLE MARGIN LOSS

As aforementioned in Sec. 3.3 and the proof of Theorem 3.5, we will encounter zero gradients when the cosine similarity $\widehat{\mathbf{w}}_y^\top \widehat{\mathbf{h}}$ is -1 or 1 , so we can adjust the loss to avoid the issue by to \mathbf{w}_y and accelerate convergence:

$$L'_{ASM}(\mathbf{W}\widehat{\mathbf{h}}, y) = \begin{cases} L_{ASM}(\mathbf{W}\widehat{\mathbf{h}}, y) & \text{if } \widehat{\mathbf{w}}_y^\top \widehat{\mathbf{h}} \geq -1 + \varepsilon, \\ -(1 + \gamma)(\mathbf{w}_y + \delta)^\top \widehat{\mathbf{h}} & \text{if } \widehat{\mathbf{w}}_y^\top \widehat{\mathbf{h}} < -1 + \varepsilon, \end{cases} \quad (\text{C.1})$$

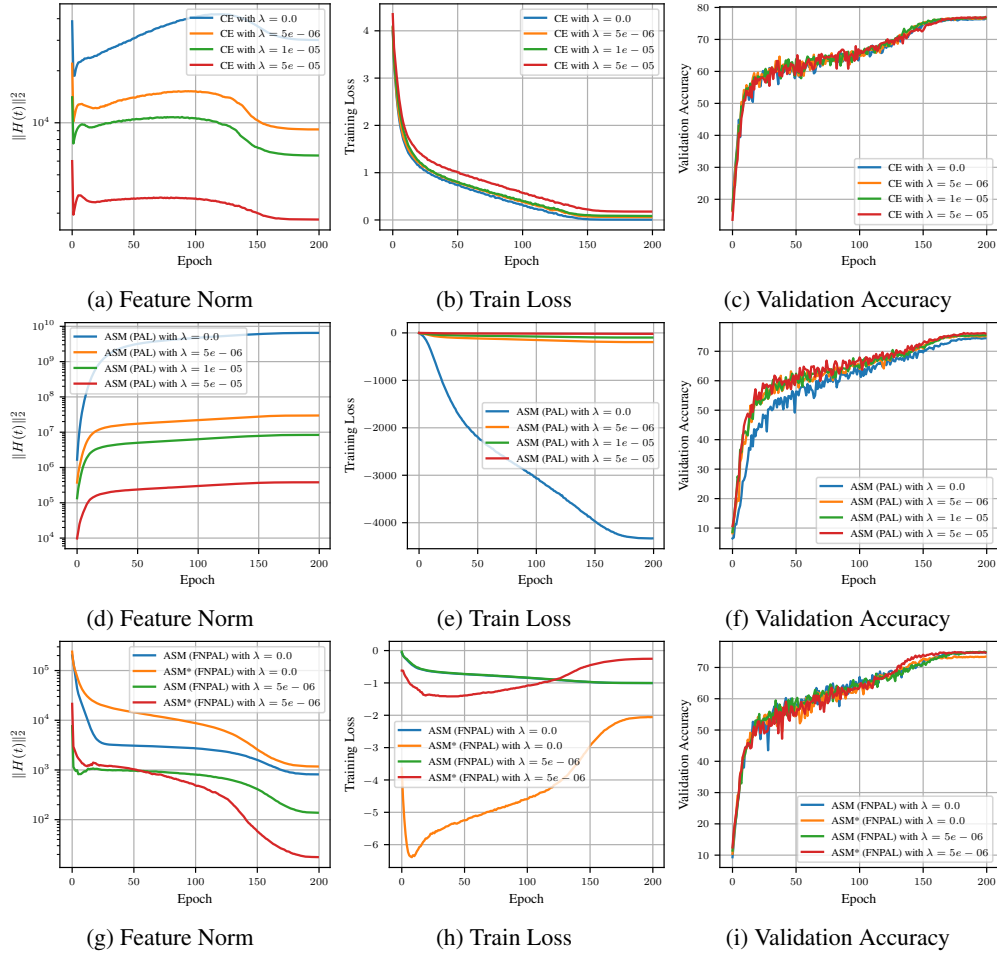


Figure 13: Behavior of visual classification on CIFAR-100 with CE, ASM (PAL) and ASM (FNPAL) under different weight decay coefficients.

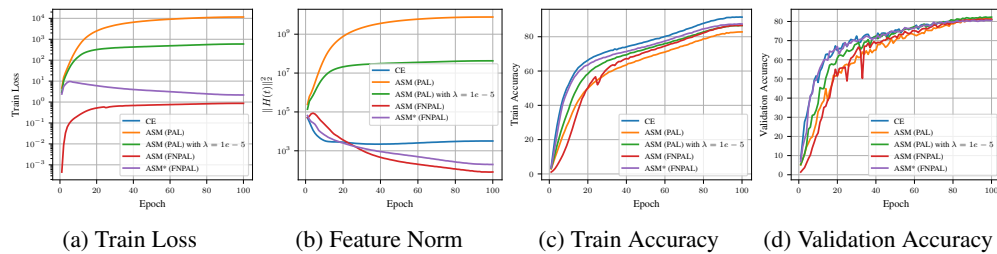


Figure 14: Behavior of visual classification on ImageNet-100 with CE, ASM (PAL) and ASM (FNPAL) under different weight decay coefficients.

where $\varepsilon \in (0, 1)$ is a hyperparameter and $\delta = - \left(1 + \frac{\hat{\mathbf{w}}_y^\top \hat{\mathbf{h}} \sqrt{1 - (1 - \varepsilon)^2}}{(1 - \varepsilon) \sqrt{1 - (\hat{\mathbf{w}}_y^\top \hat{\mathbf{h}})^2}} \right) (\mathbf{w}_y + \hat{\mathbf{h}} \hat{\mathbf{h}}^\top \mathbf{w}_y)$ (performed with a stop-gradient) satisfying $\frac{(\mathbf{w}_y + \delta)^\top \hat{\mathbf{h}}}{\|\mathbf{w}_y + \delta\|_2} = -1 + \varepsilon$.

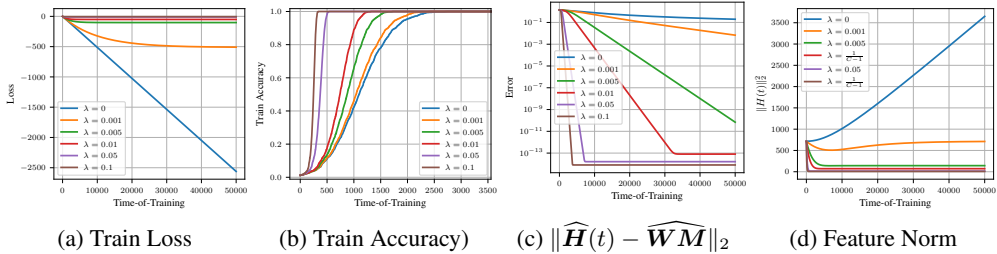


Figure 15: Behavior of gradient descent iterates of the ASM (PAL) loss in Theorem 4.1 with different explicit feature regularization ($\lambda \in \{0, 0.001, 0.005, 0.01, 0.05, 0.1\}$). We set $p = 512$, $C = 100$, $N = 10$, and $\eta = 0.1$. We randomly initialize H_0 and W , and then anchor prototypes W during training. As expected in Theorem 4.1, the error $\|\widehat{H}(t) - \widehat{WM}\|_2$ decreases as an exponential rate $O(e^{-\lambda\eta t})$, and a larger λ can accelerate the convergence.

Table 4: OOD detection performance using softmax-based (Hendrycks and Gimpel, 2016), energy-based (Liu et al., 2020), and feature norm-based approaches while model training with feature regularization ($\lambda = \{0, 1e - 6, 5e - 6, 1e - 5\}$). We use ResNet-18 and ResNet-34 to train on the in-distribution datasets CIFAR-10 and CIFAR-100, respectively. We then use SVHN (Netzer et al., 2011) as the OOD dataset to evaluate the performance of OOD detection. All values are percentages. \uparrow indicates large values are better, and \downarrow indicates smaller values are better. The best results are underlined.

Dataset \mathcal{D}_{in}^{test}	λ	FPR95 \downarrow	AUROC \uparrow	AUPR \uparrow
CIFAR-10	Softmax-based / Energy-based / Feature Norm-based			
	0	52.09 / 43.04 / 52.10	91.67 / 91.94 / 89.54	84.11 / 82.80 / 77.06
	1e-6	54.00 / 43.72 / 51.45	91.44 / 92.12 / 89.08	82.31 / 81.77 / 74.16
	5e-6	45.37 / 33.92 / 26.93	93.08 / 93.78 / 94.03	84.31 / 83.73 / 82.79
	1e-5	37.39 / 27.87 / 24.94	93.90 / 94.60 / 94.17	85.48 / 85.34 / 83.15
CIFAR-100	Softmax-based / Energy-based / Feature Norm-based			
	0	87.75 / 89.84 / 95.54	71.01 / 71.94 / 59.54	55.42 / 56.69 / 43.21
	1e-6	82.08 / 82.57 / 88.77	75.36 / 76.28 / 68.83	61.40 / 61.90 / 51.58
	5e-6	79.01 / 78.68 / 85.94	78.70 / 79.15 / 70.32	62.58 / 62.39 / 48.39
	1e-5	81.48 / 81.41 / 87.83	77.02 / 78.03 / 73.91	62.92 / 63.66 / 58.81

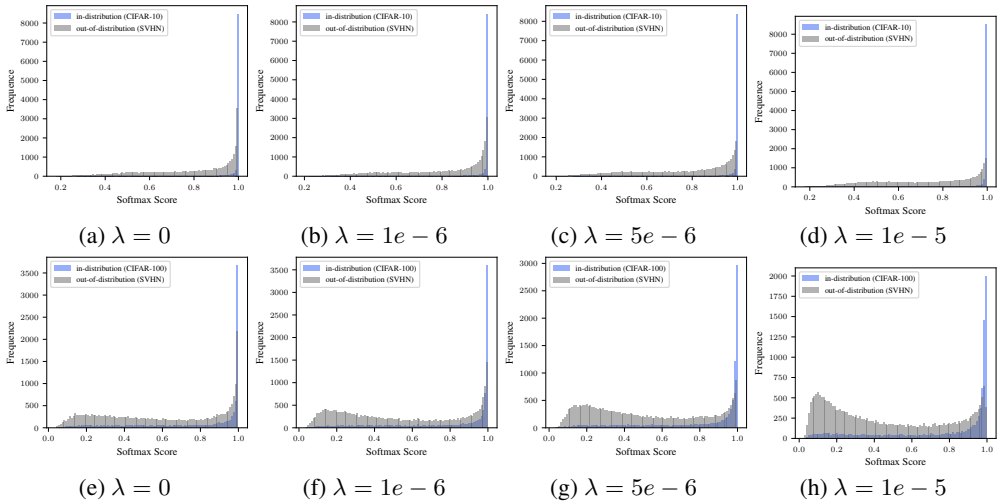


Figure 16: Distribution of softmax scores (Hendrycks and Gimpel, 2016) from models trained with different explicit feature regularization, where CE is the loss function.

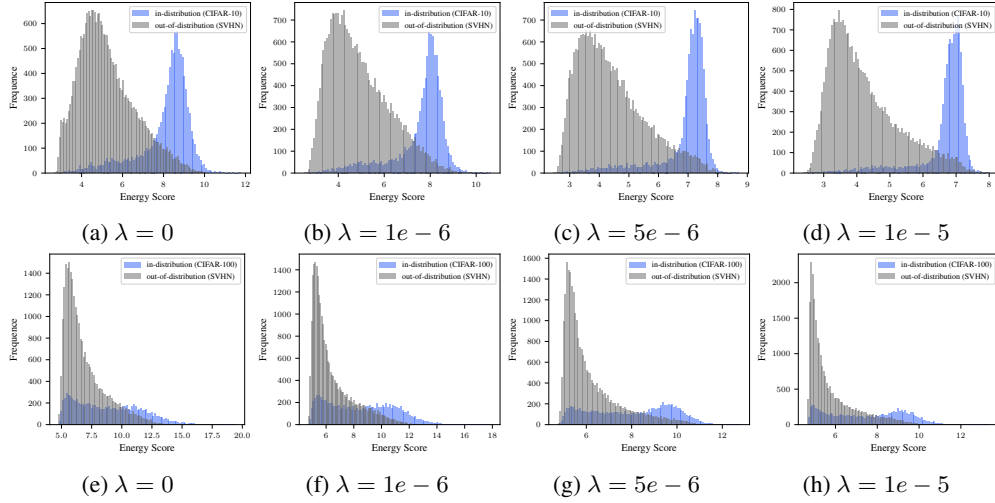


Figure 17: Distribution of energy scores (Liu et al., 2020) from models trained with different explicit feature regularization, where CE is the loss function.

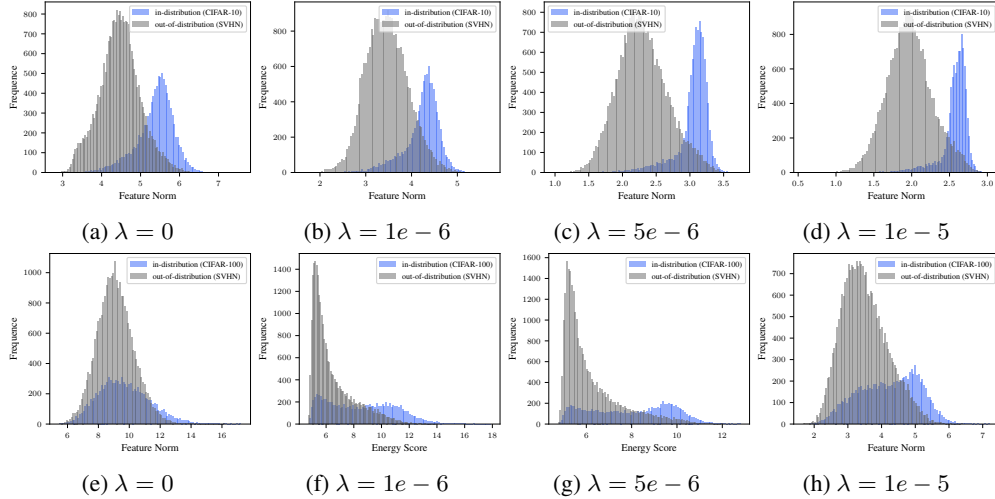


Figure 18: Distribution of feature norms from models trained with different explicit feature regularization, where CE is the loss function.

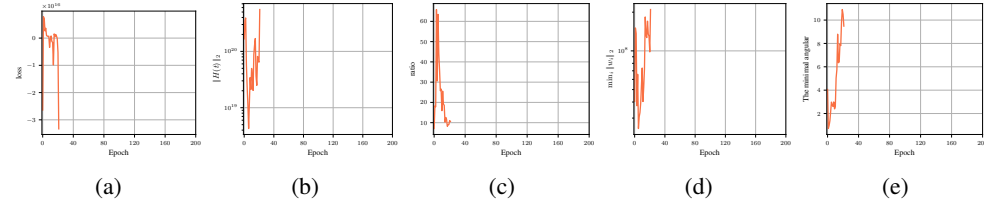


Figure 19: The behavior of features and prototypes when directly training ResNet-18 with the ASM loss 2.1 on CIFAR-10. We set the weight decay coefficient as $5e-4$. (a) The train accuracy. (b) The feature norm. (c) the ratio $\frac{\max_i \|\mathbf{w}_i\|_2}{\min_i \|\mathbf{w}_i\|_2}$. (d) $\min_i \|\mathbf{w}_i\|_2$. (e) The minimal angular between prototypes: $\arccos \max_{i \neq j} \hat{\mathbf{w}}_i^\top \hat{\mathbf{w}}_j$. In these figures, we only show the curves for the first 21 epochs, since “NaN” appears at the 22-th epoch. We can find that implicit penalization attached by other components (e.g., network architectures and weight decays) does not limit the rapid growth of the feature norm and prototype norm, indicating implicit penalization is fragile. Moreover, the ratio $\frac{\max_i \|\mathbf{w}_i\|_2}{\min_i \|\mathbf{w}_i\|_2}$ starts out very large and the minimal angular is very small, which indicates that there are two prototypes that are particular imbalanced.