

On the Optimization Landscape of Neural Collapse under MSE Loss: Global Optimality with Unconstrained Features

Jinxin Zhou^{*,#}, Xiao Li^{*,□}, Tianyu Ding[◇], Chong You[♣], Qing Qu^{†,□}, and Zihui Zhu^{†,#}

[◇]Microsoft Research, Redmond

[♣]Google Research, New York City

[#]Department of Electrical & Computer Engineering, University of Denver

[□]Department of Electrical Engineering & Computer Science, University of Michigan

March 15, 2022

Abstract

When training deep neural networks for classification tasks, an intriguing empirical phenomenon has been widely observed in the last-layer classifiers and features, where (i) the class means and the last-layer classifiers all collapse to the vertices of a Simplex Equiangular Tight Frame (ETF) up to scaling, and (ii) cross-example within-class variability of last-layer activations collapses to zero. This phenomenon is called Neural Collapse (NC), which seems to take place regardless of the choice of loss functions. In this work, we justify NC under the mean squared error (MSE) loss, where recent empirical evidence shows that it performs comparably or even better than the *de-facto* cross-entropy loss. Under a simplified unconstrained feature model, we provide the first global landscape analysis for vanilla nonconvex MSE loss and show that the (only!) global minimizers are neural collapse solutions, while all other critical points are strict saddles whose Hessian exhibit negative curvature directions. Furthermore, we justify the usage of rescaled MSE loss by probing the optimization landscape around the NC solutions, showing that the landscape can be improved by tuning the rescaling hyperparameters. Finally, our theoretical findings are experimentally verified on practical network architectures.

1 Introduction

Despite the dramatic success of modern deep neural networks (DNNs) across engineering and sciences [1–4] that we have witnessed in the past decade, the practice of deep learning has yet been shrouded with mysteries, ranging from the design of appropriate network architectures [5, 6] to the generalization and robustness properties [7–9] of the learned networks. For instance, even the right choice of training loss function has not been thoroughly justified. For classification problems, although the cross entropy (CE) loss is the standard choice for network training, recent work [10] demonstrated with extensive experiments that DNNs trained with mean-squared error (MSE) loss achieve on par or even better performance compared to those of the CE loss.

Towards demystifying DNN, a recent interesting line of work [11–18] studied and characterized the learned deep representations during the terminal phase of training, where several intriguing

*The first two authors contributed to this work equally.

†The last two authors share the corresponding authorship of this work.

phenomena have been discovered. In particular, recent seminal work of [11, 12] empirically demonstrated that last-layer features and classifiers of a trained DNN exhibit the following *Neural collapse* (\mathcal{NC}) property:

- ($\mathcal{NC}1$) **Variability collapse:** the individual features of each class concentrate to their class-means.
- ($\mathcal{NC}2$) **Convergence to simplex ETF:** the class-means have the same length and are maximally distant; they form a Simplex Equiangular Tight Frame (ETF).
- ($\mathcal{NC}3$) **Convergence to self-duality:** the last-layer linear classifiers perfectly match their class-means.
- ($\mathcal{NC}4$) **Simple decision rule:** the last-layer classifier is equivalent to a Nearest Class-Center decision rule.

It has been empirically demonstrated that the \mathcal{NC} persists across the range of canonical classification problems with the CE loss. These results imply that deep networks are essentially learning maximally separable features between classes, and a max-margin classifier in the last layer upon these learned features, touching the ceiling in terms of the training performance. Later work theoretically investigated the \mathcal{NC} based on a simplified assumption of the so-called *unconstrained feature model* [15] or *layer-peeled model* [14], where the features are viewed as free optimization variables. The underlying reasoning is that modern deep networks are often highly overparameterized with the capacity of learning any representations [19–22], so that the last-layer features can approximate, or interpolate, any point in the feature space. Under the unconstrained feature model, the work [14–16, 23–26] showed that the \mathcal{NC} solutions are the only global optimal solution for nonconvex training losses under different settings. However, given the nonconvexity of the problem, even under the unconstrained feature model these global optimality results do guarantee that the \mathcal{NC} solutions can be efficiently achieved. This has been further resolved by the recent work [18], showing that the CE loss function enjoys a benign global optimization landscape under the unconstrained feature model. It shows that every saddle point is a strict saddle with negative curvature, so that the CE loss can be efficiently optimized to the \mathcal{NC} solution regardless of the nonconvexity.

It should be noted that the \mathcal{NC} phenomenon is not *solely* pertinent to the particular choice of the CE loss. It has been recently reported [12], that DNNs trained with the MSE loss also exhibit very similar \mathcal{NC} phenomena but with even *faster* collapse in terms of training epochs and with better (adversarial) robustness. In the meanwhile, the MSE loss is not only appealing for its algebraic simplicity, but it also demonstrates on-par or even better generalization performances compared to the CE loss, as reported by recent line of work [10]. However, the theoretical study of MSE loss for \mathcal{NC} is still limited [12, 15, 26]. Under the unconstrained feature model, their work proved that the continuous gradient flow of the MSE loss converges to \mathcal{NC} solutions. In particular, the work [15] relies on linearizations of the ordinary differential equation by assuming very small initializations, which is not well aligned with the practice of deep learning where the weights are usually initialized with non-negligible magnitudes such as by the Kaiming initialization [27]. Because the choice of the loss function without balanced weight decay, the analysis in [12] only focuses on the renormalized features and studies the continually renormalized gradient flow.¹ Moreover, in practice deep networks are usually trained using iterative algorithms such as stochastic gradient descent (SGD) with nontrivial stepsizes, rather than using the continuous gradient flows.

¹The model used in [12] imposes a weight decay on the classifier, but not on the features. Thus, without renormalization, the weights of the classifier will converge to zero while the features will blow up.

The work [28–31] study deep homogeneous classification networks (without bias terms but beyond the unconstrained features model) trained with MSE loss, stochastic gradient descent, and weight decay. In particular, the solutions satisfying the so-called symmetric quasi-interpolation assumption are proved to obey \mathcal{NC} properties, but the properties of other solutions are not investigated [30, 31]

As far as we know, the work closest to ours is the concurrent work [26]. Under similar unconstrained feature models, the work studies the global optimality condition of \mathcal{NC} for the MSE loss for both two-layer and three layer networks, but *not* the global optimization landscape. Additionally, it studies special cases of the MSE loss with either no bias term, or no weight decay on the bias term. In comparison, our work not only study the MSE loss under more general setting with bias term included, but also shows the strict saddle property of the benign nonconvex landscape.

Contributions. In this work, we provide a thorough analysis of neural network by examining its last-layer features. In particular, we work under the unconstrained feature model to characterize the *global* optimization landscape of over-parameterized neural networks trained with the MSE loss. Our contributions can be highlighted as follows.

- **Characterization of global solutions.** We provide a mathematical characterization of *all* the global solutions for the last layer features and classifier, showing that they satisfy the \mathcal{NC} properties with certain choices of regularization parameters. This is in contrast to previous work [12, 15] which only characterize the solutions that are produced by a particular optimization algorithm (i.e., gradient flow). Moreover, these work only consider cases that the feature dimension is larger than the number of classes, while our analysis covers all choices of feature dimension.
- **Benign global landscape.** We prove that the loss function is a *strict saddle function* [32–34], where every critical point is either a global solution or a *strict saddle point* with negative curvature. This implies that there is *no* spurious local minimizer on the optimization landscape. Hence, our work is distinguished from previous work [14–17, 23, 24, 26] that only characterizes global minimizers. The benign global landscape implies that any method that can escape strict saddle points (e.g. stochastic gradient descent) converges to a global solution that exhibits \mathcal{NC} (see Section 4).
- **Understanding the *rescaled* MSE.** In practice, rescaling the MSE loss (see Section 2.2) is empirically demonstrated to be critical for obtaining competitive performance compared to the CE loss particularly when the number of classes is large [10, 35]. We show empirically that the \mathcal{NC} exhibits for rescaled MSE as well. To understand the benefit of the rescaling, we provide a visualization of the optimization landscape w.r.t. unconstrained features, showing that rescaling aligns the gradient direction to be perpendicular to the decision boundary between classes hence may facilitate the convergence of gradient based algorithms to more discriminative features.

Compared to the recent global landscape analysis for the CE loss [18], our result implies that both losses learn similar \mathcal{NC} features and classifiers when $d \geq K$. Hence, from the \mathcal{NC} perspective, this work provides a theoretical explanation for the observations in [10] that the DNN trained by the MSE loss achieves on par performance compared to that trained with the CE loss. Additionally, it should be noted that there are several major differences between our result and [18]. First, the work of [18] only studied the setting where the feature dimension d is larger than the number of classes K , while we characterized the global optimality for both the cases of $d < K$ and $d \geq K$. We observe dramatically different performance for DNN learned by CE and MSE when $d < K$. Second, for the MSE loss, we showed that the bias term plays an important role² for the solution

²For the MSE loss, when there is no bias term, the features (and classifier) that minimize the loss function form orthonormal matrices instead of Simplex ETFs when $d \geq K$.

to be \mathcal{NC} , while for CE loss the \mathcal{NC} solution can be achieved without bias terms.

2 The Problem Setup

The goal of deep learning is to learn a multi-layer nonlinear mapping $\psi(\cdot) : \mathbb{R}^m \mapsto \mathbb{R}^K$, that is able to fit the training data and generalize. More precisely, a deep neural network classifier can be generally written as

$$\psi_{\Theta}(\mathbf{x}) = \mathbf{W}\phi_{\theta}(\mathbf{x}) + \mathbf{b}, \quad (1)$$

where $\phi_{\theta}(\cdot) : \mathbb{R}^m \mapsto \mathbb{R}^d$ is the feature mapping, on top of which is the linear classifier (\mathbf{W}, \mathbf{b}) . $\phi_{\theta}(\mathbf{x})$ is usually referred to as the *representation* or *feature* of the input \mathbf{x} learned from the network. For convenience, we use θ to denote the network parameters in the feature mapping, and $\Theta = \{\theta, \mathbf{W}, \mathbf{b}\}$ to denote *all* the network parameters. In this way, the function implemented by a neural network classifier can also be expressed as a linear classifier acting upon $\phi_{\theta}(\mathbf{x})$.

In this work, we focus on learning deep networks for multi-class classification tasks (say, with K classes), where the class label of a sample $\mathbf{x}_{k,i}$ in the k -th class is given by a one-hot vector $\mathbf{y}_k \in \mathbb{R}^K$ with only the k th entry equal to unity ($1 \leq k \leq K$). Throughout the paper, we study the setting where the number of training samples in each class is balanced, i.e., each class has n training samples. Let $N = Kn$. During the training phase, the task is then to learn the parameters Θ so that the output of the model on an input sample $\mathbf{x}_{k,i}$ approximates the corresponding output \mathbf{y} (i.e. $\psi_{\Theta}(\mathbf{x}_{k,i}) \approx \mathbf{y}_k$). To quantify this approximation, it can be done by optimizing a simple MSE loss as follows

$$\min_{\Theta} \frac{1}{2N} \sum_{k=1}^K \sum_{i=1}^n \|\psi_{\Theta}(\mathbf{x}_{k,i}) - \mathbf{y}_k\|_2^2 + \frac{\lambda}{2} \|\Theta\|_F^2, \quad (2)$$

where $\lambda > 0$ is the regularization parameter (a.k.a., the weight decay parameter).

2.1 Basic Problem Formulation Based on Unconstrained Feature Models

Analyzing deep networks is a tremendously difficult task mainly due to the nonlinear interactions between a large number of layers. Nonetheless, as argued by a line of recent work [19–22] that modern deep networks are often highly overparameterized to approximate any continuous function, it motivates us to simplify the analysis by treating the last-layer features as *free* optimization variables $\mathbf{h}_{k,i} = \phi_{\theta}(\mathbf{x}_{k,i}) \in \mathbb{R}^d$. Such a simplification is called *unconstrained feature model* [15] (or *layer-peeled model* in [14]), which simplifies the study of the last-layer representations of the network. To simplify the notation, let us denote

$$\begin{aligned} \mathbf{W} &:= [\mathbf{w}^1 \quad \mathbf{w}^2 \quad \cdots \quad \mathbf{w}^K]^\top \in \mathbb{R}^{K \times d}, \\ \mathbf{H} &:= [\mathbf{H}_1 \quad \mathbf{H}_2 \quad \cdots \quad \mathbf{H}_K] \in \mathbb{R}^{d \times N}, \text{ and} \\ \mathbf{Y} &:= [\mathbf{Y}_1 \quad \mathbf{Y}_2 \quad \cdots \quad \mathbf{Y}_K] \in \mathbb{R}^{K \times N}, \end{aligned}$$

where \mathbf{w}^k is a row vector of \mathbf{W} , $\mathbf{H}_k := [\mathbf{h}_{k,1} \quad \cdots \quad \mathbf{h}_{k,n}] \in \mathbb{R}^{d \times n}$ contains all the k -th class features, and $\mathbf{Y}_k := [\mathbf{y}_k \quad \cdots \quad \mathbf{y}_k] \in \mathbb{R}^{K \times n}$ for all $k = 1, 2, \dots, K$. Based on the unconstrained feature model, we consider a slight variant of (2), given by

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} f(\mathbf{W}, \mathbf{H}, \mathbf{b}) := \left\{ \frac{1}{2N} \left\| \mathbf{W}\mathbf{H} + \mathbf{b}\mathbf{1}_N^\top - \mathbf{Y} \right\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2 \right\}, \quad (3)$$

where $\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}, \lambda_{\mathbf{b}} > 0$ are the penalties for \mathbf{W} , \mathbf{H} , and \mathbf{b} , respectively.

Here, because we treat the last-layer feature \mathbf{H} as a free optimization variable, we put the weight decay on \mathbf{W} and \mathbf{H} , which is different from the practice that the weight decay is enforced on all the network parameters Θ as shown in (2). Nonetheless, as discussed in [18], this idealization is reasonable since the energy of the features (i.e., $\|\mathbf{H}\|_F$) can indeed be upper bounded by the energy of the weights at every layer if the inputs are bounded (which holds in practice), implying that the norm of \mathbf{H} is *implicitly* penalized by penalizing the norm of Θ . Additionally, for the CE loss, the experiments in [18] show on-par performance for the two types of weight decay. Thus, we expect similar performances for the MSE loss.

On the other hand, the experiments in [16, 18] conducted on random labels imply that the strong assumption of unconstrained feature model is reasonable for explaining \mathcal{NC} during the training phase: when the network (1) is highly overparameterized, the learned network in practice will fit to the random labels and neural collapse, regardless of the input. Moreover, as we shall see in the following sections, both theory and experiments demonstrate that such simplification preserves the core properties of last-layer classifiers and features—the \mathcal{NC} phenomenon.

2.2 Rescaled MSE Loss under Unconstrained Features

On the other hand, it should be noted that, when training with the vanilla formulation of the MSE loss (2), empirically good performances are reported *only* when the number of classes is small (e.g., CIFAR10 [36] with $K < 100$). When training for a large number of classes such as ImageNet [37], to achieve better performance *rescaling* is often needed [10, 35]. Intuitively, the basic idea is to rescale the MSE loss (3) by a pair of positive scalars (α, M) ,

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} \frac{1}{2N} \left\| \Omega_\alpha^{\odot 1/2} \odot (\mathbf{W}\mathbf{H} + \mathbf{b}\mathbf{1}^\top - M\mathbf{Y}) \right\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2, \quad (4)$$

so that we can put more emphasize on training the correct class. Here, \odot denotes the entry-wise Hadamard product, $\Omega_\alpha^{\odot 1/2}$ means taking square root for each element, and

$$\Omega_\alpha = [\omega_1 \mathbf{1}_n^\top \quad \cdots \quad \omega_K \mathbf{1}_n^\top], \quad \text{with } \omega_k(\alpha) \in \mathbb{R}^K \text{ and } \omega_{ki}(\alpha) = \begin{cases} \alpha, & i = k, \\ 1, & \text{otherwise.} \end{cases}$$

In comparison to [12, 15, 26], our work not only studies \mathcal{NC} under the vanilla setting (3) but also investigates the more practical rescaled version of the MSE loss (4). In particular, in Section 3.3, we provide geometric intuitions on why rescaling would be a better choice for loss design. We will corroborate our reasoning via experiments on practical network training in Section 4.

3 Main Theoretical Results

In this section, we present our study on global optimality conditions as well as geometric properties of the nonconvex (rescaled) MSE loss under the unconstrained feature model.

3.1 Global Optimality Conditions

First, we study the nonconvex MSE loss (3) by characterizing its global solutions under different settings of the feature and class dimensions. We show that the only global solutions of (3) are neural collapsing, satisfying the \mathcal{NC} properties introduced at the beginning of Section 1.

Theorem 3.1 (Global Optimality Conditions) *Assume that the number of training samples in each class is balanced, $n = n_1 = \cdots = n_K$, and let $(\mathbf{W}^*, \mathbf{H}^*, \mathbf{b}^*)$ be a global minimizer of the vanilla MSE loss in (3). Let $\bar{\mathbf{H}}^* = [\bar{\mathbf{h}}_1^* \quad \cdots \quad \bar{\mathbf{h}}_K^*]$, with $\bar{\mathbf{h}}_k^*$ being the mean of the k -th class features. Then, $(\mathbf{W}^*, \mathbf{H}^*, \mathbf{b}^*)$ satisfies the following properties:*

- If $\lambda_{\mathbf{W}}\lambda_{\mathbf{H}} < \frac{1}{NK}$, then $(\mathbf{W}^*, \mathbf{H}^*)$ satisfies $\mathcal{NC}1$ and $\mathcal{NC}3$ as

$$\mathbf{h}_{k,i}^* = \bar{\mathbf{h}}_k^*, \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}n}} \mathbf{w}^{*k} = \bar{\mathbf{h}}_k^*, \quad \forall k \in [K], i \in [n].$$

Otherwise, if $\lambda_{\mathbf{W}}\lambda_{\mathbf{H}} \geq \frac{1}{NK}$, then $\mathbf{W}^* = \mathbf{0}, \mathbf{H}^* = \mathbf{0}$.

- If $\lambda_{\mathbf{W}}\lambda_{\mathbf{H}} < \frac{1}{NK}$, then $\bar{\mathbf{H}}^*$ further obeys the following properties ($\mathcal{NC}2$) for different d :
 1. If $d < K - 1$: we have $\bar{\mathbf{H}}^{*\top} \bar{\mathbf{H}}^* = C_1 \mathcal{P}_d(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top)$, where $\mathcal{P}_d(\mathbf{M})$ denotes the best rank- d approximating of \mathbf{M} ;
 2. If $d = K - 1$: we have $\bar{\mathbf{H}}^{*\top} \bar{\mathbf{H}}^* = C_2(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top)$;
 3. If $d \geq K$: we have $\bar{\mathbf{H}}^{*\top} \bar{\mathbf{H}}^* =$

$$\begin{cases} C_3(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top), & \text{if } \lambda_{\mathbf{b}} \leq \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}} \\ C_4(\mathbf{I} - \frac{\sqrt{n\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{\lambda_{\mathbf{b}}(1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}})} \mathbf{1}_K \mathbf{1}_K^\top), & \text{otherwise} \end{cases} \quad (5)$$

where $\frac{\sqrt{n\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{\lambda_{\mathbf{b}}(1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}})} \leq \frac{1}{K}$ in the second case since $\lambda_{\mathbf{b}} \geq \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}$.

Here, C_1, C_2, C_3 , and C_4 are some positive numerical constants that depend on $\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}, \lambda_{\mathbf{b}}$.

- The bias satisfies $\mathbf{b}^* = b^* \mathbf{1}_K$ with $b^* \leq \frac{1}{K}$ given by:

1. If $d < K$: we have $b^* = \frac{1}{K(\lambda_{\mathbf{b}}+1)}$;
2. Otherwise, $b^* = \begin{cases} \frac{1}{K(\lambda_{\mathbf{b}}+1)}, & \lambda_{\mathbf{b}} \leq \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}} \\ \frac{\sqrt{n\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{\lambda_{\mathbf{b}}}, & \text{otherwise.} \end{cases}$

In particular, when $\lambda_{\mathbf{b}} \rightarrow 0$, we have $b^* \rightarrow \frac{1}{K}$; when $\lambda_{\mathbf{b}} \rightarrow \infty$, we have $b^* \rightarrow 0$.

We postpone the detailed proof to Appendix B. In the following, we discuss the implications of Theorem 3.1 in detail.

- **Implications on the choice of the feature dimension d .** As we observe from Theorem 3.1, for the MSE loss (3), any global solution always exhibits variability collapse ($\mathcal{NC}1$) and self-duality ($\mathcal{NC}3$). However, the convergence of class means to simplex ETF ($\mathcal{NC}2$) critically depends on the feature dimension d . When $d \geq K - 1$, for proper choices of $\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}$, and $\lambda_{\mathbf{b}}$, the global configuration of the class mean $\bar{\mathbf{H}}^*$ is always a simplex ETF. In particular, when $d = K - 1$, the simplex ETF configuration even does not depend on $\lambda_{\mathbf{b}}$. On the other hand, if $d < K - 1$, our theory implies that the global solution for $\bar{\mathbf{H}}^*$ is only the best rank- d approximation of the simplex ETF, where the class-means of the each class are neither having equal length nor being maximally pairwise-distanced. This result is consistent with the fact that K vectors in \mathbb{R}^d cannot form a K -Simplex ETF if $K > d - 1$, and supports the practice of learning overparameterized network for choosing $d \geq K$.³
- **Comparison to the CE loss.** For the CE loss under the unconstrained feature model, when $d \geq K$ recent work [18] showed that any global solution satisfies all three \mathcal{NC} properties regardless of choices of the weight decay parameters (i.e., $\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}$, and $\lambda_{\mathbf{b}}$). Moreover, the bias term there becomes zero. In contrast, Theorem 3.1 shows that the solution with the MSE loss is dependent

³For example, the dimension of the features of a ResNet [38] is typically set to $d = 512$ for CIFAR10 [36], a dataset with $K = 10$ classes. This dimension grows to $d = 2048$ for ImageNet [37], a dataset with $K = 1000$ classes.

upon choice of regularization parameters $\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}, \lambda_{\mathbf{b}}$ and that the class mean $\overline{\mathbf{H}}^*$ may not be a simplex ETF. Moreover, the bias term is essential to achieve simplex ETF solutions for MSE loss. Without the bias term (i.e., $\lambda_{\mathbf{b}} \rightarrow \infty$), (5) implies that the class mean $\overline{\mathbf{H}}^*$ becomes an orthonormal matrix even when $d \geq K$. Thus, the analysis of global optimality conditions for the MSE loss is more complicated than for the CE loss⁴.

- **Comparison to previous work [12, 15].** As discussed in Section 1, the previous work [12, 15] only characterize the solutions to (3) that are produced by a particular optimization algorithm (i.e., gradient flow) and under specific cases such as $\lambda_{\mathbf{b}} \rightarrow 0$ and the feature dimension is larger than the number of classes. In contrast, we characterize the global optimality conditions for the MSE loss (3) and our analysis covers all choices of feature dimension and weight decay parameters.
- **Extension to the rescaled MSE.** Although our current analysis is only for the vanilla MSE loss (3), we expect that similar global optimality results should also hold for the rescaled version (4). This has been corroborated by our experimental results in Section 4. Notice that if we fix $\alpha = 1$ in (4), the analysis only with large M is simple and remain the same as Theorem 3.1. However, dealing with both α and M requires extra technicalities, that we leave for future work.

3.2 Characterizations of The Benign Global Landscape

Theorem 3.1 implies that the (only!) global minimizers to (3) are those satisfying \mathcal{NC} properties. However, the MSE loss function is *nonconvex*, hence it is not obvious whether the benign global solutions can be *efficiently* achieved even under the unconstrained feature model. To deal with this challenge, in the following we further investigate the global optimization landscape of (3). By leveraging recent advances on nonconvex optimization [32–34, 39–42], we first show that our nonconvex MSE loss (3) without bias term is a *strict saddle* function that every non-global critical point is a saddle point with negative curvature (i.e., its Hessian has at least one negative eigenvalue).

Theorem 3.2 (Benign landscape for MSE without bias term) *The following MSE loss without bias term*

$$\frac{1}{2N} \|\mathbf{W}\mathbf{H} - \mathbf{Y}\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2$$

is a strict saddle function with no spurious local minimum. That is, any of its critical point is either a global minimizer, or it is a strict saddle point whose Hessian has a strictly negative eigenvalue.

We postpone the proof to Appendix B (see Theorem B.2). By viewing \mathbf{W} and \mathbf{H} as two factors of a matrix $\mathbf{Z} = \mathbf{W}\mathbf{H}$, the formulation in (3) is closely related to nonconvex low-rank matrix problems [43–49] with the Burer-Moneirto factorization approach [50]. In particular, the work [47, 51] studied a similar problem with $\lambda_{\mathbf{W}} = \lambda_{\mathbf{H}}$, but only for particular choices of d : d is either required to be exactly the rank of the solution of the corresponding convex problem [47], or relatively large in [51]. In contrast, our Theorem 3.2 characterizes the benign landscape for all choices of feature dimension.

The following result establishes global optimization landscape of the MSE loss (3).

Theorem 3.3 (Benign landscape for MSE loss (3)) *Assume that the feature dimension d is larger than the number of classes K . The nonconvex MSE loss function $f(\mathbf{W}, \mathbf{H}, \mathbf{b})$ in (3) is a strict saddle function.*

⁴The proof of Theorem 3.1 is also dramatically different to the one for CE loss in [18]: the latter mainly shows that \mathcal{NC} solutions have small objective value than others since \mathcal{NC} solutions are the only global minimizers, while the proof of Theorem 3.1 directly analyzes the global minimizers for different scenarios.

This result is similar to that of [18, Theorem 3.2], which showed that the unconstrained feature model with CE loss is a strict saddle function. The high level proof idea for [18] is to construct the negative curvature direction for saddle points in the null space of $\mathbf{W} \in \mathbb{R}^{K \times d}$. Because the proof in [18] actually holds more generally for any smooth convex loss function with weight decay, the same technique also offers a proof for Theorem 3.3 (and potentially can extend Theorem 3.3 for the rescaled MSE in (4)). Here, it should be noted that we make the assumption $d > K$ so that the null space of $\mathbf{W} \in \mathbb{R}^{K \times d}$ always exists. However, we believe the strict saddle property holds for any d and leave it as future work.

As a consequence, if \mathbf{H} is a free optimization variable, this implies that the global solutions of the training problem in (3) can be efficiently found by many first-order and second-order optimization methods [52]. In particular, (stochastic) gradient descent with random initialization is guaranteed [32, 53] to almost surely find a global minimizer for strict saddle functions with no spurious local minima, which is the case for our problem (3). In comparison, existing results on MSE loss [12, 15] only studied the trajectory of gradient flows (3) on either the linear terms [15] or the central path component [12], which is insufficient to explain/guarantee efficient, global convergence of iterative optimization algorithms.

3.3 Delving Deeper into Optimization Landscapes: Why Rescaling Helps?

While our global landscape analysis for the vanilla MSE loss (3) in Section 3.2 implies that a gradient based algorithm converges to global \mathcal{NC} solutions *asymptotically* [53], it did not characterize the rate of convergence – in other words, how fast an optimization method converges. Often around the global solutions (i.e., the simplex ETF), we expect that the landscape has certain regularity condition which measures how well-aligned between the negative gradient direction and the direction towards the global solution. Thus, the regularity conditions in turn will characterize how fast a gradient based method converges. For better understanding the regularity properties and algorithmic convergences, we use visualization techniques to visualize the optimization landscape of MSE losses around the global ETFs solutions. In particular, our visualization sheds light on (i) why training with vanilla MSE loss performs worse than that of the CE loss, and (ii) how the rescaling techniques in Section 2.2 improves the performance of the MSE loss.

Even under the unconstrained feature model, visualization of the MSE loss landscape could still be difficult, which is due to the fact that the variables \mathbf{H} , \mathbf{W} , and \mathbf{b} are all high-dimensional. Here, we further simplify the problem by assuming $\mathbf{b} = \mathbf{0}$ and that \mathbf{W} is at the global optimum and forms a simplex ETF. Thus, we can examine the landscape only with respect to (w.r.t.) the feature vectors $\mathbf{h}_{k,i}$ for the k th class. Although $\mathbf{h}_{k,i} \in \mathbb{R}^d$ is still high-dimensional for large d , we plot the optimization landscape by restricting $\mathbf{h}_{k,i}$ to a 2D plane spanned by $\{\mathbf{w}^k, \mathbf{w}^{k'}\}$, where \mathbf{w}^k is the classifier for the k th class and $k' \neq k$ can be chosen arbitrarily because the simplex ETF is invariant to rotations. Finally, we visualize the landscape using the polar coordinates, where the s -axis denotes the ℓ_2 norm of $\mathbf{h}_{k,i}$ and the θ -axis denotes the angle between $\mathbf{h}_{k,i}$ and \mathbf{w}^k (see Figure 1 for an illustration). The predicted membership for $\mathbf{h}_{k,i}$ is determined by θ and is invariant to s . Hence, larger gradient along the θ direction may help with learning more discriminative features. See Appendix C for a formal explanation. This design choice allows us to examine the gradient in

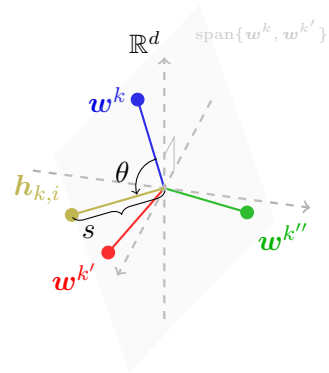


Figure 1: **An illustration of the visualization method.**

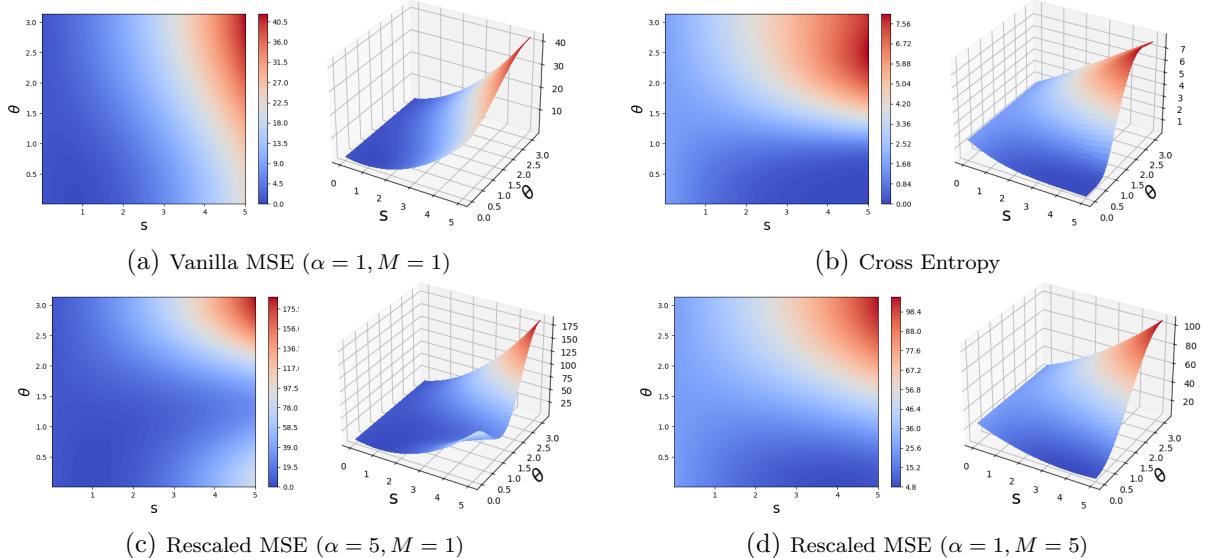


Figure 2: **Visualization of optimization landscape with different losses.** We fix \mathbf{W} as a simplex ETF and illustrate the landscape only w.r.t. a feature $\mathbf{h}_{k,i}$. For each plot, the s -axis denotes $\|\mathbf{h}_{k,i}\|_2$, and the θ -axis denotes the angle $\arccos(\langle \mathbf{h}_{k,i}, \mathbf{w}^k \rangle)$.

directions co-linear to (i.e., with varying s) and perpendicular to (i.e., with varying θ) the decision boundary separately.

In Figure 2, the visualizations of landscapes of different loss functions are provided. As we observe from Figure 2a, the landscape of vanilla MSE loss is steep w.r.t. s while it is flat w.r.t. θ . Because the size of θ determines the closeness to the right class, this implies that optimizing the vanilla MSE loss will take a longer time to converge to a desired solution with $\theta \approx 0$. In contrast, the landscape of CE loss in Figure 2b is steeper w.r.t. θ than w.r.t. s in a large region where $s > 1$ and $\theta < 1.5$. This difference of the landscapes around the global solutions potentially explains why CE is a preferred choice than the vanilla MSE, given that the features $\mathbf{h}_{k,i}$ would converge faster to the simplex ETF solutions via optimizing the CE loss. Nonetheless, the issue with the vanilla MSE can be mitigated via the rescaling approach that we discussed in Section 2.2. As shown in Figures 2c and 2d, the rescaled MSE loss (4) (with large M , in particular), leads to a “better” optimization landscape similar to that of the CE loss. Therefore, through studying the \mathcal{NC} and corresponding optimization landscapes, our work provides intuitive explanations on (i) the incompetence of the vanilla MSE loss (3), and (ii) the effectiveness of rescaling (4) for classification tasks.

4 Experiments

In this section, we conduct experiments to validate our findings from Section 3 on practical networks and standard datasets. We first introduce new metrics to better evaluate how well the \mathcal{NC} properties are satisfied in practical neural networks, in addition to the ones used in [11, 18]. Second, we verify our theoretical results in Section 3.1 by showing that the \mathcal{NC} phenomena are algorithmic independent. Third, by a similar experiment as in [18], we show that we could fix the last layer weights as a Simplex ETF while achieving comparable generalization performances as explicitly training the classifier. Finally, we examine our findings in Section 3.3 that the rescaling factors in the rescaled MSE loss is beneficial for forming benign optimization landscapes. For the details of

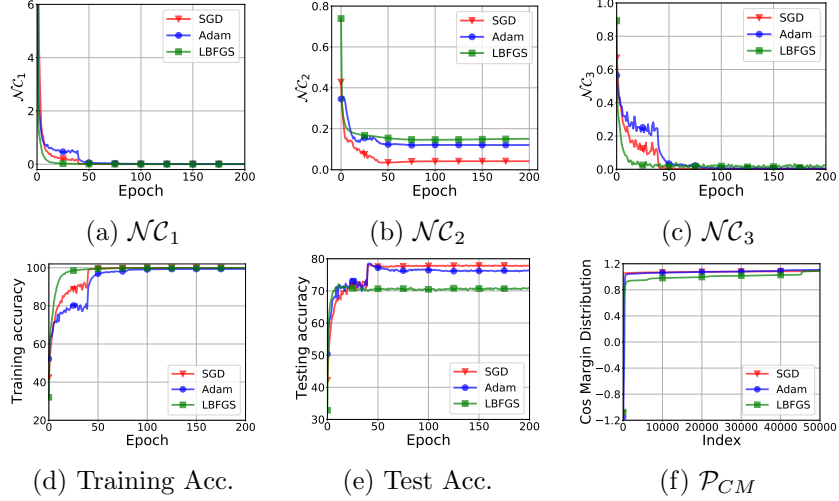


Figure 3: **Illustration of $\mathcal{N}\mathcal{C}$, training and test accuracy and cosine margin across different training algorithms with ResNet18 on CIFAR10.** The networks are trained without data augmentation.

the experimental setup, we refer readers to the Appendix A.

New metrics for evaluating $\mathcal{N}\mathcal{C}$. To evaluate the $\mathcal{N}\mathcal{C}$ properties of well-trained neural networks, we adopt the same $\mathcal{N}\mathcal{C}_1$, $\mathcal{N}\mathcal{C}_2$ and $\mathcal{N}\mathcal{C}_3$ metrics as [11, 18], which measure the within-class variability of \mathbf{H} , the convergence of \mathbf{W} to a simplex ETF, and the self-duality between \mathbf{H} and \mathbf{W} ; see Appendix A for the details.⁵ To better measure $\mathcal{N}\mathcal{C}$, this paper also introduces the following two metrics that measure the diversities and margins of the learned features:

- **Numerical rank.** The $\mathcal{N}\mathcal{C}_1$ metric measures the variability collapse through the between-class and within-class covariance matrices, which does not directly reveal the dimensionality of the features spanned for each class. Ideally, when $\mathcal{N}\mathcal{C}$ happens, for each class the feature dimension should collapse to one. To measure the dimensionality, we introduce a new metric that we call it *numerical rank*, denoted by $\widetilde{\text{rank}}(\mathbf{H}) := \frac{1}{K} \sum_{k=1}^K \frac{\|\mathbf{H}_k\|_*^2}{\|\mathbf{H}_k\|_F^2}$. Here, $\|\cdot\|_*$ represents the nuclear norm [54] (i.e., the sum of singular values), while the Frobenius norm $\|\cdot\|_F$ in the denominator serves as a normalization factor. The metric is evaluated by averaging over all the classes. Our metric is inspired by the numeral sparsity (defined as $\|\mathbf{a}\|_1^2 / \|\mathbf{a}\|_2^2$ for $\mathbf{a} \in \mathbb{R}^n$) that serves as a stable measure for sparsity of vectors [55]. For our numerical rank, we expect that the smaller $\widetilde{\text{rank}}(\mathbf{H})$ is, the more collapsed the features are to their class means.
- **Cosine margin.** All current metrics measure $\mathcal{N}\mathcal{C}$ from a panoramic view, and do not quantify the behavior of individual features. We introduce a metric based on the cosine margin of individual features. From the explanation in Section 3.3, neural network determines the class member by the direction of features rather than its length. Thus, we define the cosine margin for each sample as $CM_{k,i} = \cos \theta_{k,i;k} - \max_{j \neq k} \cos \theta_{k,i;j}$, where $\cos \theta_{k,i;j} = \frac{\langle \mathbf{w}^j - \mathbf{w}_G, \mathbf{h}_{k,i} - \mathbf{h}_G \rangle}{\|\mathbf{w}^j - \mathbf{w}_G\|_2 \|\mathbf{h}_{k,i} - \mathbf{h}_G\|_2}$ represents the cosine of the angle between the feature $\mathbf{h}_{k,i}$ and the j -th classifier \mathbf{w}^j , \mathbf{h}_G denotes the global mean of all the features, and \mathbf{w}_G denotes the mean of all the rows in \mathbf{W} . Recall that $\mathbf{h}_{k,i} \in \mathbb{R}^d$

⁵We also refer the reader to [18] for the exact definitions of these quantities. Note that for the case $d < K$, the definition of $\mathcal{N}\mathcal{C}_2$ and $\mathcal{N}\mathcal{C}_3$ will be slightly different from those in [18] based on our theoretical results in Section 3.1.

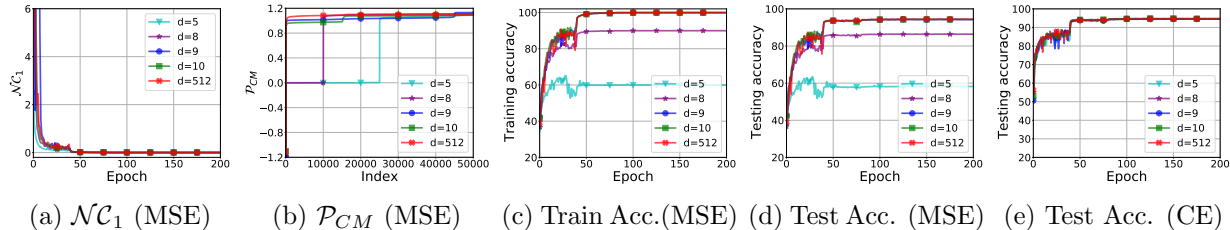


Figure 4: **Comparison of the performances on networks with different feature dimensions d for MSE and CE losses.** We compare within-class variation collapse \mathcal{NC}_1 , cosine margin distribution \mathcal{P}_{CM} , training accuracy, and test accuracy on learned classifier with different feature dimension d on CIFAR10 using ResNet18 with data augmentation. The network is trained by the SGD optimizer.

denotes the feature of i -th sample in the k -th class and $\mathbf{w}^j \in \mathbb{R}^d$ denotes the j -th row of the linear classifier weight $\mathbf{W} \in \mathbb{R}^{K \times d}$. We sort the cosine margins over the training dataset in the ascending order and denote the resulted distribution as \mathcal{P}_{CM} . We note that a similar metric has been explored by the work [56] as an alternative for the probability margin.⁶

The prevalence of \mathcal{NC} across different optimization algorithms. The benign landscape for optimization of neural networks with vanilla MSE loss suggests the existence of \mathcal{NC} regardless of specific choice of the optimizer. We validate this result by training ResNet18 on CIFAR10 with vanilla MSE loss, using three different optimization algorithms: SGD, Adam and L-BFGS. As shown in Figure 3, $\mathcal{NC}_1, \mathcal{NC}_2$, and \mathcal{NC}_3 converge to zero as training progresses, regardless of algorithm used. Similar to the observation for the CE loss in [18], although all algorithms lead to \mathcal{NC} solutions, networks trained with different algorithms have notably different generalization performances.⁷ We find the cosine distribution \mathcal{P}_{CM} consistently aligns with the test accuracy, the more and higher. This may due to the fact that different training methods have different converge rate during the terminal phase of training, and it further lead to different distribution of features.

Improving network efficiency via fixing classifiers as simplex ETFs. In Theorem 3.1, when $d \geq K - 1$ and the weight decay terms are properly chosen, we showed that the optimal classifier for the vanilla MSE loss is a simplex ETF. This implies that we can (i) fix the last-layer classifier as a simplex ETF, and (ii) reduce the feature dimension $d = K$. By doing so, we substantially reduce the number of trainable parameters without sacrificing the generalization performance as shown in Figure 5.

Choice of the feature dimension d . On the other hand, Theorem 3.1 shows that the optimal class means $\overline{\mathbf{H}}^*$ form a simplex ETF only when $d \geq K - 1$. If $d < K - 1$, then the global solution $\overline{\mathbf{H}}^*$ is only the best rank- d approximation of the simplex ETF, where the class-means of the each class neither have equal length nor are maximally distant. To demonstrate its effect, we run experiments on the CIFAR10 dataset using vanilla MSE loss and ResNet18, with both $d < K - 1$ and $d \geq K - 1$. As shown in Figure 4, even though all cases exhibit \mathcal{NC} , choosing $d \geq K - 1$ is crucial for fitting the training data and generalization to test data. This is also corroborated by observing \mathcal{P}_{CM} , which

⁶The probability margin cannot be adopted here because probability is not well-defined given that softmax is not used in the MSE loss.

⁷L-BFGS with strong Wolfe line-search strategy may result in quite small stepsize at the terminal phase of training. We think that L-BFGS with proper diminishing stepsize can improve the generalization ability.

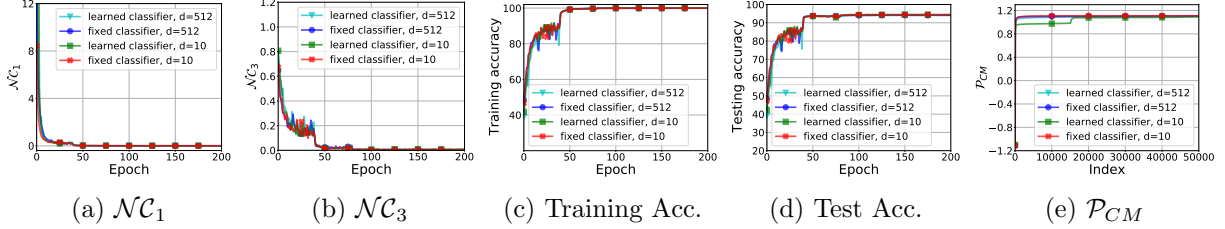


Figure 5: **Comparison of the performances on learned vs. fixed last-layer classifiers.** We compare within-class variation collapse \mathcal{NC}_1 , self-duality \mathcal{NC}_3 , training accuracy, test accuracy and cosine margin distribution \mathcal{P}_{CM} on fixed and learned classifier on CIFAR10-ResNet18 with data augmentation. The network is trained by SGD optimizer.

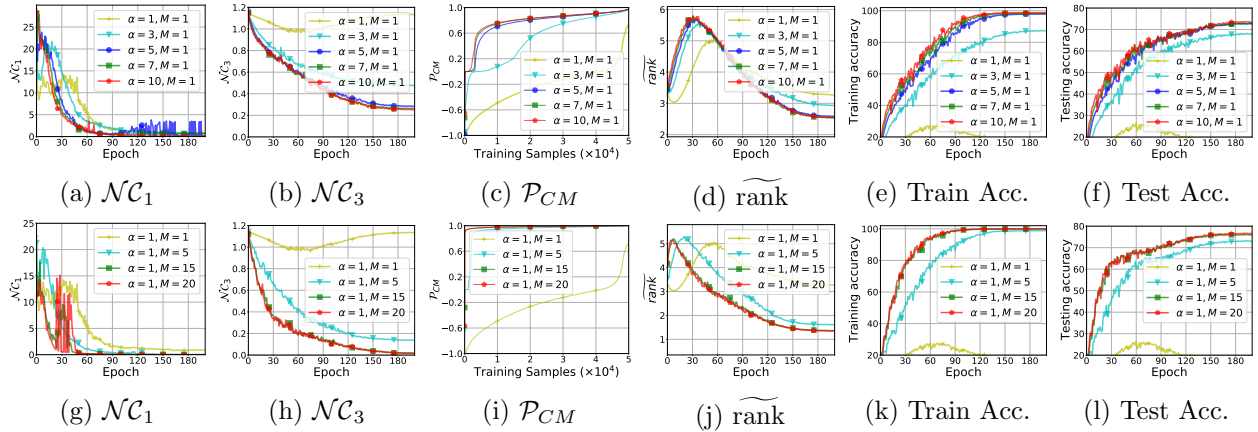


Figure 6: **Effects of rescaling parameters α and M .** Experiments are conducted on the miniImageNet dataset (MIN) with a ResNet18 backbone. Top row shows the result of varying α with fixed M . Bottom row shows the result of varying M with fixed α .

shows that more training samples lie on the decision boundary (i.e., $CM_{k,i} = 0$) as d decreases in the range of $d < K - 1$. As shown in Figure 4(e), this is in sharp contrast to CE loss which produces similar performance for different d . Note that all the existing work on CE loss [11, 13–18] only study the case when $d \geq K$. In the Appendix, we visually compare the features learned by CE and MSE, but we leave the thorough analysis for CE loss as future work.

Experiments of the rescaled MSE loss. In Section 3.3, we argued through landscape visualization that rescaling improves the optimization landscape for the MSE loss around the global solutions. Here, we corroborate our findings via experiments, showing that rescaling of MSE indeed leads to better \mathcal{NC} and hence better optimization landscapes. In particular, we empirically examine the effect of the two rescaling parameters (α, M) on the \mathcal{NC} phenomenon and the generalization performance. In Figure 6, we run experiments on the miniImageNet [57] dataset with ResNet18 [38]. We notice that when one scaling factor is fixed, the other scaling parameter has a positive correlation with the degree of \mathcal{NC} as well as the training and test performances. This observation is well-aligned with our analysis in Section 3.3.

5 Conclusion

In this work, we provide a global landscape analysis for deep neural networks trained via the MSE loss, under the unconstrained feature model. Our theoretical results reveal that all global solutions exhibit the \mathcal{NC} phenomenon, and that the global landscape is benign in the sense that it does not have spurious local minimizers. Such results extend the scope where \mathcal{NC} provably occurs with the MSE loss, which was restricted to neural networks trained via particular and unrealistic algorithms in prior work [12, 15]. More broadly, our results extend the scope of the “prevalence of neural collapse” in the seminal work [11], which was restricted to neural networks trained via the CE loss. Combined with the results in [18], the prevalence of neural collapse now subsumes (at least) that deep neural networks trained for classification tasks with both CE and MSE losses exhibit neural collapse, regardless of the training algorithm (as long as it can escape strict saddle points) and network architecture (as long as it is sufficiently expressive).

Towards designing better loss functions. As a future work, the improved understanding of \mathcal{NC} with different choices of loss functions may help us to study and demystify the role of loss design for learning more generalizable and transferable deep features [58–62]. The fact that both CE and MSE exhibit the \mathcal{NC} does not mean that they are equally good at inducing neural collapse solutions in practical neural network training. As shown in our experiments, rescaling of the MSE loss is indispensable for improving \mathcal{NC} hence producing better test performance over the vanilla MSE loss. There is, however, no reason to be satisfied with the rescaled MSE loss since it is heuristically designed and does not have any justification on its “optimality”. Even though we are able to offer insights into the benefits of rescaling for MSE loss via landscape visualization, our explanation is approximate, based on extravagant simplifications of the optimization problem (by using two parameters θ and s to summarize a very high-dimensional landscape!). In practice, all the optimization variables \mathbf{H} , \mathbf{W} and \mathbf{b} are intricately correlated, and the insights gained from the visualization via simplification may hardly be useful for the design of new loss functions. The derivation of an “optimal” loss functions for inducing \mathcal{NC} may require the development of new analysis techniques which we leave as future work.

Acknowledgements

ZZ acknowledges support from NSF grants CCF 2008460 and CCF 2106881. XL and QQ acknowledge support from NSF grant DMS 2009752 and NSF Career Award 2143904. We also acknowledge Sheng Liu (NYU CDS) and Kangning Liu (NYU CDS) for fruitful discussion during various stages of the work.

References

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [2] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- [3] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- [4] Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli Qin, Augustin Židek, Alexander WR Nelson, Alex Bridgland, et al. Improved protein structure prediction using potentials from deep learning. *Nature*, 577(7792):706–710, 2020.
- [5] Haozhi Qi, Chong You, Xiaolong Wang, Yi Ma, and Jitendra Malik. Deep isometric learning for visual recognition. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 7824–7835. PMLR, 13–18 Jul 2020.
- [6] James Martens, Andy Ballard, Guillaume Desjardins, Grzegorz Swirszcz, Valentin Dalibard, Jascha Sohl-Dickstein, and Samuel S Schoenholz. Rapid training of deep neural networks without skip connections or normalization layers using deep kernel shaping. *arXiv preprint arXiv:2110.01765*, 2021.
- [7] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *arXiv preprint arXiv:1912.02292*, 2019.
- [8] Zitong Yang, Yaodong Yu, Chong You, Jacob Steinhardt, and Yi Ma. Rethinking bias-variance trade-off for generalization of neural networks. In *International Conference on Machine Learning*, pages 10767–10777. PMLR, 2020.
- [9] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- [10] 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.
- [11] Vardan Papayan, 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.
- [12] X.Y. Han, Vardan Papayan, and David L. Donoho. Neural collapse under MSE loss: Proximity to and dynamics on the central path. In *International Conference on Learning Representations*, 2022.
- [13] Vardan Papayan. Traces of class/cross-class structure pervade deep learning spectra. *Journal of Machine Learning Research*, 21(252):1–64, 2020.
- [14] Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Layer-peeled model: Toward understanding well-trained deep neural networks. *arXiv preprint arXiv:2101.12699*, 2021.
- [15] Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *arXiv preprint arXiv:2011.11619*, 2020.
- [16] Florian Graf, Christoph Hofer, Marc Niethammer, and Roland Kwitt. Dissecting supervised contrastive learning. In *International Conference on Machine Learning*, pages 3821–3830. PMLR, 2021.
- [17] Tolga Ergen and Mert Pilanci. Revealing the structure of deep neural networks via convex duality. In *International Conference on Machine Learning*, pages 3004–3014. PMLR, 2021.
- [18] 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*, 2021.
- [19] G Cybenko. Approximation by superposition of sigmoidal functions. *Mathematics of Control, Signals and Systems*, 2(4):303–314, 1989.

- [20] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- [21] Zhou Lu, Hongming Pu, Feicheng Wang, Zhiqiang Hu, and Liwei Wang. The expressive power of neural networks: a view from the width. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6232–6240, 2017.
- [22] Uri Shaham, Alexander Cloninger, and Ronald R Coifman. Provable approximation properties for deep neural networks. *Applied and Computational Harmonic Analysis*, 44(3):537–557, 2018.
- [23] Jianfeng Lu and Stefan Steinerberger. Neural collapse with cross-entropy loss. *arXiv preprint arXiv:2012.08465*, 2020.
- [24] E Weinan and Stephan Wojtowytsch. On the emergence of tetrahedral symmetry in the final and penultimate layers of neural network classifiers. *arXiv preprint arXiv:2012.05420*, 2020.
- [25] Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J Su. An unconstrained layer-peeled perspective on neural collapse. *arXiv preprint arXiv:2110.02796*, 2021.
- [26] Tom Tirer and Joan Bruna. Extended unconstrained features model for exploring deep neural collapse. *arXiv preprint arXiv:2202.08087*, 2022.
- [27] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1026–1034, 2015.
- [28] Tomaso Poggio and Qianli Liao. Explicit regularization and implicit bias in deep network classifiers trained with the square loss. *arXiv preprint arXiv:2101.00072*, 2020.
- [29] Tomaso Poggio and Qianli Liao. Implicit dynamic regularization in deep networks. Technical report, Center for Brains, Minds and Machines (CBMM), 2020.
- [30] Akshay Rangamani, Mengjia Xu, Andrzej Banburski, Qianli Liao, and Tomaso Poggio. Dynamics and neural collapse in deep classifiers trained with the square loss. *Technical report, Center for Brains, Minds and Machines (CBMM)*, 2021.
- [31] Florian Graf, Christoph Hofer, Marc Niethammer, and Roland Kwitt. Neural collapse in deep homogeneous classifiers and the role of weight decay. In *International Conference on Machine Learning*, pages 3821–3830. PMLR, 2021.
- [32] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Proceedings of The 28th Conference on Learning Theory*, pages 797–842, 2015.
- [33] Ju Sun, Qing Qu, and John Wright. When are nonconvex problems not scary? *arXiv preprint arXiv:1510.06096*, 2015.
- [34] Yuqian Zhang, Qing Qu, and John Wright. From symmetry to geometry: Tractable nonconvex problems. *arXiv preprint arXiv:2007.06753*, 2020.
- [35] 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.
- [36] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [37] 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.
- [38] 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.

- [39] Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere i: Overview and the geometric picture. *IEEE Transactions on Information Theory*, 63(2):853–884, 2016.
- [40] Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics*, 18(5):1131–1198, 2018.
- [41] Qing Qu, Yuexiang Zhai, Xiao Li, Yuqian Zhang, and Zhihui Zhu. Geometric analysis of nonconvex optimization landscapes for overcomplete learning. In *International Conference on Learning Representations*, 2020.
- [42] Qing Qu, Zhihui Zhu, Xiao Li, Manolis C. Tsakiris, John Wright, and René Vidal. Finding the sparsest vectors in a subspace: Theory, algorithms, and applications. *arXiv preprint arXiv:2001.06970*, 2020.
- [43] Benjamin D Haeffele and René Vidal. Global optimality in tensor factorization, deep learning, and beyond. *arXiv preprint arXiv:1506.07540*, 2015.
- [44] Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *arXiv preprint arXiv:1605.07272*, 2016.
- [45] Srinadh Bhojanapalli, Behnam Neyshabur, and Nathan Srebro. Global optimality of local search for low rank matrix recovery. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 3880–3888, 2016.
- [46] Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *International Conference on Machine Learning*, pages 1233–1242. PMLR, 2017.
- [47] Qiuwei Li, Zhihui Zhu, and Gongguo Tang. The non-convex geometry of low-rank matrix optimization. *Information and Inference: A Journal of the IMA*, 8(1):51–96, 2019.
- [48] Xingguo Li, Junwei Lu, Raman Arora, Jarvis Haupt, Han Liu, Zhaoran Wang, and Tuo Zhao. Symmetry, saddle points, and global optimization landscape of nonconvex matrix factorization. *IEEE Transactions on Information Theory*, 65(6):3489–3514, 2019.
- [49] Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- [50] Samuel Burer and Renato DC Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Mathematical Programming*, 95(2):329–357, 2003.
- [51] Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Reexamining low rank matrix factorization for trace norm regularization. *arXiv preprint arXiv:1706.08934*, 2017.
- [52] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.
- [53] Jason D Lee, Max Simchowitz, Michael I Jordan, and Benjamin Recht. Gradient descent only converges to minimizers. In *Conference on learning theory*, pages 1246–1257. PMLR, 2016.
- [54] Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [55] Miles E. Lopes. Estimating unknown sparsity in compressed sensing. *arXiv preprint arXiv:1204.4227*, 2013.
- [56] Andrzej Banburski, Fernanda De La Torre, Nishka Pant, Ishana Shastri, and Tomaso Poggio. Distribution of classification margins: Are all data equal? *arXiv preprint arXiv:2107.10199*, 2021.
- [57] Oriol Vinyals, Charles Blundell, Timothy P. Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *NIPS*, 2016.
- [58] Simon Kornblith, Ting Chen, Honglak Lee, and Mohammad Norouzi. Why do better loss functions lead to less transferable features? *Advances in Neural Information Processing Systems*, 34, 2021.
- [59] Nishanth Dikkala, Gal Kaplun, and Rina Panigrahy. For manifold learning, deep neural networks can be locality sensitive hash functions. *arXiv preprint arXiv:2103.06875*, 2021.

- [60] Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding generalization in deep learning. *arXiv preprint arXiv:2202.08384*, 2022.
- [61] Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learning. In *International Conference on Learning Representations*, 2022.
- [62] Ido Ben-Shaul and Shai Dekel. Nearest class-center simplification through intermediate layers. *arXiv preprint arXiv:2201.08924*, 2022.
- [63] Thomas Strohmer and Robert W Heath Jr. Grassmannian frames with applications to coding and communication. *Applied and computational harmonic analysis*, 14(3):257–275, 2003.
- [64] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *ICLR*, 2017.

Appendices

Notations and Organizations. For a scalar function $f(\mathbf{Z})$ with a variable $\mathbf{Z} \in \mathbb{R}^{K \times N}$, its Hessian can be represented by a bilinear form defined via $[\nabla^2 f(\mathbf{Z})](\mathbf{A}, \mathbf{B}) = \sum_{i,j,k,\ell} \frac{\partial^2 f(\mathbf{Z})}{\partial z_{ij} \partial z_{k\ell}} a_{ij} b_{k\ell}$ for any $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{K \times N}$, which avoids representing the Hessian as a tensor, or vectorizing the variable \mathbf{Z} . We will use the bilinear form for the Hessian throughout the Appendix. Now we give the formal definition of Simplex ETF.

Definition .1 (K -Simplex ETF [11, 63]) A standard Simplex ETF is a collection of points in \mathbb{R}^K specified by the columns of

$$\mathbf{M} = \sqrt{\frac{K}{K-1}} \left(\mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top \right),$$

where $\mathbf{I}_K \in \mathbb{R}^{K \times K}$ is the identity matrix, and $\mathbf{1}_K \in \mathbb{R}^K$ is the all ones vector.

As in [11, 14], in this paper we consider general Simplex ETF as a collection of points in \mathbb{R}^d specified by the columns of $\sqrt{\frac{K}{K-1}} \mathbf{P} \left(\mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top \right)$, where (i) when $d \geq K$, $\mathbf{P} \in \mathbb{R}^{d \times K}$ is an orthonormal matrix, i.e., $\mathbf{P}^\top \mathbf{P} = \mathbf{I}_K$, and (ii) when $d = K - 1$, \mathbf{P} is chosen such that $\left[\mathbf{P}^\top \quad \frac{1}{\sqrt{K}} \mathbf{1}_K \right]$ is an orthonormal matrix.

The appendix is organized as follows. In Appendix A, we describe the datasets, network architectures and training settings. In Appendix B, we provide a detailed proof for Theorem 3.1, analyzing the global minimizers to our regularized MSE loss. Finally, in Appendix C we provide additional details for obtaining the visualization of rescaled MSE and CE losses presented in Section 3.3.

A Technical Details of the Experimental Setup in Section 4

In Section 4, we conduct experiments on CIFAR10 [36] and miniImageNet [57] datasets. We note that for miniImageNet dataset, since we are not doing few-shot learning where the work [57] primarily considers, we split the total 60000 images into training set (50000 images) and validation set (10000 images) such that both training and validation set include the full 100 classes. All images from the datasets are normalized by their mean and variance channel-wise. We use the ResNet18 [38] architecture throughout all the experiments. For CIFAR10, we use the same experiment setting in [18] except the replacement of CE loss by standard MSE loss for fair comparison. Specifically, we train ResNet18 for 200 epochs with three different optimizers: SGD, Adam and LBFGS. For SGD, the initial learning rate and momentum are set to 0.05 and 0.9, respectively. For Adam, the initial learning rate, β_1 and β_2 are set to 0.001, 0.9 and 0.999, respectively. We decay the learning rate by 0.1 every 40 epochs for SGD and Adam. We use LBFGS with an initial learning rate of 0.01 and strong Wolfe line search strategy for subsequent iterations. Without explicitly mentioned, we use the weight decay of 5×10^{-4} and the same data augmentation in [18] for all experiments on CIFAR10. For miniImageNet, we use the rescaled MSE loss as described in Section 2.2 with the SGD optimizer with an initial learning rate 0.01, momentum 0.9 and weight decay 0.001. We use a Cosine Annealing Warm Restarts [64] learning rate scheduler where the number of epochs before the first restart is set as 200 and the minimum learning rate is 0.0001.

Three \mathcal{NC} measures \mathcal{NC}_1 - \mathcal{NC}_3 [11, 18] For the sake of completeness, we describe the three \mathcal{NC} measures \mathcal{NC}_1 - \mathcal{NC}_3 [11, 18] used in Section 4. Towards that end, first define the global mean

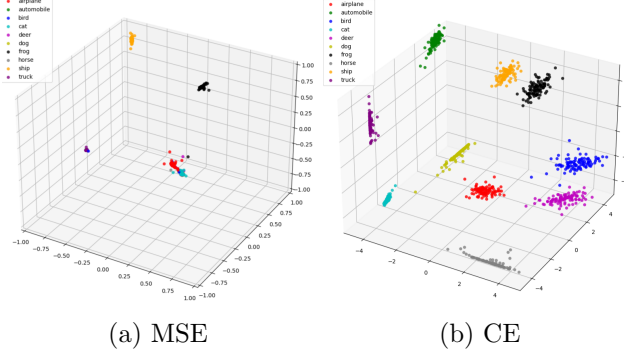


Figure 7: **Visual comparison of features learned by MSE and CE losses with feature dimension $d = 3$.** We compare the training feature distribution by setting the feature dimension $d = 3$ for ResNet18 and training it with CIFAR10. The network is trained by the SGD optimizer.

of the last-layer features $\{\mathbf{h}_{k,i}\}$ as $\mathbf{h}_G = \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \mathbf{h}_{k,i}$ and the class mean as $\bar{\mathbf{h}}_k = \frac{1}{n} \sum_{i=1}^n \mathbf{h}_{k,i}$ ($1 \leq k \leq K$).

- \mathcal{NC}_1 . We measure the within-class variability collapse by

$$\mathcal{NC}_1 := \frac{1}{K} \text{trace} \left(\boldsymbol{\Sigma}_W \boldsymbol{\Sigma}_B^\dagger \right), \quad (6)$$

where $\boldsymbol{\Sigma}_W := \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n (\mathbf{h}_{k,i} - \bar{\mathbf{h}}_k) (\mathbf{h}_{k,i} - \bar{\mathbf{h}}_k)^\top \in \mathbb{R}^{d \times d}$ denotes the within-class covariance of the features, $\boldsymbol{\Sigma}_B := \frac{1}{K} \sum_{k=1}^K (\bar{\mathbf{h}}_k - \mathbf{h}_G) (\bar{\mathbf{h}}_k - \mathbf{h}_G)^\top \in \mathbb{R}^{d \times d}$ represents the between-class covariance, and $\boldsymbol{\Sigma}_B^\dagger$ denotes the pseudo inverse of $\boldsymbol{\Sigma}_B$.

- \mathcal{NC}_2 . We measure the onvergence of the learned classifier $\mathbf{W} \in \mathbb{R}^{K \times d}$ (for $d \geq K - 1$) to a Simplex ETF by

$$\mathcal{NC}_2 := \left\| \frac{\mathbf{W}\mathbf{W}^\top}{\|\mathbf{W}\mathbf{W}^\top\|_F} - \frac{1}{\sqrt{K-1}} \left(\mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top \right) \right\|_F, \quad (7)$$

where the Simplex ETF and $\mathbf{W}\mathbf{W}^\top$ are rescaled to have unit energy (in Frobenius norm).

- \mathcal{NC}_3 . For $d \geq K - 1$, we measure the convergence to self-duality between the learned features \mathbf{H} and the learned classifier \mathbf{W} via

$$\mathcal{NC}_3 := \left\| \frac{\mathbf{W}\bar{\mathbf{H}}}{\|\mathbf{W}\bar{\mathbf{H}}\|_F} - \frac{1}{\sqrt{K-1}} \left(\mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top \right) \right\|_F, \quad (8)$$

where $\bar{\mathbf{H}} := [\bar{\mathbf{h}}_1 - \mathbf{h}_G \ \cdots \ \bar{\mathbf{h}}_K - \mathbf{h}_G] \in \mathbb{R}^{d \times K}$ are the centered class-means.

Visual comparison of features learned by MSE and CE losses with feature dimension $d = 3$. To visualize the learned features, we set the feature dimension $d = 3$ for ResNet18 and train it with CIFAR10. Figure 7 display the learned features with MSE loss and CE loss on randomly selected 100 training samples for each class. We observe that the features learned by CE loss is more diverse and discriminative than MSE loss.

B Proof of Theorem 3.1 in Section 3.1

In this part of appendices, we prove Theorem 3.1 in Section 3 that we restate as follows.

Theorem B.1 (Global Optimality Condition) Let $(\mathbf{W}^*, \mathbf{H}^*, \mathbf{b}^*)$ be a global minimizer of

$$\min_{\mathbf{W}, \mathbf{H}, \mathbf{b}} f(\mathbf{W}, \mathbf{H}, \mathbf{b}) := \frac{1}{2N} \left\| \mathbf{W}\mathbf{H} + \mathbf{b}\mathbf{1}^\top - \mathbf{Y} \right\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 + \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2. \quad (9)$$

Then $(\mathbf{W}^*, \mathbf{H}^*, \mathbf{b}^*)$ satisfies:

(NC1,3) If $\lambda_{\mathbf{W}}\lambda_{\mathbf{H}} < \frac{1}{NK}$, then $(\mathbf{W}^*, \mathbf{H}^*)$ satisfies NC1 and NC3 as

$$\mathbf{h}_{k,i}^* = \bar{\mathbf{h}}_k^*, \quad \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}n}} \mathbf{w}^{*k} = \bar{\mathbf{h}}_k^*, \quad \forall k \in [K], i \in [n].$$

Otherwise, if $\lambda_{\mathbf{W}}\lambda_{\mathbf{H}} \geq \frac{1}{NK}$, then $\mathbf{W}^* = \mathbf{0}$ and $\mathbf{H}^* = \mathbf{0}$.

(NC2) If $\lambda_{\mathbf{W}}\lambda_{\mathbf{H}} < \frac{1}{NK}$, then $\bar{\mathbf{H}}^*$ further obeys the following properties for different d :

1. If $d < K - 1$: we have $\bar{\mathbf{H}}^{*\top} \bar{\mathbf{H}}^* = C_1 \mathcal{P}_d(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top)$, where $\mathcal{P}_d(\mathbf{M})$ denotes the best rank- d approximating of \mathbf{M} ;
2. If $d = K - 1$: we have $\bar{\mathbf{H}}^{*\top} \bar{\mathbf{H}}^* = C_2(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top)$;
3. If $d \geq K$: we have

$$\bar{\mathbf{H}}^{*\top} \bar{\mathbf{H}}^* = \begin{cases} C_3 \left(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top \right), & \lambda_{\mathbf{b}} \leq \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}, \\ C_4 \left(\mathbf{I} - \frac{\sqrt{n\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{\lambda_{\mathbf{b}}(1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}})} \mathbf{1}_K \mathbf{1}_K^\top \right), & \text{otherwise,} \end{cases}$$

where $\frac{\sqrt{n\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{\lambda_{\mathbf{b}}(1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}})} \leq \frac{1}{K}$ in the second case since $\lambda_{\mathbf{b}} \geq \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}$.

Here, C_1, C_2, C_3 , and C_4 are some positive numerical constants that depend on $\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}, \lambda_{\mathbf{b}}$.

(Bias) The bias satisfies $\mathbf{b}^* = b^* \mathbf{1}_K$ with $b^* \leq \frac{1}{K}$ given by:

1. If $d < K$: we have $b^* = \frac{1}{K(\lambda_{\mathbf{b}} + 1)}$;
2. If $d \geq K$: we have $b^* = \begin{cases} \frac{1}{K(\lambda_{\mathbf{b}} + 1)}, & \lambda_{\mathbf{b}} \leq \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}, \\ \frac{\sqrt{n\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{\lambda_{\mathbf{b}}}, & \text{otherwise.} \end{cases}$

In particular, when $\lambda_{\mathbf{b}} \rightarrow 0$, we have $b^* \rightarrow \frac{1}{K}$; when $\lambda_{\mathbf{b}} \rightarrow \infty$, we have $b^* \rightarrow 0$.

B.1 Main Proof

Proof [Proof of Theorem B.1] We first characterize the solutions (\mathbf{W}, \mathbf{H}) in terms of \mathbf{b} . Denote by $\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbf{b}\mathbf{1}^\top$ and let $\tilde{\mathbf{Y}} = \mathbf{U}\Sigma\mathbf{V}^\top = \sum_{i=1}^K \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ be its SVD, where $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_K \geq 0$ are the singular values. For convenience, we denote by $\tilde{\lambda} = N\sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}$. By Theorem B.2, we know

$$f(\mathbf{W}, \mathbf{H}, \mathbf{b}) \geq \frac{\lambda_{\mathbf{b}}}{2} \|\mathbf{b}\|_2^2 + \frac{1}{N} \cdot \begin{cases} \sum_{i=1}^K \frac{1}{2} \left(\sigma_i - [\sigma_i - \tilde{\lambda}]_+ \right)^2 + \tilde{\lambda} [\sigma_i - \tilde{\lambda}]_+, & d \geq K \\ \sum_{i=1}^d \frac{1}{2} \left(\sigma_i - [\sigma_i - \tilde{\lambda}]_+ \right)^2 + \tilde{\lambda} [\sigma_i - \tilde{\lambda}]_+ + \sum_{i=d+1}^K \frac{1}{2} \sigma_i^2, & d < K \end{cases} \quad (10)$$

where the inequality becomes an equality when $\mathbf{W}\mathbf{H} = \sum_{i=1}^{\min(d,K)} [\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}]_+ \mathbf{u}_i \mathbf{v}_i^\top$.

Noting that the singular values σ_i also depend on \mathbf{b} , to minimize the right hand side (RHS) of (10) in terms of \mathbf{b} , we first rewrite each term involving the singular value as

$$\frac{1}{2} \left(\sigma_i - [\sigma_i - \tilde{\lambda}]_+ \right)^2 + \tilde{\lambda} [\sigma_i - \tilde{\lambda}]_+ = \begin{cases} \frac{1}{2} \sigma_i^2, & \sigma_i \leq \tilde{\lambda}, \\ \tilde{\lambda} \sigma_i - \frac{1}{2} \tilde{\lambda}^2, & \sigma_i \geq \tilde{\lambda}, \end{cases} \quad (11)$$

where for both cases it increases as σ_i increases. Thus, for any \mathbf{b} with the same energy, say c , minimizing the RHS of (10) is equivalent to minimizing the singular values σ_i . With this in mind, we now show that if \mathbf{b}^* is a minimizer to RHS of (10), then $\|\mathbf{b}^*\| \leq \frac{1}{\sqrt{K}}$. By Theorem B.3, we know for any \mathbf{b} we have $\sigma_2 = \sigma_3 = \dots = \sigma_{K-1} = \sqrt{n}$ and $\sigma_1 \geq \sqrt{n}$ (see (25)). On the other hand, when $\mathbf{b} = \frac{1}{K}\mathbf{1}$, we have $\sigma_1 = \sigma_2 = \dots = \sigma_{K-1} = \sqrt{n}$ and $\sigma_K = 0$, which are the smallest possible singular values that can be achieved. Thus, considering the weight decay term on (10), the minimizer \mathbf{b}^* must satisfy $\|\mathbf{b}^*\| \leq \|\frac{1}{K}\mathbf{1}\| = \frac{1}{\sqrt{K}}$.

Therefore, we only need to optimize over \mathbf{b} with $\|\mathbf{b}\| = c \leq \frac{1}{\sqrt{K}}$. In this case, it follows from Theorem B.3 that $\sigma_2 = \dots = \sigma_{K-1} = \sqrt{n}$, $\sigma_1 \geq \sqrt{n}$, $\sigma_K \geq \sqrt{n}(1 - \sqrt{K}c)$, and both inequalities become equalities *if and only if* $\mathbf{b} = \frac{c}{\sqrt{K}}\mathbf{1}$. The remaining is to optimize the RHS of (10) in terms of σ_K which depends on c . By (10) and (11), this problem reduces to

$$\min_{0 \leq c \leq \frac{1}{\sqrt{K}}} \frac{\lambda_{\mathbf{b}}}{2} c^2 + \frac{n}{2N} (1 - \sqrt{K}c)^2 \quad (12)$$

if $d < K$, and otherwise reduces to

$$\begin{cases} \min_{0 \leq c \leq \frac{1}{\sqrt{K}}} \frac{\lambda_{\mathbf{b}}}{2} c^2 + \frac{\tilde{\lambda}}{N} \left(\sqrt{n} (1 - \sqrt{K}c) - \frac{1}{2}\tilde{\lambda} \right), & \sqrt{n} (1 - \sqrt{K}c) \geq \tilde{\lambda} \\ \min_{0 \leq c \leq \frac{1}{\sqrt{K}}} \frac{\lambda_{\mathbf{b}}}{2} c^2 + \frac{n}{2N} (1 - \sqrt{K}c)^2, & \sqrt{n} (1 - \sqrt{K}c) \leq \tilde{\lambda} \end{cases} \quad (13)$$

We now consider the two cases as follows:

1. Case I: $d < K$. In this case, the problem (12) achieves its minimum at $c^* = \frac{1}{\sqrt{K}(\lambda_{\mathbf{b}}+1)}$.
2. Case II: $d \geq K$. In this case, when $c \geq \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$, problem (13) becomes (12), and thus its minimum among $c \geq \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$ is $c^* = \max\left(\frac{1}{\sqrt{K}(\lambda_{\mathbf{b}}+1)}, \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)\right)$. On the other hand, when $c \leq \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$, the problem (13) is also a quadratic function on c and achieves its minimum among $c \leq \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$ is $c^* = \min\left(\frac{\tilde{\lambda}}{\lambda_{\mathbf{b}}\sqrt{N}}, \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)\right)$.

We now find the minimum value among these two cases. When $\frac{1}{\sqrt{K}(\lambda_{\mathbf{b}}+1)} \geq \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$, i.e., $\lambda_{\mathbf{b}} \leq \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}$, we have $\frac{\tilde{\lambda}}{\lambda_{\mathbf{b}}\sqrt{N}} \geq \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$, which together with the form of the two quadratic functions implies that the minimum is achieved when $c^* = \frac{1}{\sqrt{K}(\lambda_{\mathbf{b}}+1)}$. On the other hand, when $\frac{1}{\sqrt{K}(\lambda_{\mathbf{b}}+1)} < \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$, i.e., $\lambda_{\mathbf{b}} > \frac{\sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{1 - \sqrt{KN\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}$, we have $\frac{\tilde{\lambda}}{\lambda_{\mathbf{b}}\sqrt{N}} < \frac{1}{\sqrt{K}} \left(1 - \frac{\tilde{\lambda}}{\sqrt{n}}\right)$, which together with the form of the two quadratic functions implies that the minimum is achieved when $c^* = \frac{\tilde{\lambda}}{\lambda_{\mathbf{b}}\sqrt{N}} = \frac{\sqrt{N\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}}{\lambda_{\mathbf{b}}}$. Thus, we can also conclude that $c^* \rightarrow 0$ when $\lambda_{\mathbf{b}} \rightarrow \infty$ and $c^* \rightarrow \frac{1}{\sqrt{K}}$ when $\lambda_{\mathbf{b}} \rightarrow 0$.

The proof is completed by invoking Theorem B.4 to characterize $(\mathbf{W}^*, \mathbf{H}^*)$. ■

B.2 Supporting Lemmas

We first characterize the following balance property between \mathbf{W} and \mathbf{H} for any critical point $(\mathbf{W}, \mathbf{H}, \mathbf{b})$ of our loss function:

Lemma B.2 For any K, d, N , and $\tilde{\mathbf{Y}} \in \mathbb{R}^{K \times N}$ with SVD given by $\tilde{\mathbf{Y}} = \mathbf{U}\Sigma\mathbf{V}^\top = \sum_{i=1}^K \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ where $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_K \geq 0$ are the singular values, the following problem

$$\min_{\mathbf{W} \in \mathbb{R}^{K \times d}, \mathbf{H} \in \mathbb{R}^{d \times N}} \xi(\mathbf{W}, \mathbf{H}) = \frac{1}{2} \left\| \mathbf{W}\mathbf{H} - \tilde{\mathbf{Y}} \right\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2 \quad (14)$$

is a strict saddle function with no spurious local minimizer, in the sense that

- Any local minimizer $(\mathbf{W}^*, \mathbf{H}^*)$ of (14) is a global minimizer of (14), with the following form

$$\mathbf{W}^* \mathbf{H}^* = \sum_{i=1}^{\min(d, K)} \eta_i \mathbf{u}_i \mathbf{v}_i^\top,$$

where we let $\eta_i(\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}) := [\sigma_i - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}]_+$. Correspondingly, the minimal objective value of (14) is

$$\xi_\star = \begin{cases} \sum_{i=1}^K \frac{1}{2} (\sigma_i - \eta_i)^2 + \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \eta_i, & d \geq K \\ \sum_{i=1}^d \frac{1}{2} (\sigma_i - \eta_i)^2 + \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \eta_i + \sum_{i=d+1}^K \sigma_i^2, & d < K \end{cases} \quad (15)$$

- Any critical point (\mathbf{W}, \mathbf{H}) of (14) that is not a local minimizer is a strict saddle with negative curvature, i.e. the Hessian at this critical point has at least one negative eigenvalue.

Proof [Proof of Lemma B.2] By definition, any critical point (\mathbf{W}, \mathbf{H}) of (14) satisfies the following:

$$\nabla_{\mathbf{W}} \xi(\mathbf{W}, \mathbf{H}) = (\mathbf{W}\mathbf{H} - \tilde{\mathbf{Y}}) \mathbf{H}^\top + \lambda_{\mathbf{W}} \mathbf{W} = \mathbf{0},$$

$$\nabla_{\mathbf{H}} \xi(\mathbf{W}, \mathbf{H}) = \mathbf{W}^\top (\mathbf{W}\mathbf{H} - \tilde{\mathbf{Y}}) + \lambda_{\mathbf{H}} \mathbf{H} = \mathbf{0}.$$

By left multiplying the first equation by \mathbf{W}^\top on both sides and then right multiplying second equation by \mathbf{H}^\top on both sides and combining the equations together, we obtain

$$\lambda_{\mathbf{W}} \mathbf{W}^\top \mathbf{W} = \lambda_{\mathbf{H}} \mathbf{H} \mathbf{H}^\top. \quad (16)$$

This further gives

$$\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}} \mathbf{W} \mathbf{W}^\top \mathbf{W} + \lambda_{\mathbf{W}} \mathbf{W} = \tilde{\mathbf{Y}} \mathbf{H}^\top, \quad (17)$$

$$\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}} \mathbf{H}^\top \mathbf{H} \mathbf{H}^\top + \lambda_{\mathbf{H}} \mathbf{H}^\top = \tilde{\mathbf{Y}}^\top \mathbf{W}.$$

In the following, without loss of generality, we assume that the critical point (\mathbf{W}, \mathbf{H}) satisfying the above equations has the form

$$\mathbf{W} = \begin{bmatrix} \widehat{\mathbf{W}} & \mathbf{0} \end{bmatrix}, \quad \mathbf{H} = \begin{bmatrix} \widehat{\mathbf{H}} \\ \mathbf{0} \end{bmatrix} \quad (18)$$

where the columns of $\widehat{\mathbf{W}}$ are orthogonal and the rows of $\widehat{\mathbf{H}}$ are orthogonal, and the zeros $\mathbf{0}$ in \mathbf{W} and \mathbf{H} might or might not exist depending on the rank of \mathbf{W} and \mathbf{H} . The underlying reasoning is that, for any \mathbf{W} satisfying (17), the Gram-Schmidt process implies that we can always orthogonalize \mathbf{W} by an orthonormal matrix $\mathbf{R} \in \mathbb{R}^{d \times d}$ (i.e., $\mathbf{R}^\top \mathbf{R} = \mathbf{R} \mathbf{R}^\top = \mathbf{I}$), such that $\widetilde{\mathbf{W}} = \mathbf{W} \mathbf{R} = \begin{bmatrix} \widehat{\mathbf{W}} & \mathbf{0} \end{bmatrix}$. On the other hand, let $\widetilde{\mathbf{H}} = \mathbf{R}^\top \mathbf{H}$. Because $\lambda_{\mathbf{W}} \mathbf{W}^\top \mathbf{W} = \lambda_{\mathbf{H}} \mathbf{H} \mathbf{H}^\top$, we have $\lambda_{\mathbf{W}} \widetilde{\mathbf{W}}^\top \widetilde{\mathbf{W}} = \lambda_{\mathbf{H}} \widetilde{\mathbf{H}} \widetilde{\mathbf{H}}^\top$, which implies that the rows of $\widetilde{\mathbf{H}}$ are also orthogonal. Therefore, multiply \mathbf{R} on both sides of (17), we always have

$$\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}} \widetilde{\mathbf{W}} \widetilde{\mathbf{W}}^\top \widetilde{\mathbf{W}} + \lambda_{\mathbf{W}} \widetilde{\mathbf{W}} = \widetilde{\mathbf{Y}} \widetilde{\mathbf{H}}^\top, \quad \frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}} \widetilde{\mathbf{H}}^\top \widetilde{\mathbf{H}} \widetilde{\mathbf{H}}^\top + \lambda_{\mathbf{H}} \widetilde{\mathbf{H}}^\top = \widetilde{\mathbf{Y}}^\top \widetilde{\mathbf{W}}.$$

Thus, we can verify that $(\widetilde{\mathbf{W}}, \widetilde{\mathbf{H}})$ is also a critical point with $\widetilde{\mathbf{W}} \widetilde{\mathbf{H}} = \mathbf{W} \mathbf{H}$ and has the same Hessian information as (\mathbf{W}, \mathbf{H}) . Thus, without the loss of generality, we can assume orthogonal (\mathbf{W}, \mathbf{H}) in the form (18), but with possible zero columns.

Form of the global solutions. Based on the orthogonalization, we further decompose (17) for all $i = 1, \dots, d$ columns of \mathbf{W} as

$$\begin{aligned} \left(\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}} \|\mathbf{w}_i\|^2 + \lambda_{\mathbf{W}} \right) \mathbf{w}_i &= \tilde{\mathbf{Y}} \mathbf{h}^i, \\ \left(\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}} \|\mathbf{h}^i\|^2 + \lambda_{\mathbf{H}} \right) \mathbf{h}^i &= \tilde{\mathbf{Y}}^\top \mathbf{w}_i, \end{aligned} \quad (19)$$

which implies that either (i) $\mathbf{w}_i = \mathbf{0}$ and $\mathbf{h}^i = \mathbf{0}$, or (ii) $\mathbf{w}_i, \mathbf{h}^i$ are the (scaled) left and right singular vectors of $\tilde{\mathbf{Y}}$. In particular, when $\mathbf{w}_i \neq \mathbf{0}$ and $\mathbf{h}_i \neq \mathbf{0}$, then by (16), it gives

$$\|\mathbf{h}^i\|^2 = \frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}} \|\mathbf{w}_i\|^2. \quad (20)$$

By further plugging the equation above into (19), it gives

$$\begin{aligned} \left(\sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \|\mathbf{w}_i\|^2 + \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right) \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|} &= \tilde{\mathbf{Y}} \frac{\mathbf{h}^i}{\|\mathbf{h}^i\|}, \\ \left(\sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \|\mathbf{w}_i\|^2 + \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right) \frac{\mathbf{h}^i}{\|\mathbf{h}^i\|} &= \tilde{\mathbf{Y}}^\top \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|}. \end{aligned} \quad (21)$$

Thus, when $\mathbf{w}_i \neq \mathbf{0}$ and $\mathbf{h}_i \neq \mathbf{0}$, we conclude that $\sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \|\mathbf{w}_i\|^2 + \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$ is a singular value of $\tilde{\mathbf{Y}}$, say σ_{i_j} , and $\frac{\mathbf{w}_i}{\|\mathbf{w}_i\|}$ and $\frac{\mathbf{h}^i}{\|\mathbf{h}^i\|}$ are the corresponding left and right singular vectors, respectively. In other words, when $\mathbf{w}_i \neq \mathbf{0}$ and $\mathbf{h}_i \neq \mathbf{0}$, then

$$\sigma_{i_j} = \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \|\mathbf{w}_i\|^2 + \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}, \quad \mathbf{u}_{i_j} = \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|}, \quad \mathbf{v}_{i_j} = \frac{\mathbf{h}^i}{\|\mathbf{h}^i\|} \quad (22)$$

for some i_j such that $\sigma_{i_j} > \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$. Together with (20), it further implies that

$$\mathbf{w}_i \mathbf{h}^{i\top} = \|\mathbf{w}_i\|_2^2 \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|_2} \frac{\mathbf{h}^{i\top}}{\|\mathbf{h}^i\|_2} = \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \|\mathbf{w}_i\|_2^2 \frac{\mathbf{w}_i}{\|\mathbf{w}_i\|_2} \frac{\mathbf{h}^{i\top}}{\|\mathbf{h}^i\|_2} = \left(\sigma_{i_j} - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right) \mathbf{u}_{i_j} \mathbf{v}_{i_j}^\top.$$

Next, we discuss global minimizers and global function values in two cases: (i) $d \geq K$, and (ii) $d < K$. For both cases, based on the above results, we can write

$$\begin{aligned} \mathbf{W} \mathbf{H}^\top &= \sum_{i=1}^d \mathbf{w}_i \mathbf{h}^{i\top} = \sum_{\mathbf{w}_i \neq \mathbf{0}, \mathbf{h}^i \neq \mathbf{0}} \left(\sigma_{i_j} - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right) \mathbf{u}_{i_j} \mathbf{v}_{i_j}^\top + \sum_{\mathbf{w}_i = \mathbf{0} \text{ and } \mathbf{h}^i = \mathbf{0}} \mathbf{w}_i \mathbf{h}^{i\top} \\ &= \sum_{\mathbf{w}_i \neq \mathbf{0}, \mathbf{h}^i \neq \mathbf{0}} \left(\sigma_{i_j} - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right) \mathbf{u}_{i_j} \mathbf{v}_{i_j}^\top. \end{aligned}$$

Case I: $d \geq K$. In this case, given the rank of \mathbf{W} is at most K , we know that the minimum is achieved when

$$\mathbf{W}^* \mathbf{H}^* = \sum_{i=1}^K \left[\sigma_i - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right]_+ \mathbf{u}_i \mathbf{v}_i^\top$$

with $\sigma_i \geq \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}$ for all $i = 1, \dots, K$. In this case, we have

$$\begin{aligned}\xi_{\star} &= \frac{1}{2} \sum_{i=1}^K \left(\left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+ - \sigma_i \right)^2 + \frac{\lambda_{\mathbf{W}}}{2} \sum_{i=1}^d \|\mathbf{w}_i\|_2^2 + \frac{\lambda_{\mathbf{H}}}{2} \sum_{i=1}^d \|\mathbf{h}^i\|_2^2 \\ &= \frac{1}{2} \sum_{i=1}^K \left(\left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+ - \sigma_i \right)^2 + \lambda_{\mathbf{W}} \sum_{i=1}^d \|\mathbf{w}_i\|_2^2 \\ &= \frac{1}{2} \sum_{i=1}^K \left(\left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+ - \sigma_i \right)^2 + \sqrt{\lambda_{\mathbf{H}}\lambda_{\mathbf{W}}} \sum_{i=1}^K \left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+, \end{aligned}$$

where for the second and third equality, we used (20) and (22), respectively.

Case II: $d < K$. In this case, we know that the minimum is achieved when

$$\mathbf{W}^* \mathbf{H}^* = \sum_{i=1}^d \left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+ \mathbf{u}_i \mathbf{v}_i^\top$$

with $\sigma_i \geq \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}$ for all $i = 1, \dots, d$. Similarly, we have

$$\xi_{\star} = \frac{1}{2} \sum_{i=1}^d \left(\left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+ - \sigma_i \right)^2 + \sqrt{\lambda_{\mathbf{H}}\lambda_{\mathbf{W}}} \sum_{i=1}^d \left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+ + \sum_{i=d+1}^K \sigma_i^2,$$

where the extra term $\sum_{i=d+1}^K \sigma_i^2$ is coming from the singular values of $\hat{\mathbf{Y}}$ and the decomposition of $\frac{1}{2} \left\| \mathbf{W} \mathbf{H} - \tilde{\mathbf{Y}} \right\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2$.

In summary, the minimum function value is obtained when

$$\mathbf{W}^* \mathbf{H}^* = \sum_{i=1}^{\min\{d, K\}} \left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+ \mathbf{u}_i \mathbf{v}_i^\top = \sum_{i=1}^{\min\{d, K\}} \eta_i \mathbf{u}_i \mathbf{v}_i^\top, \quad (23)$$

with $\eta_i(\lambda_{\mathbf{W}}, \lambda_{\mathbf{H}}) := \left[\sigma_i - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}} \right]_+$, and the minimum function value is attained as in (15).

Showing negative curvature for strict saddles. In the remaining part, we show those critical point (\mathbf{W}, \mathbf{H}) that does not satisfy the condition in (23) are strict saddle points, by showing that the Hessian of (14) has negative eigenvalues. First, we derive the directional Hessian of (14), which has the following form

$$\begin{aligned} 2[\nabla^2 \xi(\mathbf{W}, \mathbf{H})](\Delta, \Delta) &= \|\Delta_{\mathbf{W}} \mathbf{H} + \mathbf{W} \Delta_{\mathbf{H}}\|_F^2 + 2 \left\langle \mathbf{W} \mathbf{H} - \tilde{\mathbf{Y}}, \Delta_{\mathbf{W}} \Delta_{\mathbf{H}} \right\rangle \\ &\quad + \lambda_{\mathbf{W}} \|\Delta_{\mathbf{W}}\|_F^2 + \lambda_{\mathbf{H}} \|\Delta_{\mathbf{H}}\|_F^2. \end{aligned} \quad (24)$$

Given that a critical point (\mathbf{W}, \mathbf{H}) is not a global minimizer, then (23) is not satisfied. This implies that there must exist a singular value of $\hat{\mathbf{Y}}$ with $\sigma_j > \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}$, which cannot be covered by any $(\mathbf{w}_i, \mathbf{h}_i)$ in the sense that $\mathbf{w}_j \mathbf{h}_j^\top \neq (\sigma_j - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}) \mathbf{u}_j \mathbf{v}_j^\top$ for some j . We now discuss this situation separately in two cases: (i) $d \geq K$, and (ii) $d < K$.

Case I: $d \geq K$. In this case, since each column of \mathbf{W} is either zero or corresponds to the left singular vectors of $\tilde{\mathbf{Y}}$, it implies that the column space of \mathbf{W} has a non-trivial null space, i.e., there must exist a unit vector $\boldsymbol{\alpha} \in \mathbb{R}^d$ such that $\mathbf{W} \boldsymbol{\alpha} = \mathbf{0}$. Since $\lambda_{\mathbf{W}} \mathbf{W} \mathbf{W}^\top \mathbf{W} = \lambda_{\mathbf{H}} \mathbf{H} \mathbf{H}^\top$, we also have $\boldsymbol{\alpha}^\top \mathbf{H} = \mathbf{0}$. With this property, for the index j with $\mathbf{w}_j \mathbf{h}_j^\top \neq (\sigma_j - \sqrt{\lambda_{\mathbf{W}}\lambda_{\mathbf{H}}}) \mathbf{u}_j \mathbf{v}_j^\top$, we construct

$\Delta_{\mathbf{W}} = \left(\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}\right)^{1/4} \mathbf{u}_j \boldsymbol{\alpha}^\top$, $\Delta_{\mathbf{H}} = \left(\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}\right)^{1/4} \boldsymbol{\alpha} \mathbf{v}_j^\top$. Given that $\Delta_{\mathbf{W}} \mathbf{H} = \mathbf{0}$ and $\mathbf{W} \Delta_{\mathbf{H}} = \mathbf{0}$

$$\begin{aligned} \|\Delta_{\mathbf{W}} \mathbf{H} + \mathbf{W} \Delta_{\mathbf{H}}\|_F^2 &= 0, \\ \langle \mathbf{W} \mathbf{H} - \tilde{\mathbf{Y}}, \Delta_{\mathbf{W}} \Delta_{\mathbf{H}} \rangle &= -\sigma_j, \\ \lambda_{\mathbf{W}} \|\Delta_{\mathbf{W}}\|_F^2 + \lambda_{\mathbf{H}} \|\Delta_{\mathbf{H}}\|_F^2 &= 2\sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}. \end{aligned}$$

Plugging this into the Hessian (24), it gives

$$2[\nabla^2 \xi(\mathbf{W}, \mathbf{H})](\Delta, \Delta) = -2\sigma_j + 2\sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} = -2(\sigma_j - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}) < 0.$$

This implies that there exists a negative curvature for the Hessian, and the saddle point must be strict saddle.

Case II: $d < K$. Recall from (18) and (22) that $\sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \mathbf{W}^\top \mathbf{W}$ is a diagonal matrix with the values of diagonal entry from $\left\{[\sigma_1 - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}]_+, \dots, [\sigma_K - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}]_+, 0\right\}$, but here it excludes $[\sigma_j - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}]_+$ which equals $\sigma_j - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$ by our assumption. Thus, $\sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \mathbf{W}^\top \mathbf{W}$ has at least one diagonal entry which is strictly smaller than $\sigma_j - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$. Now let $\boldsymbol{\alpha} \in \mathbb{R}^d$ be the eigenvector associated with the smallest eigenvalue of $\mathbf{W}^\top \mathbf{W}$, so that

$$\nu := \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \boldsymbol{\alpha}^\top \mathbf{W}^\top \mathbf{W} \boldsymbol{\alpha} < \sigma_j - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}.$$

Since $\lambda_{\mathbf{W}} \mathbf{W}^\top \mathbf{W} = \lambda_{\mathbf{H}} \mathbf{H} \mathbf{H}^\top$, we also have $\sqrt{\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}} \boldsymbol{\alpha}^\top \mathbf{H}^\top \mathbf{H} \boldsymbol{\alpha} = \nu$. With this property, we construct $\Delta_{\mathbf{W}} = \left(\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}\right)^{1/4} \mathbf{u}_j \boldsymbol{\alpha}^\top$, $\Delta_{\mathbf{H}} = \left(\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}\right)^{1/4} \boldsymbol{\alpha} \mathbf{v}_j^\top$, which satisfies

$$\begin{aligned} \|\Delta_{\mathbf{W}} \mathbf{H} + \mathbf{W} \Delta_{\mathbf{H}}\|_F^2 &= \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \boldsymbol{\alpha}^\top \mathbf{W}^\top \mathbf{W} \boldsymbol{\alpha} + \sqrt{\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}} \boldsymbol{\alpha}^\top \mathbf{H}^\top \mathbf{H} \boldsymbol{\alpha} = 2\nu, \\ \langle \mathbf{W} \mathbf{H} - \tilde{\mathbf{Y}}, \Delta_{\mathbf{W}} \Delta_{\mathbf{H}} \rangle &= -\sigma_j, \\ \lambda_{\mathbf{W}} \|\Delta_{\mathbf{W}}\|_F^2 + \lambda_{\mathbf{H}} \|\Delta_{\mathbf{H}}\|_F^2 &= 2\sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}. \end{aligned}$$

Plugging this into the Hessian quadratic form gives

$$2[\nabla^2 \xi(\mathbf{W}, \mathbf{H})](\Delta, \Delta) = 2\nu - 2\sigma_j + 2\sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} = -2(\sigma_j - \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} - \nu) < 0.$$

Therefore, we prove (\mathbf{W}, \mathbf{H}) is a strict saddle for both cases. This completes the proof. \blacksquare

Lemma B.3 *Assume the number of training samples in each class is balanced, i.e., $n = n_1 = \dots = n_K$, and let $\mathbf{Y} = [\mathbf{y}_1 \ \dots \ \mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_K] \in \mathbb{R}^{K \times nK}$ be the matrix that contains the one-hot vectors for all the training samples. Then $\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbf{b} \mathbf{1}^\top$ has at least $K - 2$ singular values being \sqrt{n} . The rest of the two singular values, without loss of generality, denoted by σ_1 and σ_K , depend on \mathbf{b} . Then, we have the following lower bounds for σ_1 and σ_K .*

1. For any \mathbf{b} , the largest singular value σ_1 can be lower bounded by

$$\sigma_1 \geq \sqrt{n} \max \left(\sqrt{1 + K \left(\|\mathbf{b}\|^2 - \frac{1}{K} (\mathbf{1}^\top \mathbf{b})^2 \right)}, \left| 1 - \mathbf{1}^\top \mathbf{b} \right| \right). \quad (25)$$

2. For any \mathbf{b} on the sphere $\{\mathbf{b} \in \mathbb{R}^K : \|\mathbf{b}\| = c\}$ with $c \leq \frac{1}{\sqrt{K}}$, we have

$$\sigma_1 \geq \sqrt{n}, \quad \sigma_K \geq \sqrt{n} \left(1 - \sqrt{K} c \right) \quad (26)$$

and both inequalities become equalities if and only if $\mathbf{b} = \frac{c}{\sqrt{K}}\mathbf{1}$.

Proof [Proof of Theorem B.3] To study the singular values of \mathbf{Y} , it is equivalent to look at the eigenvalues of the Gram matrix of $\tilde{\mathbf{Y}}$:

$$\mathbf{G} = \tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top = (\mathbf{Y} - \mathbf{b}\mathbf{1}^\top) (\mathbf{Y} - \mathbf{b}\mathbf{1}^\top)^\top = n(\mathbf{I} - \mathbf{b}\mathbf{1}^\top - \mathbf{1}\mathbf{b}^\top + K\mathbf{b}\mathbf{b}^\top).$$

If \mathbf{b} is aligned with $\mathbf{1}$, i.e., they live in the same line, then $-\mathbf{b}\mathbf{1}^\top - \mathbf{1}\mathbf{b}^\top + K\mathbf{b}\mathbf{b}^\top$ is a rank-1 matrix and \mathbf{G} has $K - 1$ eigenvalues being n and the rest eigenvalue being $n(1 - \mathbf{1}^\top\mathbf{b})^2$. On the other hand, if \mathbf{b} is not aligned with $\mathbf{1}$, then $-\mathbf{b}\mathbf{1}^\top - \mathbf{1}\mathbf{b}^\top + K\mathbf{b}\mathbf{b}^\top$ is a rank-2 matrix and \mathbf{G} has $K - 2$ eigenvalues being n . In this case, the rest of the two eigenvalues, denoted by π_1 and π_K , correspond to the eigenvectors within the subspace spanned by $\mathbf{1}$ and \mathbf{b} .

To estimate the largest eigenvalues π_1 , we construct two orthonormal vectors within this subspace spanned by $\mathbf{1}$ and \mathbf{b} and compute the corresponding Rayleigh quotient. Specifically, we first compute the Rayleigh quotient along the direction $\mathbf{1}$ as

$$\frac{\mathbf{1}^\top \mathbf{G} \mathbf{1}}{\mathbf{1}^\top \mathbf{1}} = \frac{n}{K} \left(K - 2K\mathbf{1}^\top\mathbf{b} + K(\mathbf{1}^\top\mathbf{b})^2 \right) = n \left(1 - 2\mathbf{1}^\top\mathbf{b} + (\mathbf{1}^\top\mathbf{b})^2 \right) = n(1 - \mathbf{1}^\top\mathbf{b})^2.$$

Use Gram-Schmidt orthonormalization to obtain the other direction as $\mathbf{a} = \mathbf{b} - \frac{1}{K}\mathbf{1}^\top\mathbf{b}\mathbf{1}$, which gives the following Rayleigh quotient:

$$\frac{\mathbf{a}^\top \mathbf{G} \mathbf{a}}{\mathbf{a}^\top \mathbf{a}} = \frac{n}{\|\mathbf{a}\|^2} \left(\|\mathbf{a}\|^2 + K \left(\|\mathbf{b}\|^2 - \frac{1}{K} (\mathbf{1}^\top\mathbf{b})^2 \right)^2 \right) = n + nK \left(\|\mathbf{b}\|^2 - \frac{1}{K} (\mathbf{1}^\top\mathbf{b})^2 \right),$$

where the last equality follows because $\|\mathbf{a}\|^2 = \|\mathbf{b}\|^2 - \frac{1}{K} (\mathbf{1}^\top\mathbf{b})^2$. Thus, by the min-max theorem (i.e., Courant–Fischer–Weyl min-max principle), we have

$$\pi_1 \geq \max \left(\frac{\mathbf{1}^\top \mathbf{G} \mathbf{1}}{\mathbf{1}^\top \mathbf{1}}, \frac{\mathbf{a}^\top \mathbf{G} \mathbf{a}}{\mathbf{a}^\top \mathbf{a}} \right) \geq \max \left(n(1 - \mathbf{1}^\top\mathbf{b})^2, n + nK \left(\|\mathbf{b}\|^2 - \frac{1}{K} (\mathbf{1}^\top\mathbf{b})^2 \right) \right) \geq n,$$

where the last inequality becomes an inequality *if and only if* \mathbf{b} is a scaled version of the vector $\mathbf{1}$, i.e., $\mathbf{b} = \frac{\|\mathbf{b}\|}{\sqrt{K}}\mathbf{1}$.

To obtain a lower bound for π_K whenever $\|\mathbf{b}\| \leq \frac{1}{\sqrt{K}}$, we again use the the min-max theorem as

$$\begin{aligned} \frac{1}{n}\pi_K &\geq \min_{\|\mathbf{u}\|=1} \frac{1}{n} \mathbf{u}^\top \mathbf{G} \mathbf{u} = \min_{\|\mathbf{u}\|=1} 1 - 2\mathbf{u}^\top\mathbf{b}\mathbf{1}^\top\mathbf{u} + K(\mathbf{u}^\top\mathbf{b})^2 \\ &\geq \min_{\|\mathbf{u}\|=1} 1 - 2\sqrt{K} \left| \mathbf{u}^\top\mathbf{b} \right| + K(\mathbf{u}^\top\mathbf{b})^2 \\ &\geq 1 - 2\sqrt{K} \|\mathbf{b}\| + K \|\mathbf{b}\|^2 = \left(1 - \sqrt{K} \|\mathbf{b}\| \right)^2, \end{aligned}$$

where the first inequality achieves equality when \mathbf{u} is restricted to the subspace spanned by $\mathbf{1}$ and \mathbf{b} , the second inequality becomes an equality only when $\mathbf{u} = \mathbf{1}/\sqrt{K}$ and $\mathbf{u}^\top\mathbf{b} \geq 0$ or $\mathbf{u} = -\mathbf{1}/\sqrt{K}$ and $\mathbf{u}^\top\mathbf{b} \leq 0$, and the last inequality achieves equality if and only if \mathbf{u} is aligned with \mathbf{b} , i.e., $|\mathbf{u}^\top\mathbf{b}| = \|\mathbf{b}\|$. Thus, for any \mathbf{b} on the sphere $\{\mathbf{b} \in \mathbb{R}^K : \|\mathbf{b}\| = c\}$ with $c \leq \frac{1}{\sqrt{K}}$, π_K achieves its minimum possible value $n \left(1 - \sqrt{K} \|\mathbf{b}\| \right)^2$ if and only if $\mathbf{b} = \pm \frac{c}{\sqrt{K}}\mathbf{1}$. This completes the proof. \blacksquare

Lemma B.4 Assume the number of training samples in each class is balanced, i.e., $n = n_1 = \dots = n_K$, and let $\mathbf{Y} = [\mathbf{y}_1 \ \dots \ \mathbf{y}_1 \ \mathbf{y}_2 \ \dots \ \mathbf{y}_K] \in \mathbb{R}^{K \times nK}$ be the matrix that contains the one-hot

vectors for all the training samples. Suppose $b^* \leq \frac{1}{K}$. Then any global minimizer $(\mathbf{W}^*, \mathbf{H}^*)$ of

$$\min_{\mathbf{W}, \mathbf{H}} \frac{1}{2N} \left\| \mathbf{W}\mathbf{H} + b^* \mathbf{1}\mathbf{1}^\top - \mathbf{Y} \right\|_F^2 + \frac{\lambda_{\mathbf{W}}}{2} \|\mathbf{W}\|_F^2 + \frac{\lambda_{\mathbf{H}}}{2} \|\mathbf{H}\|_F^2. \quad (27)$$

satisfies the self-duality

$$\mathbf{h}_{k,i}^* = \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}} n}} \mathbf{w}^{*k}, \quad \forall k \in [K], i \in [n].$$

Moreover, if $\lambda_{\mathbf{W}} \lambda_{\mathbf{H}} \geq \frac{1}{NK}$, then $\mathbf{W}^* = \mathbf{0}$ and $\mathbf{H}^* = \mathbf{0}$. On the other hand, if $\lambda_{\mathbf{W}} \lambda_{\mathbf{H}} < \frac{1}{NK}$, $(\mathbf{W}^*, \mathbf{H}^*)$ further obeys the following properties for different d :

1. $d < K - 1$: $\mathbf{W}^* \mathbf{W}^{*\top} \sim \mathcal{P}_d(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top)$ where \mathcal{P}_d denotes the best rank- d approximating and $\mathbf{A} \sim \mathbf{B}$ means that there is a constant c such that $\mathbf{A} = c\mathbf{B}$;
2. $d = K - 1$: In this case, $\mathbf{W}^* \mathbf{W}^{*\top} \sim \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top$;
3. $d \geq K$ and $b^* \geq \frac{1}{K} - \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$: $\mathbf{W}^* \mathbf{W}^{*\top} \sim \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top$;
4. $d \geq K$ and $b^* < \frac{1}{K} - \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$: $\mathbf{W}^* \mathbf{W}^{*\top} \sim \mathbf{I} - \frac{b^*}{1 - K \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}} \mathbf{1}_K \mathbf{1}_K^\top$;

Proof [Proof of Theorem B.4] For convenience, let $\mathbf{1}_{K \times L}$ represents an all-ones matrix of size $K \times L$. Since $\mathbf{Y} - b^* \mathbf{1}_{K \times nK}$ contains many repeated columns, we first consider $\bar{\mathbf{Y}} = \mathbf{I}_K - b^* \mathbf{1}_{K \times K}$ that contains the non-repeated columns of $\mathbf{Y} - b^* \mathbf{1}_{K \times nK}$. Let $\bar{\mathbf{Y}} = \mathbf{U} \bar{\boldsymbol{\Sigma}} \mathbf{U}^\top$ be the eigenvalue decomposition, where $\mathbf{U} \in \mathbb{R}^{K \times K}$ is an orthonormal matrix and $\bar{\boldsymbol{\Sigma}} \in \mathbb{R}^{K \times K}$ is a diagonal matrix with eigenvalues $\bar{\sigma}_1 \geq \dots \geq \bar{\sigma}_K$ along the diagonals. Since $b^* \leq \frac{1}{K}$, the eigenvalues are given by $\bar{\sigma}_1 = \dots = \bar{\sigma}_{K-1} = 1 \geq \bar{\sigma}_K = 1 - b^* K$, and the eigenvector corresponding to $\bar{\sigma}_K$ is $\mathbf{u}_K = \frac{1}{\sqrt{K}} \mathbf{1}$, which implies that $[\mathbf{U}]_{K-1} [\mathbf{U}]_{K-1}^\top = \mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top$, where $[\mathbf{U}]_r$ means a $K \times r$ submatrix of \mathbf{U} by taking the first r columns.

Let $\boldsymbol{\Sigma} = \sqrt{n} \bar{\boldsymbol{\Sigma}}$ and $\mathbf{V}^\top = \frac{1}{\sqrt{n}} [\mathbf{u}^1 \ \dots \ \mathbf{u}^1 \ \mathbf{u}^2 \ \dots \ \mathbf{u}^K] \in \mathbb{R}^{K \times nK}$ that repeats the rescaled version of the column of \mathbf{U} n times so that $\mathbf{V}^\top \mathbf{V} = \mathbf{U}^\top \mathbf{U} = \mathbf{I}$. By noting the relation between $\mathbf{Y} - b^* \mathbf{1}_{K \times nK}$ and $\bar{\mathbf{Y}}$, we know $\mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$ is the SVD of $\mathbf{Y} - b^* \mathbf{1}_{K \times nK}$. When $\lambda_{\mathbf{W}} \lambda_{\mathbf{H}} \geq \frac{1}{NK}$, by applying Theorem B.2 and Theorem B.3, we conclude that $\mathbf{W}^* = \mathbf{0}$ and $\mathbf{H}^* = \mathbf{0}$ since $\sqrt{n} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \leq 0$. We now assume $\lambda_{\mathbf{W}} \lambda_{\mathbf{H}} < \frac{1}{NK}$ and utilize Theorem B.2 and Theorem B.3 again for the following cases:

1. $d < K - 1$: In this case, we have

$$\begin{aligned} \mathbf{W}^* &= \sqrt{\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}} \left(\sqrt{n} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right)^{1/2} \mathbf{U}(:, 1:d) \mathbf{R}, \\ \mathbf{H}^* &= \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \left(\sqrt{n} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right)^{1/2} \mathbf{R}^\top \mathbf{V}(:, 1:d)^\top, \forall \mathbf{R} \in \mathbb{R}^{d \times d}, \mathbf{R}^\top \mathbf{R} = \mathbf{I}. \end{aligned}$$

Thus, $\mathbf{h}_{k,i}^* = \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}} n}} \mathbf{w}^{*k}$ and $\mathbf{W}^* \mathbf{W}^{*\top} \sim \mathbf{U}(:, 1:d) \mathbf{U}(:, 1:d)^\top = \mathcal{P}_d(\mathbf{I} - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top)$.

2. $d = K - 1$: In this case, we have

$$\begin{aligned} \mathbf{W}^* &= \sqrt{\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}} \left(\sqrt{n} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right)^{1/2} \mathbf{U}(:, 1:K-1) \mathbf{R}, \\ \mathbf{H}^* &= \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \left(\sqrt{n} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right)^{1/2} \mathbf{R}^\top \mathbf{V}(:, 1:K-1)^\top, \forall \mathbf{R} \in \mathbb{R}^{(K-1) \times (K-1)}, \mathbf{R}^\top \mathbf{R} = \mathbf{I}. \end{aligned}$$

Thus, $\mathbf{h}_{k,i}^* = \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}n}} \mathbf{w}^{*k}$ and $\mathbf{W}^* \mathbf{W}^{*\top} \sim [\mathbf{U}]_{K-1} [\mathbf{U}]_{K-1}^\top = \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top$.

3. $d = K$: In this case, we have

$$\begin{aligned} \mathbf{W}^* &= \sqrt{\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}} \mathbf{U} \left[\boldsymbol{\Sigma} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right]_+^{1/2} \mathbf{R}, \\ \mathbf{H}^* &= \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \mathbf{R}^\top \left[\boldsymbol{\Sigma} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right]_+^{1/2} \mathbf{V}^\top, \quad \forall \mathbf{R} \in \mathbb{R}^{K \times K}, \mathbf{R}^\top \mathbf{R} = \mathbf{I} \end{aligned}$$

Thus, $\mathbf{h}_{k,i}^* = \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}n}} \mathbf{w}^{*k}$. Moreover, if $N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \geq \sqrt{n}(1 - b^* K)$, i.e., $b^* \geq \frac{1}{K} - \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$, then $\mathbf{W}^* \mathbf{W}^{*\top} \sim [\mathbf{U}]_{K-1} [\mathbf{U}]_{K-1}^\top = \mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top$. On the other hand, if $b^* < \frac{1}{K} - \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}$, then

$$\begin{aligned} \mathbf{W}^* \mathbf{W}^{*\top} &\sim \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^\top - K \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \mathbf{U} \mathbf{U}^\top = \bar{\sigma}_K = \mathbf{I} - b^* \mathbf{1}_K \mathbf{1}_K^\top - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \mathbf{I} \\ &= (1 - K \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}) \mathbf{I} - b^* \mathbf{1}_K \mathbf{1}_K^\top \sim \mathbf{I} - \frac{b^*}{1 - K \sqrt{n \lambda_{\mathbf{W}} \lambda_{\mathbf{H}}}} \mathbf{1}_K \mathbf{1}_K^\top. \end{aligned}$$

4. $d > K$: In this case, we have

$$\begin{aligned} \mathbf{W}^* &= \sqrt{\frac{\lambda_{\mathbf{H}}}{\lambda_{\mathbf{W}}}} \left[\mathbf{U} \left[\boldsymbol{\Sigma} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right]_+^{1/2} \quad \mathbf{0} \right] \mathbf{R}, \\ \mathbf{H}^* &= \sqrt{\frac{\lambda_{\mathbf{W}}}{\lambda_{\mathbf{H}}}} \mathbf{R}^\top \left[\begin{array}{c} \left[\boldsymbol{\Sigma} - N \sqrt{\lambda_{\mathbf{W}} \lambda_{\mathbf{H}}} \right]_+^{1/2} \mathbf{V}^\top \\ \mathbf{0} \end{array} \right], \quad \forall \mathbf{R} \in \mathbb{R}^{d \times d}, \mathbf{R}^\top \mathbf{R} = \mathbf{I}. \end{aligned}$$

One can verify that $(\mathbf{W}^*, \mathbf{H}^*)$ satisfies the same properties as in the case of $d = K$. ■

C Visualizations of Optimization Landscapes in Section 3.3

C.1 Details of the Visualization Technique

We provide the technical details on how the visualization in Section 3.3 is obtained.

The following result expresses the output of the classifier layer for a feature vector \mathbf{h} as a function of the norm of \mathbf{h} and its angle to a classifier weight vector \mathbf{w}^k .

Proposition C.1 *Given any $d \geq K - 1 > 1$, take the classifier weights \mathbf{W}, \mathbf{b} to be such that \mathbf{W} is an arbitrary K -Simplex ETF (see Definition .1) and $\mathbf{b} = \mathbf{0}$. Take any $k, k' \in \{1, \dots, K\}$, and consider a vector \mathbf{h} on the two-dimensional plane $\text{span}\{\mathbf{w}^k, \mathbf{w}^{k'}\}$ parameterized in the polar coordinate system with polar axis being \mathbf{w}^k . Denote s and θ the radial and angular (in radians) coordinates of \mathbf{h} , respectively (positive angular direction of the polar coordinate system is taken so that $\mathbf{w}^{k'}$'s angular coordinate is in $(0, \pi)$). We have*

- The feature \mathbf{h} can be expressed as a linear combination of \mathbf{w}^k and $\mathbf{w}^{k'}$:

$$\mathbf{h} = s \left(\frac{\sin \theta}{\sqrt{K^2 - 2K}} + \cos \theta \right) \mathbf{w}^k + s(K - 1) \frac{\sin \theta}{\sqrt{K^2 - 2K}} \mathbf{w}^{k'}; \quad (28)$$

- The output of the classifier layer (\mathbf{W}, \mathbf{b}) is given by

$$\langle \mathbf{w}^{k''}, \mathbf{h} \rangle + b_{k''} = \begin{cases} s \cos \theta, & \text{if } k'' = k; \\ s \frac{\sqrt{K^2 - 2K}}{K-1} \sin \theta - \frac{s}{K-1} \cos \theta, & \text{if } k'' = k'; \\ -s \sqrt{\frac{K}{K-2}} \frac{1}{K-1} \sin \theta - \frac{s}{K-1} \cos \theta, & \text{otherwise.} \end{cases} \quad (29)$$

Note that (29) is invariant to the arbitrary rotation in K -Simplex ETF.

We omit the proof to Proposition C.1 as it can be obtained via simple algebra.

Based on Proposition C.1, we can obtain the (rescaled) MSE and CE losses as a function of (s, θ) . Assuming that \mathbf{h} belongs to class k , the rescaled MSE loss defined in (4) w.r.t. \mathbf{h} is given by

$$\text{Loss}_{\text{MSE}}(\mathbf{h}; \alpha, M) = \frac{\alpha}{2} \left(\langle \mathbf{w}^k, \mathbf{h} \rangle + b_k - M \right)^2 + \frac{1}{2} \sum_{k'' \neq k} \left(\langle \mathbf{w}^{k''}, \mathbf{h} \rangle + b_{k''} - 1 \right)^2, \quad (30)$$

where α, M are rescaling parameters. Plugging in the results in (29), we obtain

$$\begin{aligned} \text{Loss}_{\text{MSE}}(s, \theta; \alpha, M) = & \frac{\alpha}{2} \cdot (s \cos \theta - M)^2 + \frac{s^2}{2} \cdot \left(\frac{\sqrt{K^2 - 2K} \sin \theta - \cos \theta}{K-1} \right)^2 \\ & + \frac{s^2}{2} \cdot (K-2) \cdot \left(\frac{\sqrt{\frac{K}{K-2}} \sin \theta + \cos \theta}{K-1} \right)^2. \end{aligned} \quad (31)$$

Similarly, we may obtain the CE loss as

$$\text{Loss}_{\text{CE}}(s, \theta) = -\log \left(\frac{e^{s \cos \theta}}{e^{s \cos \theta} + e^{s \frac{\sqrt{K^2 - 2K} \sin \theta - \cos \theta}{K-1}} + (K-2) e^{-s \frac{\sqrt{\frac{K}{K-2}} \sin \theta + \cos \theta}{K-1}}} \right). \quad (32)$$

Figure 2 is obtained by plotting the loss functions in (31) and (32).

C.2 Visualization of the Gradient Vector Field

We consider the regime of $K \rightarrow \infty$ in which the rescaled MSE loss (31) becomes

$$\lim_{K \rightarrow \infty} \text{Loss}_{\text{MSE}}(s, \theta; \alpha, M) = \frac{\alpha}{2} (s \cos \theta - M)^2 + \frac{1}{2} s^2 \sin^2 \theta. \quad (33)$$

Taking the derivative w.r.t. s and θ , we obtain

$$\begin{aligned} \frac{\partial}{\partial s} \lim_{K \rightarrow \infty} \text{Loss}_{\text{MSE}}(s, \theta; \alpha, M) &= s + (\alpha - 1) s \cos^2 \theta - \alpha M \cos \theta, \\ \frac{\partial}{\partial \theta} \lim_{K \rightarrow \infty} \text{Loss}_{\text{MSE}}(s, \theta; \alpha, M) &= \alpha M s \sin \theta - (\alpha - 1) s^2 \sin \theta \cos \theta. \end{aligned} \quad (34)$$

Similarly, we may obtain the gradient for CE as

$$\begin{aligned} \frac{\partial}{\partial s} \lim_{K \rightarrow \infty} \text{Loss}_{\text{CE}}(s, \theta; \alpha, M) &= \frac{e^{\sin \theta} (\sin \theta - \cos \theta)}{e^{\sin \theta} + e^{\cos \theta}}, \\ \frac{\partial}{\partial \theta} \lim_{K \rightarrow \infty} \text{Loss}_{\text{CE}}(s, \theta; \alpha, M) &= \frac{s e^{\sin \theta} (\sin \theta + \cos \theta)}{e^{\sin \theta} + e^{\cos \theta}}. \end{aligned} \quad (35)$$

In Figure 8, we visualize the gradient of MSE (in (34)) and CE (in (35)) losses by plotting their gradient vector fields. It shows that rescaling of the MSE loss by either increasing M or increasing α helps to align the gradient along the direction of minimizing θ . Recall that θ determines the classifier's prediction of the class membership for \mathbf{h} while s is irrelevant.

When restricting our attention to a feature \mathbf{h} with $\theta = \frac{\pi}{2}$, the gradient w.r.t. s and θ becomes

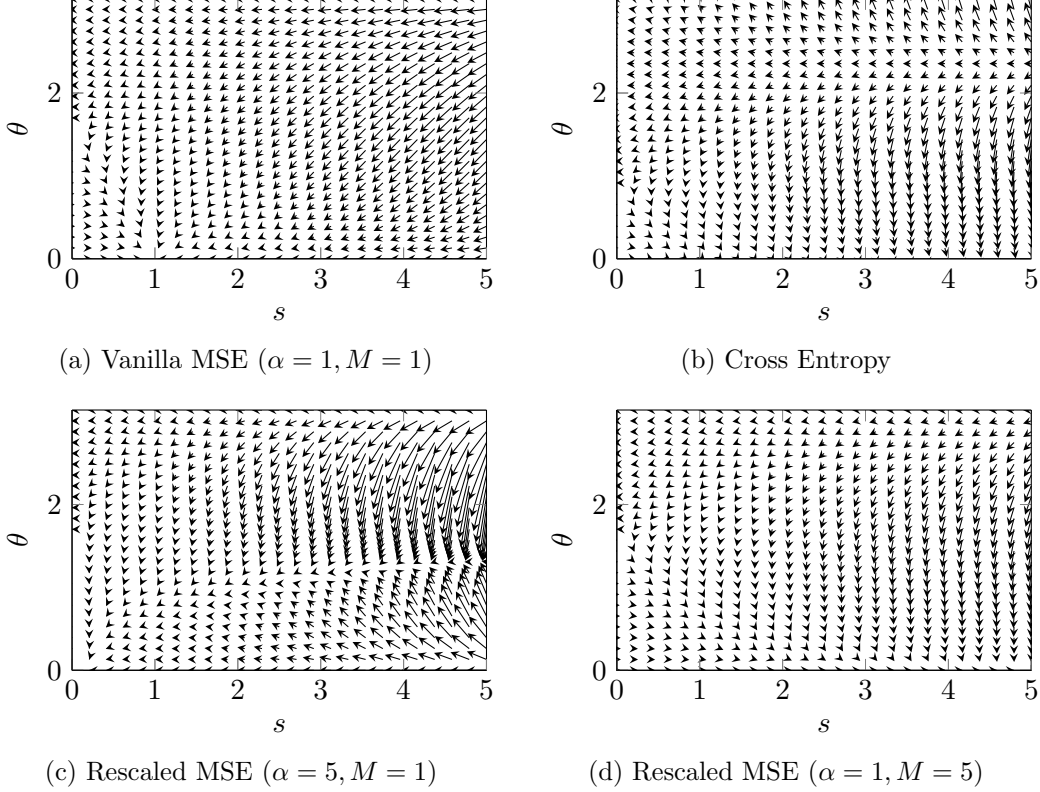


Figure 8: **Visualization of the gradient vector fields with different losses.** We fix \mathbf{W} as a simplex ETF and illustrate the landscape only w.r.t. a feature $\mathbf{h}_{k,i}$. For each plot, the s -axis denotes $\|\mathbf{h}_{k,i}\|_2$, and the θ -axis denotes the angle $\arccos(\langle \mathbf{h}_{k,i}, \mathbf{w}^k \rangle)$. The arrows point to gradient descent directions with length proportional to the gradient norm.

s and αMs , respectively. Here, increasing the rescaling parameters α or M in the range of $(1, \infty)$ has the effect of increasing the component of the gradient along the θ direction while keeping the component along the s direction fixed.