

UNDERSTANDING THE ROLE OF NONLINEARITY IN TRAINING DYNAMICS OF CONTRASTIVE LEARNING

Yuandong Tian

Meta AI (FAIR)

yuandong@meta.com

ABSTRACT

While the empirical success of self-supervised learning (SSL) heavily relies on the usage of deep nonlinear models, existing theoretical works on SSL understanding still focus on linear ones. In this paper, we study the role of nonlinearity in the training dynamics of contrastive learning (CL) on one and two-layer nonlinear networks with homogeneous activation $h(x) = h'(x)x$. We have two major theoretical discoveries. First, the presence of nonlinearity can lead to many local optima even in 1-layer setting, each corresponding to certain patterns from the data distribution, while with linear activation, only one major pattern can be learned. This suggests that models with lots of parameters can be regarded as a *brute-force* way to find these local optima induced by nonlinearity. Second, in the 2-layer case, linear activation is proven not capable of learning specialized weights into diverse patterns, demonstrating the importance of nonlinearity. In addition, for 2-layer setting, we also discover *global modulation*: those local patterns discriminative from the perspective of global-level patterns are prioritized to learn, further characterizing the learning process. Simulation verifies our theoretical findings.

1 INTRODUCTION

Over the last few years, deep models have demonstrated impressive empirical performance in many disciplines, not only in supervised but also in recent self-supervised setting (SSL), in which models are trained with a surrogate loss (e.g., predictive (Devlin et al., 2018; He et al., 2021), contrastive (Chen et al., 2020; Caron et al., 2020; He et al., 2020) or noncontrastive loss (Grill et al., 2020; Chen & He, 2020)) and its learned representation is then used for downstream tasks.

From the theoretical perspective, understanding the roles of nonlinearity in deep neural networks is one critical part of understanding how modern deep models work. Currently, most works focus on linear variants of deep models (Jacot et al., 2018; Arora et al., 2019a; Kawaguchi, 2016; Jing et al., 2022; Tian et al., 2021; Wang et al., 2021). When nonlinearity is involved, deep models are often treated as richer families of black-box functions than linear ones (Arora et al., 2019b; HaoChen et al., 2021). The role played by nonlinearity is also studied, mostly on model expressibility (Gühring et al., 2020; Raghu et al., 2017; Lu et al., 2017) in which specific weights are found to fit the complicated structure of the data well, regardless of the training algorithm. However, many questions remain open: if model capacity is the key, why traditional models like k -NN (Fix & Hodges, 1951) or kernel SVM (Cortes & Vapnik, 1995) do not achieve comparable empirical performance, even if theoretically they can also fit any functions (Hammer & Gersmann, 2003; Devroye et al., 1994). Moreover, while traditional ML theory suggests carefully controlling model capacity to avoid overfitting, large neural models often generalize well in practice (Brown et al., 2020; Chowdhery et al., 2022).

In this paper, we study the critical role of nonlinearity in the training dynamics of contrastive learning (CL). Specifically, by extending the recent α -CL framework (Tian, 2022) and linking it to kernels (Paulsen & Raghupathi, 2016), we show that even with 1-layer nonlinear networks, nonlinearity plays a critical role by creating many local optima. As a result, the more nonlinear nodes in 1-layer networks with different initialization, the more local optima are likely to be collected as learned patterns in the trained weights, and the richer the resulting representation becomes. Moreover, popular loss functions like InfoNCE tends to have more local optima than quadratic ones. In contrast, in the linear setting, contrastive learning becomes PCA under certain conditions (Tian, 2022), and

only the most salient pattern (i.e., the maximal eigenvector of the data covariance matrix) is learned while other less salient ones are lost, regardless of the number of hidden nodes.

Based on this finding, we extend our analysis to 2-layer ReLU setting with non-overlapping receptive fields. In this setting, we prove the fundamental limitation of linear networks: the gradients of multiple weights at the same receptive field are always co-linear, preventing diverse pattern learning.

Finally, we also characterize the interaction between layers in 2-layer network: while in each receptive field, many patterns exist, those contributing to global patterns are prioritized to learn by the training dynamics. This *global modulation* changes the eigenstructure of the low-level covariance matrix so that relevant patterns are learned with higher probability.

In summary, through the lens of training dynamics, we discover unique roles played by nonlinearity which linear activation cannot do: (1) nonlinearity creates many local optima for different patterns of the data, and (2) nonlinearity enables weight specialization to diverse patterns. In addition, we also discover a mechanism for how global pattern prioritizes which local patterns to learn, shedding light on the role played by network depth. Preliminary experiments on simulated data verify our findings.

Related works. Many works analyze network at initialization (Hayou et al., 2019; Roberts et al., 2021) and avoid the complicated training dynamics. Previous works (Wilson et al., 1997; Li & Yuan, 2017; Tian et al., 2019; Tian, 2017; Allen-Zhu & Li, 2020) that analyze training dynamics mostly focus on supervised learning. Different from Saunshi et al. (2022); Ji et al. (2021) that analyzes feature learning process in linear models of CL, we focus on the critical role played by nonlinearity. Our analysis is also more general than Li & Yuan (2017) that focuses on 1-layer ReLU network with symmetric weight structure trained on sparse linear models. Along the line of studying dynamics of contrastive learning, Jing et al. (2022) analyzes dimensional collapsing on 1 and 2 layer linear networks. Tian (2022) proves that such collapsing happens in linear networks of any depth and further analyze ReLU scenarios but with strong assumptions (e.g., one-hot positive input). Our work uses much more relaxed assumptions and performs in-depth analysis for homogeneous activations.

2 PROBLEM SETUP

Notation. In this section, we introduce our problem setup of contrastive learning. Let $\mathbf{x}_0 \sim p_{\mathcal{D}}(\cdot)$ be a sample drawn from the dataset, and $\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)$ be an augmentation view of the sample \mathbf{x}_0 . Here both \mathbf{x}_0 and \mathbf{x} are random variables. Let $\mathbf{f} = \mathbf{f}(\mathbf{x}; \boldsymbol{\theta})$ be the output of a deep neural network that maps input \mathbf{x} into some representation space with parameter $\boldsymbol{\theta}$ to be optimized. Given a batch of size N , $\mathbf{x}_0[i]$ represent i -th sample (i.e., instantiation) of corresponding random variables, and $\mathbf{x}[i]$ and $\mathbf{x}[i']$ are two of its augmented views. Here $\mathbf{x}[\cdot]$ has $2N$ samples, $1 \leq i \leq N$ and $N+1 \leq i' \leq 2N$.

Contrastive learning (CL) aims to learn the parameter $\boldsymbol{\theta}$ so that the representation \mathbf{f} are distinct from each other: we want to maximize squared distance $d_{ij}^2 := \|\mathbf{f}[i] - \mathbf{f}[j]\|_2^2/2$ between samples $i \neq j$ and minimize $d_i^2 := \|\mathbf{f}[i] - \mathbf{f}[i']\|_2^2/2$ between two views $\mathbf{x}[i]$ and $\mathbf{x}[i']$ from the same sample $\mathbf{x}_0[i]$.

Many objectives in contrastive learning have been proposed to combine these two goals into one. For example, InfoNCE (Oord et al., 2018) minimizes the following (here τ is the temperature):

$$\mathcal{L}_{nce} := -\tau \sum_{i=1}^N \log \frac{\exp(-d_i^2/\tau)}{\epsilon \exp(-d_i^2/\tau) + \sum_{j \neq i} \exp(-d_{ij}^2/\tau)} \quad (1)$$

In this paper, we follow α -CL (Tian, 2022) that proposes a general CL framework that covers a broad family of existing CL losses. α -CL maximizes an energy function $\mathcal{E}_\alpha(\boldsymbol{\theta})$ using gradient ascent:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \eta \nabla_{\boldsymbol{\theta}} \mathcal{E}_{\text{sg}(\alpha(\boldsymbol{\theta}_t))}(\boldsymbol{\theta}), \quad (2)$$

where η is the learning rate, $\text{sg}(\cdot)$ is the stop gradient operator, the *energy* function $\mathcal{E}_\alpha(\boldsymbol{\theta}) := \frac{1}{2} \text{tr} \mathbb{C}_\alpha[\mathbf{f}, \mathbf{f}]$ and $\mathbb{C}_\alpha[\cdot, \cdot]$ is the *contrastive covariance* (Tian, 2022; Jing et al., 2022)¹:

$$\mathbb{C}_\alpha[\mathbf{a}, \mathbf{b}] := \frac{1}{2N^2} \sum_{i,j=1}^N \alpha_{ij} [(\mathbf{a}[i] - \mathbf{a}[j])(\mathbf{b}[i] - \mathbf{b}[j])^\top - (\mathbf{a}[i] - \mathbf{a}[i'])(\mathbf{b}[i] - \mathbf{b}[i'])^\top] \quad (3)$$

¹Compared to Tian (2022), our \mathbb{C}_α definition has an additional constant term $1/2N^2$ to simply the notation.

One important quantity is the *pairwise importance* $\alpha(\boldsymbol{\theta}) = [\alpha_{ij}(\boldsymbol{\theta})]_{i,j=1}^N$, which are N^2 weights on pairwise pairs of N samples in a batch. Intuitively, these weights make the training focus more on *hard negative pairs*, i.e., distinctive sample pairs that are similar in the representation space but are supposed to be separated away. Many existing CL losses (InfoNCE, triplet loss, etc) are special cases of α -CL (Tian, 2022) by choosing different $\alpha(\boldsymbol{\theta})$, e.g., quadratic loss corresponds to $\alpha_{ij} := \text{const}$ and InfoNCE (with $\epsilon = 0$) corresponds to $\alpha_{ij} := \exp(-d_{ij}^2/\tau) / \sum_{j \neq i} \exp(-d_{ij}^2/\tau)$.

For brevity $\mathbb{C}_\alpha[\mathbf{x}] := \mathbb{C}_\alpha[\mathbf{x}, \mathbf{x}]$. For the energy function $\mathcal{E}_\alpha(\boldsymbol{\theta}) := \text{tr} \mathbb{C}_\alpha[\mathbf{f}(\mathbf{x}; \boldsymbol{\theta})]$, in this work we mainly study its landscape, i.e., existence of local optima, their local properties and overall distributions, where \mathbf{f} is a nonlinear network with parameters $\boldsymbol{\theta}$. Note that in Eqn. 2, the stop gradient operator $\text{sg}(\alpha)$ means that while the value of α may depend on $\boldsymbol{\theta}$, when studying the local property of $\boldsymbol{\theta}$, α makes no contribution to the gradient and should be treated as an independent variable.

Since \mathbb{C}_α is an abstract mathematical object with complicated definitions, as the first contribution, we give its connection to regular variance $\mathbb{V}[\cdot]$, if the pairwise importance α has certain *kernel structures* (Ghojogh et al., 2021; Paulsen & Raghupathi, 2016):

Definition 1 (Kernel structure of pairwise importance α). *There exists a (kernel) function $\mathcal{K}(\cdot, \cdot)$ so that $\alpha_{ij} = \mathcal{K}(\mathbf{x}_0[i], \mathbf{x}_0[j])$. Here \mathcal{K} satisfies the decomposition $\mathcal{K}(\mathbf{a}, \mathbf{b}) = \boldsymbol{\phi}^\top(\mathbf{a})\boldsymbol{\phi}(\mathbf{b}) = \sum_{l=0}^{+\infty} \phi_l(\mathbf{a})\phi_l(\mathbf{b})$ with non-negative high-dimensional mapping $\boldsymbol{\phi}(\cdot) = [\phi_l(\cdot)] \geq 0$.*

Definition 2 (Adjusted PDF $\tilde{p}_l(\mathbf{x})$). *For l -th component ϕ_l of the mapping $\boldsymbol{\phi}$, we define the adjusted density $\tilde{p}_l(\mathbf{x}; \alpha) := \frac{1}{z_l(\alpha)} \phi_l(\mathbf{x}; \alpha) p_{\mathbb{D}}(\mathbf{x})$, where $z_l(\alpha) := \int \phi_l(\mathbf{x}) p_{\mathbb{D}}(\mathbf{x}) d\mathbf{x} \geq 0$ is the normalizer.*

Obviously $\alpha_{ij} \equiv 1$ (uniform α corresponding to quadratic loss) satisfies Def. 1 with 1D mapping $\boldsymbol{\phi} \equiv 1$. Here we show a non-trivial case, *Gaussian* α , whose normalized version leads to InfoNCE:

Lemma 1 (Gaussian α). *For any function $\mathbf{g}(\cdot)$ that is bounded below, if we use $\alpha_{ij} := \exp(-\|\mathbf{g}(\mathbf{x}_0[i]) - \mathbf{g}(\mathbf{x}_0[j])\|_2^2/2\tau)$ as the pairwise importance, then it has kernel structure (Def. 1).*

Note that Gaussian α computes N^2 pairwise distances using *un-augmented* samples \mathbf{x}_0 , while InfoNCE (and most of CL losses) uses augmented views \mathbf{x} and \mathbf{x}' and normalizes along one dimension to yield asymmetric α_{ij} . Here Gaussian α is a convenient tool for analysis. We now show \mathbb{C}_α is a summation of regular variances but with different probability of data, adjusted by the pairwise importance α that has kernel structures. Please check Appendix A.1 for detailed proofs.

Lemma 2 (Relationship between Contrastive Covariance and Variance in large batch size). *If α satisfies Def. 1, then for any function $\mathbf{g}(\cdot)$, $\mathbb{C}_\alpha[\mathbf{g}(\mathbf{x})]$ is asymptotically PSD when $N \rightarrow +\infty$:*

$$\mathbb{C}_\alpha[\mathbf{g}(\mathbf{x})] \rightarrow \sum_l z_l^2 \mathbb{V}_{\mathbf{x}_0 \sim \tilde{p}_l(\cdot; \alpha)} [\mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot | \mathbf{x}_0)}[\mathbf{g}(\mathbf{x}) | \mathbf{x}_0]] \quad (4)$$

Corollary 1 (No augmentation and large batchsize). *With the condition of Lemma 2, if we further assume there is no augmentation (i.e., $p_{\text{aug}}(\mathbf{x} | \mathbf{x}_0) = \delta(\mathbf{x} - \mathbf{x}_0)$), then $\mathbb{C}_\alpha[\mathbf{g}] \rightarrow \sum_l z_l^2 \mathbb{V}_{\tilde{p}_l}[\mathbf{g}]$.*

3 ONE-LAYER CASE

Now let us first consider 1-layer network with K hidden nodes: $\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}) = h(W\mathbf{x})$, where $W = [\mathbf{w}_1, \dots, \mathbf{w}_K]^\top \in \mathbb{R}^{K \times d}$, $\boldsymbol{\theta} = \{W\}$ and $h(x)$ is the activation. The k -th row of W is a weight \mathbf{w}_k and its output is $f_k := h(\mathbf{w}_k^\top \mathbf{x})$. In this case, $\text{tr} \mathbb{C}_\alpha[\mathbf{f}] = \sum_{k=1}^K \mathbb{C}_\alpha[f_k]$. We consider per-filter normalization $\|\mathbf{w}_k\|_2 = 1$, which can be achieved by imposing BatchNorm (Ioffe & Szegedy, 2015) at each node k (Tian, 2022). In this case, optimization can be decoupled into each filter \mathbf{w}_k :

$$\max_{\boldsymbol{\theta}} \mathcal{E}_\alpha(\boldsymbol{\theta}) = \frac{1}{2} \max_{\|\mathbf{w}_k\|_2=1, 1 \leq k \leq K} \text{tr} \mathbb{C}_\alpha[\mathbf{f}] = \frac{1}{2} \sum_{k=1}^K \max_{\|\mathbf{w}_k\|_2=1} \mathbb{C}_\alpha[h(\mathbf{w}_k^\top \mathbf{x})] \quad (5)$$

Now let's think about, which parameters \mathbf{w}_k maximizes the summation? For the linear case, since $\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})] = \mathbb{C}_\alpha[\mathbf{w}^\top \mathbf{x}] = \mathbf{w}^\top \mathbb{C}_\alpha[\mathbf{x}] \mathbf{w}$, all \mathbf{w}_k converge to the maximal eigenvector of $\mathbb{C}_\alpha[\mathbf{x}]$ (a constant matrix), regardless of how they are initialized and what the distribution of \mathbf{x} is. Therefore, the linear case will only learn the most salient single pattern due to the (overly-smooth) landscape of the objective function, a winner-take-all effect that neglects many patterns in the data.

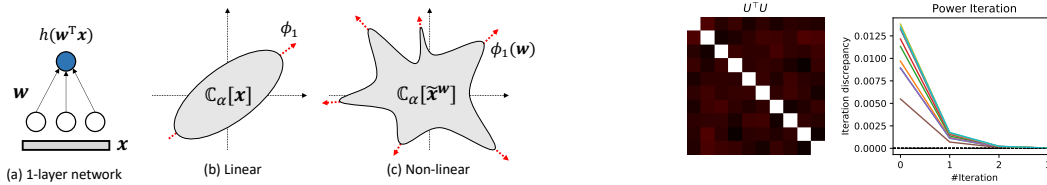


Figure 1: **Left:** Summary of Sec. 3. (a) We analyze the dynamics of one-layer network $h(\mathbf{w}^\top \mathbf{x})$ under CL loss (Eqn. 2). (b) With linear activation $h(x) = x$, then there is only one fixed point (PCA direction). (c) Non-linear activation $h(x)$ creates many critical points and a proper choice of pairwise importance α can make them local optima, enabling learning of diverse features. **Right:** Convergence patterns (iteration t versus iteration discrepancy $\|\mathbf{w}(t+1) - \mathbf{w}(t)\|_2$) of Power Iteration (Eqn. PI) in latent summation models, when $\|U^\top U - I\|_2$ is small but non-zero. In this case, Theorem 3 tells there still exist local optima close to each \mathbf{u}_m .

In contrast, nonlinearity can change the landscape and create more local optima in $\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})]$, each capturing one pattern. In this paper, we consider a general category of nonlinearity activations:

Assumption 1 (Homogeneity (Du et al., 2018)/Reversibility (Tian et al., 2020)). *The activation satisfies $h(x) = h'(x)x$.*

Many activations satisfy this assumption, including linear, ReLU, LeakyReLU and monomial activations like $h(x) = x^p$ (with an additional global constant). In this case we have:

$$h(\mathbf{w}^\top \mathbf{x}) = \mathbf{w}^\top h'(\mathbf{w}^\top \mathbf{x}) \mathbf{x} = \mathbf{w}^\top \tilde{\mathbf{x}}^{\mathbf{w}}, \quad (6)$$

where $\tilde{\mathbf{x}}^{\mathbf{w}} := \mathbf{x} \cdot h'(\mathbf{w}^\top \mathbf{x})$ is the input data after nonlinear gating. When there is no ambiguity, we just write $\tilde{\mathbf{x}}^{\mathbf{w}}$ as $\tilde{\mathbf{x}}$ and omit the weight superscript. One property is $\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})] = \mathbf{w}^\top \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] \mathbf{w}$.

Now let $A(\mathbf{w}) := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}]$. With the constraint $\|\mathbf{w}\|_2 = 1$, the learning dynamics is:

Lemma 3 (Training dynamics of 1-layer network with homogeneous activation in contrastive learning). *The gradient dynamics of Eqn. 5 is (note that α is treated as an independent variable):*

$$\dot{\mathbf{w}}_k = P_{\mathbf{w}_k}^\perp A(\mathbf{w}_k) \mathbf{w}_k \quad (7)$$

Here $P_{\mathbf{w}_k}^\perp := I - \mathbf{w}_k \mathbf{w}_k^\top$ projects a vector into the complementary subspace spanned by \mathbf{w}_k .

See Appendix B.2 for derivations. Now the question is that: what is the critical point of the dynamics and whether they are attractive (i.e., local optima). In linear case, the maximal eigenvector is the one fixed point; in nonlinear case, we are looking for *locally* maximal eigenvectors, called *LME*.

Definition 3 (Locally maximal eigenvector (LME)). *\mathbf{w}_* is a locally maximal eigenvector of $A(\mathbf{w})$, if $A(\mathbf{w}_*) \mathbf{w}_* = \lambda_* \mathbf{w}_*$, where $\lambda_* = \lambda_{\max}(A(\mathbf{w}_*))$ is the distinct maximal eigenvalue of $A(\mathbf{w}_*)$.*

It is easy to see each LME is a critical point of the dynamics, since $P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*) \mathbf{w}_* = \lambda P_{\mathbf{w}_*}^\perp \mathbf{w}_* = 0$.

3.1 EXAMPLES WITH MULTIPLE LMEs IN RELU SETTING

To see why the nonlinear activation leads to many LMEs in Eqn. 7, we first give two exemplar generative models of the input \mathbf{x} that show Eqn. 7 has multiple critical points, then introduce more general cases. To make the examples simple and clear, we assume the condition of Corollary 1 (no augmentation and large batchsize), and let $\alpha_{ij} \equiv 1$. Notice that $\tilde{\mathbf{x}}^{\mathbf{w}}$ is a deterministic function of \mathbf{x} , therefore $A(\mathbf{w}) := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}]$. We also use ReLU activation $h(x) = \max(x, 0)$.

Let $U = [\mathbf{u}_1, \dots, \mathbf{u}_M]$ be orthonormal bases ($\mathbf{u}_m^\top \mathbf{u}_{m'} = \mathbb{I}(m = m')$). Here are two examples:

Latent categorical model. Suppose y is a categorical random variable taking M possible values, $\mathbb{P}[\mathbf{x}|y = m] = \delta(\mathbf{x} - \mathbf{u}_m)$. Then we have (see Appendix B.1 for detailed steps):

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{P}[y = m](1 - \mathbb{P}[y = m]) \mathbf{u}_m \mathbf{u}_m^\top \quad (8)$$

Now it is clear that $\mathbf{w} = \mathbf{u}_m$ is an LME for any m .

Latent summation model. Suppose there is a latent variable \mathbf{y} so that $\mathbf{x} = U\mathbf{y}$, where $\mathbf{y} := [y_1, y_2, \dots, y_M]$. Each y_m is a standardized Bernoulli random variable: $\mathbb{E}[y_m] = 0$

and $\mathbb{E}[y_m^2] = 1$. This means that $y_m = y_m^+ := \sqrt{(1 - q_m)/q_m}$ with probability q_m and $y_m = y_m^- := -\sqrt{q_m/(1 - q_m)}$ with probability $1 - q_m$. For $m_1 \neq m_2$, y_{m_1} and y_{m_2} are independent. Then we have:

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}] = (1 - q_m)^2 \mathbf{u}_m \mathbf{u}_m^\top + q_m (I - \mathbf{u}_m \mathbf{u}_m^\top) \quad (9)$$

which has a maximal and distinct eigenvector of \mathbf{u}_m with a unique eigenvalue $(1 - q_m)^2$, when $q_m < \frac{1}{2}(3 - \sqrt{5}) \approx 0.382$. Therefore, different \mathbf{w} leads to different LMEs.

In both cases, the presence of ReLU removes the ‘‘redundant energy’’ so that $A(\mathbf{w})$ can focus on specific directions, creating multiple LMEs that correspond to multiple learnable patterns. The two examples can be computed analytically due to our specific choices on nonlinearity h and α .

3.2 RELATE LMEs TO LOCAL OPTIMA

Once LMEs are identified, the next step is to check whether they are attractive, or *stable* critical points, or *local optima*. That is, whether the weights converge into them and stay there during training. For this, some notations are introduced below.

Notations. Let $\lambda_i(\mathbf{w})$ be the i -th largest eigenvalue of $A(\mathbf{w})$, and $\phi_i(\mathbf{w})$ the corresponding unit eigenvector, $\lambda_{\text{gap}}(\mathbf{w}) := \lambda_1(\mathbf{w}) - \lambda_2(\mathbf{w})$ the eigenvalue gap. Let $\rho(\mathbf{w})$ be the *local roughness measure*: $\rho(\mathbf{w})$ is the smallest scalar to satisfy $\|(A(\mathbf{v}) - A(\mathbf{w}))\mathbf{w}\|_2 \leq \rho(\mathbf{w})\|\mathbf{v} - \mathbf{w}\|_2 + \mathcal{O}(\|\mathbf{v} - \mathbf{w}\|_2^2)$ in a local neighborhood of \mathbf{w} . The following theorem gives a sufficient condition for stability of \mathbf{w}_* :

Theorem 1 (Stability of \mathbf{w}_*). *If \mathbf{w}_* is a LME of $A(\mathbf{w}_*)$ and $\lambda_{\text{gap}}(\mathbf{w}_*) > \rho(\mathbf{w}_*)$, then \mathbf{w}_* is stable.*

This shows that lowering roughness measure $\rho(\mathbf{w}_*)$ at critical point \mathbf{w}_* could lead to more local optima and more patterns to be learned. To characterize such a behavior, we bound $\rho(\mathbf{w}_*)$:

Theorem 2 (Bound of local roughness $\rho(\mathbf{w})$ in ReLU setting). *If input $\|\mathbf{x}\|_2 \leq C_0$ is bounded, α has kernel structure (Def. 1) and batchsize $N \rightarrow +\infty$, then $\rho(\mathbf{w}_*) \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} r(\mathbf{w}_*, \alpha)$, where $r(\mathbf{w}, \alpha) := \sum_{l=0}^{+\infty} z_l^2(\alpha) \max_{\mathbf{w}^\top \mathbf{x}=0} \tilde{p}_l(\mathbf{x}; \alpha)$.*

From Thm. 2, the bound critically depends on $r(\alpha)$ that contains the *adjusted density* $\tilde{p}_l(\mathbf{x}; \alpha)$ (Def. 2) at the plane $\mathbf{w}_*^\top \mathbf{x} = 0$. This is because a local perturbation of \mathbf{w}_* leads to data inclusion/exclusion close to the plane, and thus changes $\rho(\mathbf{w}_*)$. Different α leads to different $\tilde{p}_l(\mathbf{x}; \alpha)$, and thus different upper bound of $\rho(\mathbf{w}_*)$, creating fewer or more local optima (i.e., patterns) to learn. Here is an example that shows Gaussian α (see Lemma 1), whose normalized version is used in InfoNCE, can lead to more local optima than uniform α , by lowering roughness bound characterized by $r(\mathbf{w}_*, \alpha)$:

Corollary 2 (Effect of different α). *For uniform α_u ($\alpha_{ij} := 1$) and 1-D Gaussian α_g ($\alpha_{ij} := \exp(-\|h(\mathbf{w}^\top \mathbf{x}_0[i]) - h(\mathbf{w}^\top \mathbf{x}_0[j])\|_2^2/2\tau)$), we have $r(\mathbf{w}_*, \alpha_g) = z_0(\alpha_g)r(\mathbf{w}_*, \alpha_u)$ with $z_0(\alpha_g) := \int \exp(-h^2(\mathbf{w}_*^\top \mathbf{x})/2\tau) p_D(\mathbf{x}) d\mathbf{x} \leq 1$. As a result, $z_0(\alpha_g) \ll 1$ leads to $r(\mathbf{w}_*, \alpha_g) \ll r(\mathbf{w}_*, \alpha_u)$.*

In practice, $z_0(\alpha_g)$ can be exponentially small (e.g., when most data appear on the positive side of the weight \mathbf{w}_*) and the roughness with Gaussian α can be much smaller than that of uniform α , which is presumably the reason why InfoNCE outperforms quadratic CL loss (Tian, 2022).

3.3 FINDING CRITICAL POINTS WITH INITIAL GUESS

In the following, we focus on how can we find an LME, when $A(\mathbf{w})$ does not have analytic form. We show that if there is an ‘‘approximate eigenvector’’ of $A(\mathbf{w}) := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}]$, then a real one is nearby.

Let L be the Lipschitz constant of $A(\mathbf{w})$: $\|A(\mathbf{w}) - A(\mathbf{w}')\|_2 \leq L\|\mathbf{w} - \mathbf{w}'\|_2$ for any \mathbf{w}, \mathbf{w}' on the unit sphere $\|\mathbf{w}\|_2 = 1$, and the *correlation function* $c(\mathbf{w}) := \mathbf{w}^\top \phi_1(\mathbf{w})$ be the inner product between \mathbf{w} and the maximal eigenvector of $A(\mathbf{w})$. We can construct a fixed point using *Power Iteration* (PI) (Golub & Van Loan, 2013), starting from initial value $\mathbf{w} = \mathbf{w}(0)$:

$$\mathbf{w}(t+1) \leftarrow A(\mathbf{w}(t))\mathbf{w}(t)/\|A(\mathbf{w}(t))\mathbf{w}(t)\|_2 \quad (\text{PI})$$

We show that even $A(\mathbf{w})$ varies over $\|\mathbf{w}\|_2 = 1$, the iteration can still converge to a fixed point \mathbf{w}_* , if the following quantity $\omega(\mathbf{w})$, called *irregularity*, is small enough.

Definition 4 (Irregularity $\omega(\mathbf{w})$ in the neighborhood of fixed points). Let $\mu(\mathbf{w}) := .5(1 + c(\mathbf{w}))c^{-2}(\mathbf{w}) [1 - \lambda_{\text{gap}}(\mathbf{w})/\lambda_1(\mathbf{w})]^2$ and $\omega(\mathbf{w}) := \omega(c(\mathbf{w}), \lambda_{\text{gap}}(\mathbf{w}), \lambda_1(\mathbf{w}), L, \kappa) \geq 0$ defined as

$$\omega(\mathbf{w}) := \mu(\mathbf{w}) + 2\kappa L^2(1 + \mu(\mathbf{w})c(\mathbf{w})) + 2L\lambda_{\text{gap}}^{-1}(\mathbf{w})\sqrt{\mu(\mathbf{w})(1 + \mu(\mathbf{w})c(\mathbf{w}))}, \quad (10)$$

here κ is the high-order eigenvector bound defined in Appendix (Lemma 9).

Intuitively, when $\mathbf{w}(0)$ is sufficiently close to any LME \mathbf{w}_* , i.e., $\mathbf{w}(0)$ is an ‘‘approximate’’ LME, we have $\omega(\mathbf{w}(0)) \ll 1$. In such a case, $\mathbf{w}(0)$ can be used to find \mathbf{w}_* using power iteration (Eqn. PI).

Theorem 3 (Existence of critical points). Let $c_0 := c(\mathbf{w}(0)) \neq 0$. If there exists $\gamma < 1$ so that:

$$\sup_{\mathbf{w} \in B_\gamma} \omega(\mathbf{w}) \leq \gamma, \quad (11)$$

where $B_\gamma := \left\{ \mathbf{w} : \mathbf{w}^\top \mathbf{w}(0) \geq \frac{c_0 - c_\gamma}{1 - c_\gamma}, c_\gamma := \frac{2\sqrt{\gamma}}{1 + \gamma} \right\}$ is the neighborhood of initial value $\mathbf{w}(0)$. Then Power Iteration (Eqn. PI) converges to a critical point $\mathbf{w}_* \in B_\gamma$ of Eqn. 7.

See proof in Appendix B.4. Intuitively, with L and κ small, c_0 close to 1, and λ_{gap} large, Eqn. 11 can always hold with $\gamma < 1$ and the fixed point exists. For example, for the two cases in Sec. 3.1, if $U = [\mathbf{u}_1, \dots, \mathbf{u}_M]$ is only approximately orthogonal (i.e., $\|U^\top U - I\|$ is not zero but small), and/or the conditions of Corollary 1 hold roughly, then Theorem 3 tells that multiple local optima close to \mathbf{u}_m still exist for each m (Fig. 1). We leave it for future work to further relax the condition.

Possible relation to empirical observations. Since there exist many local optima in the dynamics (Eqn. 7), even if objective involving \mathbf{w}_k are identical (Eqn. 5), each \mathbf{w}_k may still converge to different local optima due to initialization. We suspect that this can be a tentative explanation why larger model performs better: more local optima are collected and some can be *useful*. Other empirical observations like lottery ticket hypothesis (LTH) (Frankle & Carbin, 2019; Morcos et al., 2019; Tian et al., 2019; Yu et al., 2020), recently also verified in CL (Chen et al., 2021), may also be explained similarly. In LTH, first a large network is trained and pruned to be a small subnetwork \mathcal{S} , then retraining \mathcal{S} using its original initialization yields comparable or even better performance, while retraining \mathcal{S} with a different initialization performs much worse. For LTH, our explanation is that \mathcal{S} contains weights that are initialized *luckily*, i.e., close to useful local optima and converge to them during training. We leave a thorough empirical study to justify this line of thought for future work.

Given this intuition, it is tempting to study the *distribution* of local optima of Eqn. 7, their *attractive basin* $\text{Basin}(\mathbf{w}_*) := \{\mathbf{w} : \mathbf{w}(0) = \mathbf{w}, \lim_{t \rightarrow +\infty} \mathbf{w}(t) = \mathbf{w}_*\}$ for each local optimum \mathbf{w}_* , and the probability of random initialized weights fall into them. Interestingly, *data augmentation* may play an important role, by removing unnecessary local optima with symmetry (see Appendix B.5), focusing the learning on important patterns. Theorem 3 also gives hints. A formal study is left for future work.

4 TWO-LAYER SETTING

Now we understand how 1-layer nonlinearity learns in contrastive learning setting. In practice, many patterns exist and most of them may not be relevant for the downstream tasks. A natural question arises: how does the network prioritizes which patterns to learn? To answer this question, we analyze the behavior of 2-layer nonlinear networks with non-overlapping receptive fields (Fig. 2(a)).

Setting and Notations. In the lower layer, there are K disjoint *receptive fields* (abbreviated as **RF**) $\{R_k\}$, each has input \mathbf{x}_k and M weight $\mathbf{w}_{km} \in \mathbb{R}^d$ where $m = 1 \dots M$. The output of the bottom-layer is denoted as \mathbf{f}_1, f_{1km} for its km -th component, and $\mathbf{f}_1[i]$ for i -th sample. The top layer has weight $V \in \mathbb{R}^{d_{\text{out}} \times KM}$. Define $S := V^\top V$. As the $(km, k'm')$ entry of the matrix S , $s_{km, k'm'} := [S]_{km, k'm'} = \mathbf{v}_{km}^\top \mathbf{v}_{k'm'}$.

At each RF R_k , define $\tilde{\mathbf{x}}_{km}$ as an brief notation of gated input $\tilde{\mathbf{x}}_k^{w_{km}} := \mathbf{x}_k \cdot h'(\mathbf{w}_{km}^\top \mathbf{x}_k)$. Define $\tilde{\mathbf{x}}_k := [\tilde{\mathbf{x}}_{k1}; \tilde{\mathbf{x}}_{k2}; \dots, \tilde{\mathbf{x}}_{kM}] \in \mathbb{R}^{Md}$ as the concatenation of $\tilde{\mathbf{x}}_{km}$ and finally $\tilde{\mathbf{x}} := [\tilde{\mathbf{x}}_1; \dots; \tilde{\mathbf{x}}_K] \in \mathbb{R}^{KMd}$ is the concatenation of all $\tilde{\mathbf{x}}_k$. Similarly, let $\mathbf{w}_k := [\mathbf{w}_{km}]_{m=1}^M \in \mathbb{R}^{Md}$ be a concatenation of all \mathbf{w}_{km} in the same RF R_k , and $\mathbf{w} := [\mathbf{w}_k]_{k=1}^K \in \mathbb{R}^{KMd}$ be a column concatenation of all \mathbf{w}_k . Finally, $P_{\mathbf{w}}^\perp := \text{diag}_{km} [P_{\mathbf{w}_{km}}^\perp]$ is a block-diagonal matrix putting all projections together.

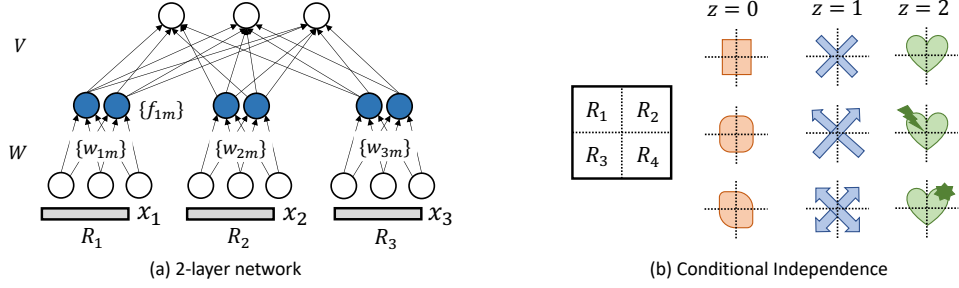


Figure 2: Our setting for 2-layer network. **(a)** We use W for low-layer weights and V for top-layer weights. There are K disjoint receptive fields (abbreviated as **RF**) R_k , each with M weight vectors in the low-layer, denoted as w_{km} . The activation function of hidden layer nodes is $h(x)$ and can be linear or nonlinear. **(b)** Conditional independence in Assumption 2: there exists a global categorical variable z . Given z , variation in different RFs are assumed to be independent.

Lemma 4 (Dynamics of 2-layer nonlinear network with contrastive loss).

$$\dot{V} = VC_\alpha[\mathbf{f}_1], \quad \dot{\mathbf{w}} = P_w^\perp [(S \otimes \mathbf{1}_d \mathbf{1}_d^\top) \circ C_\alpha[\tilde{\mathbf{x}}]] \mathbf{w} \quad (12)$$

where $\mathbf{1}_d$ is d -dimensional all-one vector, \otimes is Kronecker product and \circ is Hadamard product.

See Appendix C.1 for the proof. Now we analyze the stationary points of the equations. If $C_\alpha[\mathbf{f}_1]$ has unique maximal eigenvector \mathbf{s} , then following similar analysis as in Tian (2022), a necessary condition for (W, V) to be a stationary point is that $V = \mathbf{v} \mathbf{s}^\top$, where \mathbf{v} is any arbitrary unit vector. Therefore, we have $S = V^\top V = \mathbf{s} \mathbf{s}^\top$ as a rank-1 matrix and $s_{km, k'm'} = s_{km} s_{k'm'}$. Note that \mathbf{s} , as a unique maximal eigenvector of $C_\alpha[\mathbf{f}_1]$, is a function of the low-level feature computed by W .

On the other hand, the stationary point of W can be much more complicated, since it has the feedback term S from the top level. A more detailed analysis requires further assumptions, as we list below:

Assumption 2. For analysis of two-layer networks, we assume:

- **Uniform α , large batchsize and no augmentation.** Then $C_\alpha[g(\mathbf{x})] = \mathbb{V}[g(\mathbf{x})]$ for any function $g(\cdot)$ following Corollary 1.
- **Fast top-level training.** V undergoes fast training and has always converged to its stationary point given $C_\alpha[\mathbf{f}_1]$. That is, $S = \mathbf{s} \mathbf{s}^\top$ is a rank-1 matrix;
- **Conditional Independence.** The input in each R_k are conditional independent given a latent global random variable z taking C different values:

$$\mathbb{P}[\mathbf{x}|z] = \prod_{k=1}^K \mathbb{P}[\mathbf{x}_k|z] \quad (13)$$

Explanation of the assumptions. The *uniform α* condition is mainly for notation simplicity. For kernel-like α , the analysis is similar by combining multiple variance terms using Lemma 1. The *no augmentation* condition is mainly technical. Conclusion still holds if $\mathbb{E}_{p_{\text{aug}}}[\mathbf{g}(\mathbf{x})|\mathbf{x}_0] \approx \mathbf{g}(\mathbb{E}_{p_{\text{aug}}}[\mathbf{x}|\mathbf{x}_0])$ for $\mathbf{g}(\mathbf{x}) := \tilde{\mathbf{x}}^w$, i.e., augmentation swaps with nonlinear gating. For *conditional independence*, intuitively z can be regarded as different type of global patterns that determines what input \mathbf{x} can be perceived (Fig. 2(b)). Once z is given, the remaining variation resides within each RF R_k and independent across different R_k . Note that there exists many patterns in each RF R_k . Some are parts of the global pattern z , and others may come from noise. We study how each weight w_{km} captures distinct and useful patterns after training.

With all the assumptions, we can compute the term $A_k(\mathbf{w}_k) := C_\alpha[\tilde{\mathbf{x}}_k] = \mathbb{V}[\tilde{\mathbf{x}}_k]$. Our Assumption 2 is weaker than orthogonal mixture condition in Tian (2022) that is used to analyze CL, which requires the instance of input $\mathbf{x}_k[i]$ to have only one positive component.

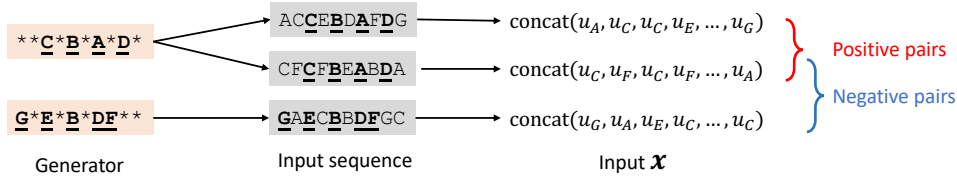


Figure 3: Experimental setting (Sec. 5). When generating input, we first randomly pick one generator (e.g., $**C*B*A*D*$) from a pool of G generators, generate the sequence by instantiating wildcard $*$ with an arbitrary token, and then replace the token a of sequence with an embedding vector u_a to form the input \mathbf{x} . Inputs from the same generator are treated as positive pairs, otherwise negative pairs for contrastive loss.

4.1 WHY NONLINEARITY IS CRITICAL: LINEAR ACTIVATION FAILS

Since in each RF R_k , there are M filters $\{\mathbf{w}_{km}\}$, it would be ideal to have one filter to capture one distinct pattern in the covariance matrix A_k . However, with linear activation, $\tilde{\mathbf{x}}_k = \mathbf{x}_k$ and as a result, learning of diverse features never happens, no matter how large M is (proof in Appendix C.4):

Theorem 4 (Gradient Colinearity in linear networks). *With linear activation, W follows the dynamics:*

$$\dot{\mathbf{w}}_{km} = s_{km} \mathbf{b}_k(W, V) \quad (14)$$

where $\mathbf{b}_k(W, V) := \mathbb{C}_\alpha \left[\mathbf{x}_k, \sum_{k', m'} s_{k'm'} \mathbf{w}_{k'm'}^\top \mathbf{x}_{k'} \right]$ is a linear function w.r.t. W . As a result, (1) $\dot{\mathbf{w}}_{km}$ are co-linear over m , and (2) If $s_{km} \neq 0$, from any critical point with distinct $\{\mathbf{w}_{km}\}$, there exists a path of critical points to identical weights ($\mathbf{w}_{km} = \mathbf{w}_k$).

This brings about the weakness of linear activation. First, the gradient of \mathbf{w}_{km} within one RF R_k during CL training all points towards the same direction \mathbf{b}_k ; Second, even if the critical points \mathbf{w}_{km} have any diversity within RF R_k , there exist a path for them to converge to identical weights. Therefore, diverse features, even they reside in the data, cannot be learned by the linear models.

4.2 THE EFFECT OF GLOBAL MODULATION IN THE SPECIAL CASE OF $C = 2$ AND $M = 1$

When z is binary ($C = 2$) with a single weight per RF ($M = 1$), \mathbf{w}_k 's dynamics has close form. Let \mathbf{w}_k represent \mathbf{w}_{k1} , the only weight at each R_k , $\Delta_k := \mathbb{E}[\tilde{\mathbf{x}}_k | z = 1] - \mathbb{E}[\tilde{\mathbf{x}}_k | z = 0]$. We have:

Theorem 5 (Dynamics of \mathbf{w}_k under conditional independence). *When $C = 2$ and $M = 1$, the dynamics of \mathbf{w}_k is given by (s_k^2 and $\delta_k \geq 0$ are scalars defined in the proof):*

$$\dot{\mathbf{w}}_k = P_{\mathbf{w}_k}^\perp (s_k^2 A_k(\mathbf{w}_k) + \delta_k \Delta_k \Delta_k^\top) \mathbf{w}_k \quad (15)$$

See proof in Appendix C.3. There are several interesting observations. First, the dynamics are decoupled (i.e., $\dot{\mathbf{w}}_k = A_k(W) \mathbf{w}_k$) and other $\mathbf{w}_{k'}$ with $k' \neq k$ only affects the dynamics of \mathbf{w}_k through the matrix $A_k(W)$. Second, while $A_k(\mathbf{w}_k)$ contains multiple patterns (i.e., local optima) in R_k , the additional term $\Delta_k \Delta_k^\top$, as the *global modulation* from the top level, encourages the model to learn the pattern like Δ_k which is a discriminative feature that separates the event of $z = 0$ and $z = 1$. Quantitatively:

Theorem 6 (Global modulation of attractive basin). *If the structural assumption holds: $A_k(\mathbf{w}_k) = \sum_l g(\mathbf{u}_l^\top \mathbf{w}_k) \mathbf{u}_l \mathbf{u}_l^\top$ with $g(\cdot) > 0$ a linear increasing function and $\{\mathbf{u}_l\}$ orthonormal bases, then for $A_k + c \mathbf{u}_l \mathbf{u}_l^\top$, its attractive basin of $\mathbf{w}_k = \mathbf{u}_l$ is larger than A_k 's for $c > 0$.*

Therefore, if Δ_k is a LME of A_k and \mathbf{w}_k is randomly initialized, Thm 5 tells that $\mathbb{P}[\mathbf{w}_k \rightarrow \Delta_k]$ is higher than the probability that \mathbf{w}_k goes to other patterns of A_k , i.e., the global variable z *modulates* the training of the lower layer. This is similar to ‘‘Backward feature correction’’ (Allen-Zhu & Li, 2020) and ‘‘top-down modulation’’ (Tian et al., 2019) in supervised learning, here we show it in CL.

We also analyze how BatchNorm helps alleviates diverse variances among RFs (see Appendix D).

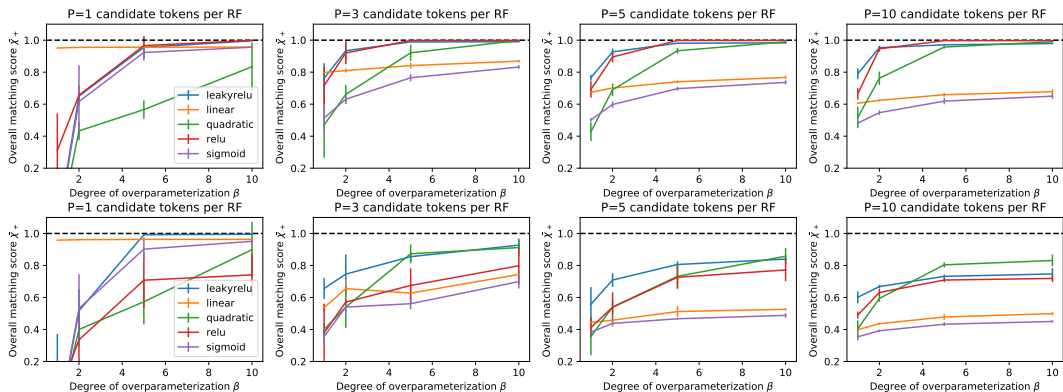


Figure 4: Overall matching score $\bar{\chi}_+$ (Eqn. 16) with InfoNCE (**top row**) and quadratic loss (**bottom row**). When $P = 1$, linear model works well regardless of the degree of over-parameterization β , while ReLU requires large over-parameterization to perform well. When each R_k has multiple patterns ($P > 1$) related to generators, ReLU models can capture diverse patterns better than linear ones in the over-parameterization region $\beta > 1$. We found similar trend for other homogeneous activations such as LeakyReLU (with negative slope 0.05) and quadratic. In contrast, linear models are much less affected by over-parameterization. While the trends are similar, quadratic loss is not as effective as InfoNCE in feature learning. Each setting is repeated 3 times and mean/std are reported. See Appendix (Fig. 9 and Fig. 10) for $\bar{\chi}_-$.

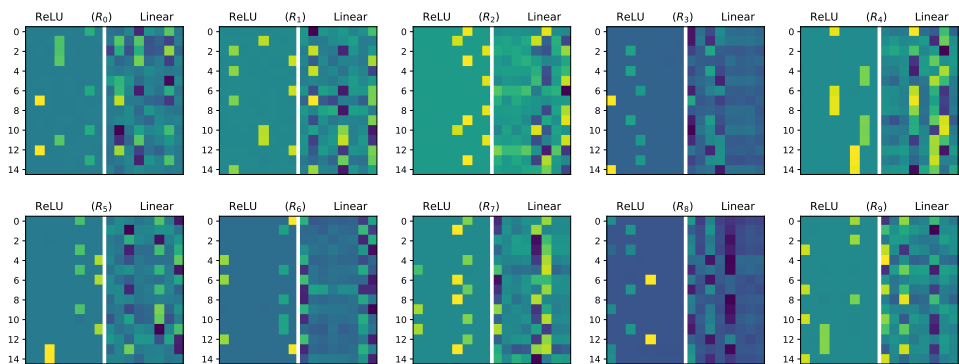


Figure 5: Visualization of learned weights with $P = 3$ (3 local patterns related to generators at each RF) and $\beta = 5$ (5x over-parameterization). Each of the $K = 10$ subfigures corresponds to a RF (R_0 - R_9). In each subfigure, the left panel is the learned weight by ReLU, while the right panel is from linear activations. 15 rows corresponds to $M = \beta P = 15$ weights and each weight is $d = 8$ dimensional. With ReLU activation, learned weights clearly capture the 3 candidate tokens within R_k^E at each RF R_k , while linear activation cannot.

5 EXPERIMENTS

Setup. To verify our finding, we perform contrastive learning with a 2-layer network on a synthetic dataset containing token sequences, generated as follows. From a pool of $G = 40$ generators, we pick a generator of length K in the form of $**C*B*A*D*$ (here $K = 10$) and generate $EF_{CDB}A_{AC}DB$ by sampling from $d = 20$ tokens for each wildcard $*$. The final input \mathbf{x} is then constructed by replacing each token a with the pre-defined embedding $\mathbf{u}_a \in \mathbb{R}^d$. $\{\mathbf{u}_a\}$ forms a orthonormal bases (see Fig. 3). The data augmentation is achieved by generating another sequence from the same generator.

While there exists $d = 20$ tokens, in each RF R_k we pick a subset R_k^g of $P < d$ tokens as the candidates used in the generator, to demonstrate the effect of global modulation. Before training, each generator is created by first randomly picking 5 receptive fields, then picking one of the P tokens from R_k^g at each RF R_k and filling the remaining RFs with wildcard $*$. Therefore, if a token appears at R_k but $a \notin R_k^g$, then a must be instantiated from the wildcard. Any $a \notin R_k^g$ is noise and should not to be learned in the weights of R_k since it is not part of any global pattern from the generator.

We train a 2-layer network on this dataset. The 2-layer network has $K = 10$ disjoint RFs, within each RF, there are $M = \beta P$ filters. Here $\beta \geq 1$ is a hyper-parameter that controls the degree of *over-parameterization*. The network is trained with InfoNCE loss and SGD with learning rate 2×10^{-3} , momentum 0.9, and weight decay 5×10^{-3} for 5000 minibatches and batchsize 128. Code is in PyTorch runnable on a single modern GPU.

Evaluation metric. We check whether the weights corresponding to each token is learned in the lower layer. At each RF R_k , we know R_k^g , the subsets of tokens it contains, as well as their embeddings $\{\mathbf{u}_a\}_{a \in R_k^g}$ due to the generation process, and verify whether these embeddings are learned after the model is trained. Specifically, for each token $a \in R_k^g$, we look for its best match on the learned filter $\{\mathbf{w}_{km}\}$, as formulated by the following per-RF score $\chi_+(R_k)$ and overall matching score $\bar{\chi}_+ \in [-1, 1]$ as the average over all RFs (similarly we can also define $\bar{\chi}_-$ for $a \notin R_k^g$):

$$\chi_+(R_k) = \frac{1}{P} \sum_{a \in R_k^g} \max_m \frac{\mathbf{w}_{km}^\top \mathbf{u}_a}{\|\mathbf{w}_{km}\|_2 \|\mathbf{u}_a\|_2}, \quad \bar{\chi}_+ = \frac{1}{K} \sum_k \chi_+(R_k) \quad (16)$$

5.1 RESULTS

Linear v.s ReLU activation and the effect of over-parameterization (Sec. 4.1). From Fig. 4, we can clearly see that ReLU (and other homogeneous) activations achieve better reconstruction of the input patterns, when each RF contains many patterns ($P > 1$) and specialization of filters in each RF is needed. On the other hand, when $P = 1$, linear activation works better. ReLU activation clearly benefits from over-parameterization ($\beta > 1$): the larger β is, the better $\bar{\chi}_+$ becomes. In contrast, for linear activation, over-parameterization does not quite affect the performance, which is consistent with our theoretical analysis.

Quadratic versus InfoNCE. Fig. 4 shows that quadratic CL loss underperforms InfoNCE, while the trend of linear/ReLU and over-parameterization remains similar. According to Corollary 2, non-uniform α (e.g., Gaussian α , Lemma 1) creates more and deeper local optima that better accommodate local patterns, yielding better performance. This provides a novel landscape point of view on why non-uniform α is better, expanding the intuition that it focuses more on important sample pairs.

Global modulation (Sec. 4.2). As shown in Fig. 5, the learned weights indeed focus on the token subset R_k^g that receives top-down support from the generators and no noise token is learned. We also verify that quantitatively by computing $\bar{\chi}_-$ over multiple runs, provided in Appendix (Fig. 9-10).

REFERENCES

- Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2019a. URL <https://openreview.net/forum?id=SkMQg3C5K7>.

- Sanjeev Arora, Hrishikesh Khandeparkar, Mikhail Khodak, Orestis Plevrakis, and Nikunj Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *arXiv preprint arXiv:1902.09229*, 2019b.
- Mary L Boas and Philip Peters. *Mathematical methods in the physical sciences*, 1984.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- James R Bunch, Christopher P Nielsen, and Danny C Sorensen. Rank-one modification of the symmetric eigenproblem. *Numerische Mathematik*, 31(1):31–48, 1978.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020.
- Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia Liu, Yang Zhang, Michael Carbin, and Zhangyang Wang. The lottery tickets hypothesis for supervised and self-supervised pre-training in computer vision models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16306–16316, 2021.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2020.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Aleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *CoRR*, abs/2204.02311, 2022. doi: 10.48550/arXiv.2204.02311. URL <https://doi.org/10.48550/arXiv.2204.02311>.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Luc Devroye, Laszlo Györfi, Adam Krzyżak, and Gábor Lugosi. On the strong universal consistency of nearest neighbor regression function estimates. *The Annals of Statistics*, 22(3):1371–1385, 1994.
- Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. *arXiv preprint arXiv:1806.00900*, 2018.
- Francisco M Fernández. *Introduction to perturbation theory in quantum mechanics*. CRC press, 2000.
- Evelyn Fix and Joseph Lawson Hodges. Nonparametric discrimination: consistency properties. *Randolph Field, Texas, Project*, pp. 21–49, 1951.

- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rJl-b3RcF7>.
- Benyamin Ghogh, Ali Ghodsi, Fakhri Karray, and Mark Crowley. Reproducing kernel hilbert space, mercer’s theorem, eigenfunctions, nyström method, and use of kernels in machine learning: Tutorial and survey. *arXiv preprint arXiv:2106.08443*, 2021.
- Mike B Giles. Collected matrix derivative results for forward and reverse mode algorithmic differentiation. In *Advances in Automatic Differentiation*, pp. 35–44. Springer, 2008.
- Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- Alain Goriely and Craig Hyde. Finite-time blow-up in dynamical systems. *Physics Letters A*, 250(4-6):311–318, 1998.
- Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *NeurIPS*, 2020.
- Ingo Gühring, Mones Raslan, and Gitta Kutyniok. Expressivity of deep neural networks. *arXiv preprint arXiv:2007.04759*, 2020.
- Barbara Hammer and Kai Gersmann. A note on the universal approximation capability of support vector machines. *neural processing letters*, 17(1):43–53, 2003.
- Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-supervised deep learning with spectral contrastive loss. *NeurIPS*, 2021.
- Soufiane Hayou, Arnaud Doucet, and Judith Rousseau. On the impact of the activation function on deep neural networks training. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2672–2680. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/hayou19a.html>.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. PMLR, 2015.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Wenlong Ji, Zhun Deng, Ryumei Nakada, James Zou, and Linjun Zhang. The power of contrast for feature learning: A theoretical analysis. *arXiv preprint arXiv:2110.02473*, 2021.
- Li Jing, Pascal Vincent, Yann LeCun, and Yuandong Tian. Understanding dimensional collapse in contrastive self-supervised learning. *ICLR*, 2022.
- Kenji Kawaguchi. Deep learning without poor local minima. *NeurIPS*, 2016.
- Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. *Advances in neural information processing systems*, 30, 2017.
- Zhou Lu, Hongming Pu, Feicheng Wang, Zhiqiang Hu, and Liwei Wang. The expressive power of neural networks: A view from the width. *Advances in neural information processing systems*, 30, 2017.

- Ari Morcos, Haonan Yu, Michela Paganini, and Yuandong Tian. One ticket to win them all: generalizing lottery ticket initializations across datasets and optimizers. *Advances in neural information processing systems*, 32, 2019.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Vern I Paulsen and Mrinal Raghupathi. *An introduction to the theory of reproducing kernel Hilbert spaces*, volume 152. Cambridge university press, 2016.
- Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *international conference on machine learning*, pp. 2847–2854. PMLR, 2017.
- Daniel A. Roberts, Sho Yaida, and Boris Hanin. The principles of deep learning theory. *CoRR*, abs/2106.10165, 2021. URL <https://arxiv.org/abs/2106.10165>.
- Nikunj Saunshi, Jordan Ash, Surbhi Goel, Dipendra Misra, Cyril Zhang, Sanjeev Arora, Sham Kakade, and Akshay Krishnamurthy. Understanding contrastive learning requires incorporating inductive biases. *arXiv preprint arXiv:2202.14037*, 2022.
- John Michael Tutill Thompson, H Bruce Stewart, and Rick Turner. Nonlinear dynamics and chaos. *Computers in Physics*, 4(5):562–563, 1990.
- Yuandong Tian. An analytical formula of population gradient for two-layered relu network and its applications in convergence and critical point analysis. In *International Conference on Machine Learning*, pp. 3404–3413. PMLR, 2017.
- Yuandong Tian. Understanding deep contrastive learning via coordinate-wise optimization. *NeurIPS*, 2022.
- Yuandong Tian, Tina Jiang, Qucheng Gong, and Ari Morcos. Luck matters: Understanding training dynamics of deep relu networks. *arXiv preprint arXiv:1905.13405*, 2019.
- Yuandong Tian, Lantao Yu, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning with dual deep networks. *arXiv preprint arXiv:2010.00578*, 2020.
- Yuandong Tian, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning dynamics without contrastive pairs. In *International Conference on Machine Learning*, pp. 10268–10278. PMLR, 2021.
- Xiang Wang, Xinlei Chen, Simon S Du, and Yuandong Tian. Towards demystifying representation learning with non-contrastive self-supervision. *arXiv preprint arXiv:2110.04947*, 2021.
- Charles L Wilson, James L Blue, and Omid M Omidvar. Training dynamics and neural network performance. *Neural Networks*, 10(5):907–923, 1997.
- Haonan Yu, Sergey Edunov, Yuandong Tian, and Ari S. Morcos. Playing the lottery with rewards and multiple languages: lottery tickets in rl and nlp. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=S1xnXRVFwH>.

A PROOFS

A.1 PROBLEM SETUP (SEC. 2)

Lemma 1 (Gaussian α). *For any function $\mathbf{g}(\cdot)$ that is bounded below, if we use $\alpha_{ij} := \exp(-\|\mathbf{g}(\mathbf{x}_0[i]) - \mathbf{g}(\mathbf{x}_0[j])\|_2^2/2\tau)$ as the pairwise importance, then it has kernel structure (Def. 1).*

Proof. Since $\mathbf{g}(\cdot)$ is bounded below, there exists a vector \mathbf{v} so that each component of $\mathbf{g}(\mathbf{x}) - \mathbf{v}$ is always nonnegative for any \mathbf{x} . Let $\mathbf{y}[i] := \mathbf{g}(\mathbf{x}_0[i]) - \mathbf{v} \in \mathbb{R}^d$, then $\mathbf{y}[i] \geq 0$ and we have:

$$\alpha_{ij} = \exp\left(-\frac{\|\mathbf{y}[i] - \mathbf{y}[j]\|_2^2}{2\tau}\right) \quad (17)$$

$$= \exp\left(-\frac{\|\mathbf{y}[i]\|_2^2}{2\tau}\right) \exp\left(-\frac{\|\mathbf{y}[j]\|_2^2}{2\tau}\right) \exp\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right) \quad (18)$$

And using Taylor expansion, we have

$$\exp\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right) = 1 + \frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau} + \frac{1}{2}\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right)^2 + \dots + \frac{1}{k!}\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right)^k + \dots \quad (19)$$

Let

$$\tilde{\phi}(\mathbf{y}) := \begin{bmatrix} 1 \\ \tau^{-1/2}\mathbf{y} \\ \frac{1}{\sqrt{2!}}\text{AllChoose}(\tau^{-1/2}\mathbf{y}, 2) \\ \dots \\ \frac{1}{\sqrt{k!}}\text{AllChoose}(\tau^{-1/2}\mathbf{y}, k) \\ \dots \end{bmatrix} \geq 0 \quad (20)$$

be an infinite dimensional vector, where $\text{AllChoose}(\mathbf{y}, k)$ is a d^k -dimensional column vector that enumerates all possible d^k products $y_{i_1}y_{i_2}\dots y_{i_k}$, where $1 \leq i_k \leq d$ and y_i is the i -th component of \mathbf{y} . Then it is clear that $\exp(\mathbf{y}^\top[i]\mathbf{y}[j]/\tau) = \tilde{\phi}^\top(\mathbf{y}[i])\tilde{\phi}(\mathbf{y}[j])$ and thus

$$\alpha_{ij} = \phi^\top(\mathbf{x}_0[i])\phi(\mathbf{x}_0[j]) = \sum_{l=0}^{+\infty} \phi_l(\mathbf{x}_0[i])\phi_l(\mathbf{x}_0[j]) \quad (21)$$

which satisfies Def. 1. Here

$$\phi(\mathbf{x}) := \exp\left(-\frac{\|\mathbf{y}\|_2^2}{2\tau}\right) \tilde{\phi}(\mathbf{y}) = \exp\left(-\frac{\|\mathbf{g}(\mathbf{x}) - \mathbf{v}\|_2^2}{2\tau}\right) \tilde{\phi}(\mathbf{g}(\mathbf{x}) - \mathbf{v}) \quad (22)$$

is the infinite dimensional feature mapping for input \mathbf{x} , and $\phi_l(\mathbf{x})$ is its l -th component. \square

Lemma 2 (Relationship between Contrastive Covariance and Variance in large batch size). *If α satisfies Def. 1, then for any function $\mathbf{g}(\cdot)$, $\mathbb{C}_\alpha[\mathbf{g}(\mathbf{x})]$ is asymptotically PSD when $N \rightarrow +\infty$:*

$$\mathbb{C}_\alpha[\mathbf{g}(\mathbf{x})] \rightarrow \sum_l z_l^2 \mathbb{V}_{\mathbf{x}_0 \sim \tilde{p}_l(\cdot; \alpha)} [\mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot | \mathbf{x}_0)}[\mathbf{g}(\mathbf{x}) | \mathbf{x}_0]] \quad (4)$$

Proof. First let

$$\mathbb{C}_\alpha^{\text{inter}}[\mathbf{a}, \mathbf{b}] := \frac{1}{2N^2} \sum_{i=1}^N \sum_{j \neq i} \alpha_{ij} (\mathbf{a}[i] - \mathbf{a}[j])(\mathbf{b}[i] - \mathbf{b}[j])^\top \quad (23)$$

$$\mathbb{C}_\alpha^{\text{intra}}[\mathbf{a}, \mathbf{b}] := \frac{1}{2N} \sum_{i=1}^N \left(\frac{1}{N} \sum_{j \neq i} \alpha_{ij} \right) (\mathbf{a}[i] - \mathbf{a}[i'])(\mathbf{b}[i] - \mathbf{b}[i'])^\top \quad (24)$$

and $\mathbb{C}_\alpha^{\text{inter}}[\mathbf{a}] := \mathbb{C}_\alpha^{\text{inter}}[\mathbf{a}, \mathbf{a}]$, $\mathbb{C}_\alpha^{\text{intra}}[\mathbf{a}] := \mathbb{C}_\alpha^{\text{inter}}[\mathbf{a}, \mathbf{a}]$. Then we have

$$\mathbb{C}_\alpha[\mathbf{g}] = \mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}] - \mathbb{C}_\alpha^{\text{intra}}[\mathbf{g}]. \quad (25)$$

With the condition, for the first term $\mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}]$, we have

$$\mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}] = \frac{1}{2N^2} \sum_{ij} \mathcal{K}(\mathbf{x}_0[i], \mathbf{x}_0[j]) (\mathbf{g}(\mathbf{x}[i]) - \mathbf{g}(\mathbf{x}[j])) (\mathbf{g}(\mathbf{x}[i]) - \mathbf{g}(\mathbf{x}[j]))^\top \quad (26)$$

When $N \rightarrow +\infty$, we have:

$$\mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}] \rightarrow \frac{1}{2} \int \mathcal{K}(\mathbf{x}_0, \mathbf{y}_0) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y}))^\top \mathbb{P}(\mathbf{x}, \mathbf{x}_0) \mathbb{P}(\mathbf{y}, \mathbf{y}_0) d\mathbf{x} d\mathbf{y} d\mathbf{x}_0 d\mathbf{y}_0$$

We integrate over \mathbf{x}_0 and \mathbf{y}_0 first:

$$\int (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y}))^\top \mathbb{P}(\mathbf{x}|\mathbf{x}_0) \mathbb{P}(\mathbf{y}|\mathbf{y}_0) d\mathbf{x} d\mathbf{y} \quad (27)$$

$$= \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] + \mathbb{E}_{\cdot|\mathbf{y}_0}[\mathbf{g}\mathbf{g}^\top] - \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}]\mathbb{E}_{\cdot|\mathbf{y}_0}[\mathbf{g}^\top] - \mathbb{E}_{\cdot|\mathbf{y}_0}[\mathbf{g}]\mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}^\top] \quad (28)$$

We now compute the four terms separately. With the condition that $\mathcal{K}(\mathbf{x}_0, \mathbf{y}_0) = \sum_l \phi_l(\mathbf{x}_0) \phi_l(\mathbf{y}_0)$, and the definition of adjusted probability $\tilde{p}_l(\mathbf{x}) := \frac{1}{z_l} \phi_l(\mathbf{x}) \mathbb{P}(\mathbf{x})$ where $z_l := \int \phi_l(\mathbf{x}) \mathbb{P}(\mathbf{x}) d\mathbf{x}$, for the first term, we have:

$$\begin{aligned} & \int \phi_l(\mathbf{x}_0) \phi_l(\mathbf{y}_0) \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] \mathbb{P}(\mathbf{x}_0) \mathbb{P}(\mathbf{y}_0) d\mathbf{x}_0 d\mathbf{y}_0 \\ &= z_l^2 \int \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] \tilde{p}_l(\mathbf{x}_0) d\mathbf{x}_0 \end{aligned} \quad (29)$$

$$= z_l^2 \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] \quad (30)$$

So we have:

$$\mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}] \rightarrow \sum_l z_l^2 (\mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] - \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}] \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}^\top]) \quad (31)$$

$$= \sum_l z_l^2 \mathbb{V}_{\mathbf{x}_0 \sim \tilde{p}_l, \mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \quad (32)$$

On the other hand, for $\mathbb{C}_\alpha^{\text{intra}}[\mathbf{g}]$, when $N \rightarrow +\infty$, we have:

$$\frac{1}{N} \sum_{j \neq i} \alpha_{ij} = \frac{1}{N} \sum_{j \neq i} \mathcal{K}(\mathbf{x}_0[i], \mathbf{x}_0[j]) \rightarrow \int \mathcal{K}(\mathbf{x}_0, \mathbf{y}_0) \mathbb{P}(\mathbf{y}_0) d\mathbf{y}_0 \quad (33)$$

$$= \sum_l \phi_l(\mathbf{x}_0) \int \phi_l(\mathbf{y}_0) \mathbb{P}(\mathbf{y}_0) d\mathbf{y}_0 = \sum_l z_l \phi_l(\mathbf{x}_0) \quad (34)$$

Therefore, we have:

$$\mathbb{C}_\alpha^{\text{intra}}[\mathbf{g}] \rightarrow \frac{1}{2} \sum_l z_l \int \phi_l(\mathbf{x}_0) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}')) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}'))^\top \mathbb{P}(\mathbf{x}, \mathbf{x}'|\mathbf{x}_0) \mathbb{P}(\mathbf{x}_0) d\mathbf{x} d\mathbf{x}' d\mathbf{x}_0 \quad (35)$$

Similarly,

$$\int (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}')) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}'))^\top \mathbb{P}(\mathbf{x}, \mathbf{x}'|\mathbf{x}_0) d\mathbf{x} d\mathbf{x}' \quad (36)$$

$$= 2 \int \mathbf{g}(\mathbf{x}) \mathbf{g}^\top(\mathbf{x}) \mathbb{P}(\mathbf{x}|\mathbf{x}_0) d\mathbf{x} - 2 \int \mathbf{g}(\mathbf{x}) \mathbb{P}(\mathbf{x}|\mathbf{x}_0) d\mathbf{x} \int \mathbf{g}^\top(\mathbf{x}') \mathbb{P}(\mathbf{x}'|\mathbf{x}_0) d\mathbf{x}' \quad (37)$$

$$= 2 \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}\mathbf{g}^\top] - 2 \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}^\top] \quad (38)$$

$$= 2 \mathbb{V}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \quad (39)$$

So we have:

$$\mathbb{C}_\alpha^{\text{intra}}[\mathbf{g}] \rightarrow \frac{1}{2} \sum_l z_l \int \phi_l(\mathbf{x}_0) 2 \mathbb{V}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \mathbb{P}(\mathbf{x}_0) d\mathbf{x}_0 \quad (40)$$

$$= \sum_l z_l^2 \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{V}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \quad (41)$$

Using the law of total variation, finally we have:

$$\mathbb{C}_\alpha[\mathbf{g}] \rightarrow \sum_l z_l^2 \mathbb{V}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \quad (42)$$

□

B ONE-LAYER MODEL (SEC. 3)

B.1 COMPUTATION OF THE TWO EXAMPLE MODELS

Here we assume ReLU activation $h(x) := \max(x, 0)$, which is a homogeneous activation $h(x) = h'(x)x$. Note that we consider $h'(0) = 0$. Therefore, for any sample \mathbf{x} , if $\mathbf{w}^\top \mathbf{x} = 0$, then we don't consider it to be included in the active region of ReLU, i.e., $\tilde{\mathbf{x}}^{\mathbf{w}} = \mathbf{x} \cdot h'(\mathbf{w}^\top \mathbf{x}) = 0$.

Let z be a hidden binary variable and we could compute $A(\mathbf{w})$ (here $p_0 := \mathbb{P}[z = 0]$ and $p_1 := \mathbb{P}[z = 1]$):

$$\mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}_z[\mathbb{E}[\tilde{\mathbf{x}}^{\mathbf{w}}|z]] + \mathbb{E}_z[\mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}|z]] = p_0 p_1 \Delta(\mathbf{w}) \Delta^\top(\mathbf{w}) + p_0 \Sigma_0(\mathbf{w}) + p_1 \Sigma_1(\mathbf{w}) \quad (43)$$

where $\Delta(\mathbf{w}) := \mathbb{E}[\tilde{\mathbf{x}}|z = 1] - \mathbb{E}[\tilde{\mathbf{x}}|z = 0]$ and $\Sigma_z(\mathbf{w}) := \mathbb{V}[\tilde{\mathbf{x}}|z]$.

Latent categorical model. If $\mathbf{w} = \mathbf{u}_m$, let $z := \mathbb{I}(y = m)$. This leads to $\Sigma_1(\mathbf{u}_m) = \Sigma_0(\mathbf{u}_m) = 0$ and $\Delta(\mathbf{u}_m) = \mathbf{u}_m$. Therefore, we have:

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{P}[y = m] (1 - \mathbb{P}[y = m]) \mathbf{u}_m \mathbf{u}_m^\top \quad (44)$$

Latent summation model. If $\mathbf{w} = \mathbf{u}_m$, first notice that due to orthogonal constraints we have $\mathbf{w}^\top \mathbf{x} = \sum_{m'} y_{m'} \mathbf{u}_{m'}^\top \mathbf{x} = y_m$. Let $z := \mathbb{I}(y_m > 0)$, then we can compute $\Delta(\mathbf{u}_m) = y_m^+ \mathbf{u}_m$, $\Sigma_1(\mathbf{u}_m) = I - \mathbf{u}_m \mathbf{u}_m^\top$ and $\Sigma_0(\mathbf{u}_m) = 0$. Therefore, we have:

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}] = (1 - q_m)^2 \mathbf{u}_m \mathbf{u}_m^\top + q_m (I - \mathbf{u}_m \mathbf{u}_m^\top) \quad (45)$$

B.2 DERIVATION OF TRAINING DYNAMICS

Lemma 3 (Training dynamics of 1-layer network with homogeneous activation in contrastive learning). *The gradient dynamics of Eqn. 5 is (note that α is treated as an independent variable):*

$$\dot{\mathbf{w}}_k = P_{\mathbf{w}_k}^\perp A(\mathbf{w}_k) \mathbf{w}_k \quad (7)$$

Here $P_{\mathbf{w}_k}^\perp := I - \mathbf{w}_k \mathbf{w}_k^\top$ projects a vector into the complementary subspace spanned by \mathbf{w}_k .

Proof. First of all, it is clear that from Eqn. 5, each \mathbf{w}_k evolves independently. Therefore, we omit the subscript k and derive the dynamics of one node \mathbf{w} .

To compute the training dynamics, we only need to compute the differential of $\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})]$. We use matrix differential form (Giles, 2008) to make the derivation easier to understand.

Note that for one-layer network with $K = 1$ nodes, $\mathcal{E}(\mathbf{w}) := \frac{1}{2} \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})] = \frac{1}{2} \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h(\mathbf{w}^\top \mathbf{x})]$ be the objective function to be maximized. Using the fact that

- $\mathbb{C}_\alpha[\mathbf{x}, \mathbf{y}]$ is a bilinear form (linear w.r.t \mathbf{x} and \mathbf{y}) given fixed α ,
- for any vector \mathbf{a} and \mathbf{b} , we have $\mathbf{a}^\top \mathbb{C}_\alpha[\mathbf{x}, \mathbf{y}] \mathbf{b} = \mathbb{C}_\alpha[\mathbf{a}^\top \mathbf{x}, \mathbf{b}^\top \mathbf{y}]$,
- for scalar x and y , $\mathbb{C}_\alpha[x, y] = \mathbb{C}_\alpha[y, x]$,

and by the product rule $d(x \cdot y) = dx \cdot y + x \cdot dy$, we have:

$$\begin{aligned} d\mathcal{E} &= \frac{1}{2} \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h'(\mathbf{w}^\top \mathbf{x}) d\mathbf{w}^\top \mathbf{x}] + \frac{1}{2} \mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x}) d\mathbf{w}^\top \mathbf{x}, h(\mathbf{w}^\top \mathbf{x})] \\ &= \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h'(\mathbf{w}^\top \mathbf{x}) \mathbf{x}] d\mathbf{w} \end{aligned} \quad (46)$$

Now use the homogeneous condition (Assumption 1) for activation h : $h(x) = h'(x)x$, which gives $h(\mathbf{w}^\top \mathbf{x}) = h'(\mathbf{w}^\top \mathbf{x}) \mathbf{w}^\top \mathbf{x}$, therefore, we have:

$$d\mathcal{E} = \mathbf{w}^\top \mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x}) \mathbf{x}, h'(\mathbf{w}^\top \mathbf{x}) \mathbf{x}] d\mathbf{w} = \mathbf{w}^\top A(\mathbf{w}) d\mathbf{w} \quad (47)$$

where $A(\mathbf{w}) := \mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x}) \mathbf{x}, h'(\mathbf{w}^\top \mathbf{x}) \mathbf{x}] = \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}, \tilde{\mathbf{x}}^{\mathbf{w}}]$. Therefore, by checking the coefficient associated with the differential form $d\mathbf{w}$, we know $\frac{\partial \mathcal{E}}{\partial \mathbf{w}} = A(\mathbf{w}) \mathbf{w}$. By gradient ascent, we have $\dot{\mathbf{w}} = A(\mathbf{w}) \mathbf{w}$. Since \mathbf{w} has the additional constraint $\|\mathbf{w}\|_2 = 1$, the final dynamics is $\dot{\mathbf{w}} = P_{\mathbf{w}}^\perp A(\mathbf{w}) \mathbf{w}$ where $P_{\mathbf{w}}^\perp := I - \mathbf{w} \mathbf{w}^\top$ is a projection matrix that projects a vector into the orthogonal complement subspace of the subspace spanned by \mathbf{w} . \square

Remarks. Note that an alternative route is to use homogeneous condition first: $\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})] = \mathbf{w}^\top A(\mathbf{x})\mathbf{w}$, then taking the differential. This involves an additional term $\frac{1}{2}\mathbf{w}^\top (dA)\mathbf{w}$. In the following we will show it is zero. For this we first compute dA :

$$dA = d\mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x})] \quad (48)$$

$$= \mathbb{C}_\alpha[h''(\mathbf{w}^\top \mathbf{x})(d\mathbf{w}^\top \mathbf{x})\mathbf{x}, h'(\mathbf{w}^\top \mathbf{x})\mathbf{x}] + \mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x})\mathbf{x}, h''(\mathbf{w}^\top \mathbf{x})(d\mathbf{w}^\top \mathbf{x})\mathbf{x}] \quad (49)$$

Therefore, since $\mathbf{a}^\top \mathbb{C}_\alpha[\mathbf{x}, \mathbf{y}]\mathbf{b} = \mathbb{C}_\alpha[\mathbf{a}^\top \mathbf{x}, \mathbf{b}^\top \mathbf{y}]$, we have:

$$\begin{aligned} \mathbf{w}^\top (dA)\mathbf{w} &= \mathbb{C}_\alpha[(d\mathbf{w}^\top \mathbf{x})h''(\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x}, h(\mathbf{w}^\top \mathbf{x})] + \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h''(\mathbf{w}^\top \mathbf{x})(d\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x}] \\ &= 2\mathbb{C}_\alpha[(d\mathbf{w}^\top \mathbf{x})h''(\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x}, h(\mathbf{w}^\top \mathbf{x})] \end{aligned} \quad (50)$$

Note that we now see the term $h''(\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x}$. For ReLU activation, its second derivative $h''(x) = \delta(x)$, where $\delta(x)$ is Direct delta function (Boas & Peters, 1984). From the property of delta function, we have $xh''(x) = x\delta(x) = 0$ even evaluated at $x = 0$. Therefore, $h''(\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x} = 0$ and $\mathbf{w}^\top (dA)\mathbf{w} = 0$. This is similar for LeakyReLU as well.

B.3 LOCAL STABILITY

Theorem 1 (Stability of \mathbf{w}_*). *If \mathbf{w}_* is a LME of $A(\mathbf{w}_*)$ and $\lambda_{\text{gap}}(\mathbf{w}_*) > \rho(\mathbf{w}_*)$, then \mathbf{w}_* is stable.*

Proof. For any unit direction $\|\mathbf{u}\|_2 = 1$ so that $\mathbf{u}^\top \mathbf{w}_* = 0$, consider the perturbation $\mathbf{v} = \sqrt{1 - \epsilon^2}\mathbf{w}_* + \epsilon\mathbf{u}$. Since $\|\mathbf{w}_*\|_2 = 1$ we have $\|\mathbf{v}\|_2 = 1$.

Now let's compute $P_v^\perp A(\mathbf{v})\mathbf{v}$. First, we have:

$$P_v^\perp = I - \mathbf{v}\mathbf{v}^\top = I - \left(\sqrt{1 - \epsilon^2}\mathbf{w}_* + \epsilon\mathbf{u}\right) \left(\sqrt{1 - \epsilon^2}\mathbf{w}_* + \epsilon\mathbf{u}\right)^\top \quad (51)$$

$$= I - \mathbf{w}_*\mathbf{w}_*^\top - \epsilon(\mathbf{u}\mathbf{w}_*^\top + \mathbf{w}_*\mathbf{u}^\top) + \mathcal{O}(\epsilon^2) \quad (52)$$

$$= P_{\mathbf{w}_*}^\perp - \epsilon(\mathbf{u}\mathbf{w}_*^\top + \mathbf{w}_*\mathbf{u}^\top) + \mathcal{O}(\epsilon^2) \quad (53)$$

So we have:

$$P_v^\perp A(\mathbf{w}_*)\mathbf{v} = P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*)\mathbf{v} - \epsilon(\mathbf{u}\mathbf{w}_*^\top + \mathbf{w}_*\mathbf{u}^\top)A(\mathbf{w}_*)\mathbf{v} + \mathcal{O}(\epsilon^2) \quad (54)$$

$$= P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*)\epsilon\mathbf{u} - \epsilon\lambda_*\mathbf{u} + \mathcal{O}(\epsilon^2) \quad (55)$$

$$= P_{\mathbf{w}_*}^\perp (A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} + \mathcal{O}(\epsilon^2) \quad (56)$$

The previous derivation is due to the fact that $P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*)\mathbf{w}_* = 0$, $\mathbf{u}^\top A(\mathbf{w}_*)\mathbf{w}_* = 0$ and $P_{\mathbf{w}_*}^\perp \mathbf{u} = \mathbf{u}$. Therefore, for $P_v^\perp A(\mathbf{v})\mathbf{v}$, we can decompose it to two parts:

$$P_v^\perp A(\mathbf{v})\mathbf{v} = P_v^\perp A(\mathbf{w}_*)\mathbf{v} + P_v^\perp (A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v} \quad (57)$$

$$= P_{\mathbf{w}_*}^\perp (A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} + P_v^\perp (A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v} + \mathcal{O}(\epsilon^2) \quad (58)$$

Therefore, since $\mathbf{u}^\top \mathbf{w}_* = 0$, we have:

$$\mathbf{u}^\top P_{\mathbf{w}_*}^\perp (A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} = \mathbf{u}^\top (I - \mathbf{w}_*\mathbf{w}_*^\top)(A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} \quad (59)$$

$$= \epsilon\mathbf{u}^\top (A(\mathbf{w}_*) - \lambda_*I)\mathbf{u} \leq -\lambda_{\text{gap}}(\mathbf{w}_*)\epsilon + \mathcal{O}(\epsilon^2) \quad (60)$$

and since $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$ and $\|P_v^\perp\|_2 = 1$, we have:

$$|\mathbf{u}^\top P_v^\perp (A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v}| \leq \|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v}\|_2 \quad (61)$$

By the definition of local roughness measure $\rho(\mathbf{w}_*)$, we have:

$$\|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{w}_*\|_2 \leq \rho(\mathbf{w}_*)\|\mathbf{v} - \mathbf{w}_*\|_2 + \mathcal{O}(\|\mathbf{v} - \mathbf{w}_*\|_2^2) = \rho(\mathbf{w}_*)\epsilon + \mathcal{O}(\epsilon^2) \quad (62)$$

This leads to

$$\|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v}\|_2 \leq \|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{w}_*\|_2 + \|(A(\mathbf{v}) - A(\mathbf{w}_*))(\mathbf{v} - \mathbf{w}_*)\|_2 \quad (63)$$

$$\leq \rho(\mathbf{w}_*)\epsilon + \mathcal{O}(\epsilon^2) \quad (64)$$

Therefore, we have:

$$\mathbf{u}^\top P_v^\perp A(\mathbf{v})\mathbf{v} \leq -(\lambda_{\text{gap}}(\mathbf{w}_*) - \rho(\mathbf{w}_*))\epsilon + \mathcal{O}(\epsilon^2) \quad (65)$$

When $\lambda_{\text{gap}}(\mathbf{w}_*) > \rho(\mathbf{w}_*)$ and we have $\mathbf{u}^\top P_v^\perp A(\mathbf{v})\mathbf{v} < 0$ for any $\mathbf{u} \perp \mathbf{w}_*$ and sufficiently small ϵ . Therefore, the critical point \mathbf{w}_* is stable. \square

Theorem 2 (Bound of local roughness $\rho(\mathbf{w})$ in ReLU setting). *If input $\|\mathbf{x}\|_2 \leq C_0$ is bounded, α has kernel structure (Def. 1) and batchsize $N \rightarrow +\infty$, then $\rho(\mathbf{w}_*) \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} r(\mathbf{w}_*, \alpha)$, where $r(\mathbf{w}, \alpha) := \sum_{l=0}^{+\infty} z_l^2(\alpha) \max_{\mathbf{w}^\top \mathbf{x}=0} \tilde{p}_l(\mathbf{x}; \alpha)$.*

Proof. Suppose \mathbf{w}_* and its local perturbation \mathbf{w} are on the unit sphere $\|\mathbf{w}\|_2 = \|\mathbf{w}_*\|_2 = 1$. Since \mathbf{w} is a local perturbation, we have $\mathbf{w}^\top \mathbf{w}_* \geq 1 - \epsilon$ for $\epsilon \ll 1$.

In the following we will check how we bound $\|(A(\mathbf{w}) - A(\mathbf{w}_*))\mathbf{w}_*\|_2$ in terms of $\|\mathbf{w} - \mathbf{w}_*\|_2$ and then we can get the upper bound of local roughness metric $\rho(\mathbf{w}_*)$.

Let the function $\mathbf{g}(\mathbf{x}) := \tilde{\mathbf{x}}^{\mathbf{w}}$, apply Corollary 1 with no augmentation and the large batch limits, we have

$$A(\mathbf{w}) := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \sum_l z_l^2 \mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}]. \quad (66)$$

where $\tilde{p}_l(\mathbf{x}) = \frac{1}{z_l} \mathbb{P}(\mathbf{x}) \phi_l(\mathbf{x})$ is the probability distribution of the input \mathbf{x} , adjusted by the mapping of the kernel function determined by the pairwise importance α_{ij} (Def. 1). z_l is its normalization constant.

To study $(A(\mathbf{w}) - A(\mathbf{w}_*))\mathbf{w}_*$, we will study each component $(\mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}] - \mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}_*}])\mathbf{w}_*$.

Note that since $\tilde{\mathbf{x}}^{\mathbf{w}} := \mathbf{x} \mathbb{I}(\mathbf{w}^\top \mathbf{x} \geq 0)$, we have $\mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{E}_{\tilde{p}_l}[\mathbf{x} \mathbf{x}^\top \mathbb{I}(\mathbf{w}^\top \mathbf{x} \geq 0)] - \mathbb{E}_{\tilde{p}_l}[\mathbf{x} \mathbb{I}(\mathbf{w}^\top \mathbf{x} \geq 0)] \mathbb{E}_{\tilde{p}_l}[\mathbf{x}^\top \mathbb{I}(\mathbf{w}^\top \mathbf{x} \geq 0)]$. Let

$$\mathbf{e} := \int_{\mathbf{w}^\top \mathbf{x} \geq 0} \mathbf{x} \tilde{p}_l(\mathbf{x}) d\mathbf{x}, \quad \mathbf{e}_* := \int_{\mathbf{w}_*^\top \mathbf{x} \geq 0} \mathbf{x} \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (67)$$

$$E := \int_{\mathbf{w}^\top \mathbf{x} \geq 0} \mathbf{x} \mathbf{x}^\top \tilde{p}_l(\mathbf{x}) d\mathbf{x}, \quad E_* := \int_{\mathbf{w}_*^\top \mathbf{x} \geq 0} \mathbf{x} \mathbf{x}^\top \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (68)$$

So we can write

$$\mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}] = E - \mathbf{e} \mathbf{e}^\top, \quad \mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}_*}] = E_* - \mathbf{e}_* \mathbf{e}_*^\top \quad (69)$$

and $\mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}] - \mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}_*}] = (E - E_*) + (\mathbf{e}_* \mathbf{e}_*^\top - \mathbf{e} \mathbf{e}^\top)$.

Define the following regions

$$\Omega_+ := \{\mathbf{x} : \mathbf{w}_*^\top \mathbf{x} \geq 0, \mathbf{w}^\top \mathbf{x} \leq 0\} \quad (70)$$

$$\Omega_- := \{\mathbf{x} : \mathbf{w}_*^\top \mathbf{x} \leq 0, \mathbf{w}^\top \mathbf{x} \geq 0\} \quad (71)$$

$$\Omega := \Omega_+ \cup \Omega_- \quad (72)$$

Now let's bound $(E - E_*)\mathbf{w}_*$ and $(\mathbf{e}_* \mathbf{e}_*^\top - \mathbf{e} \mathbf{e}^\top)\mathbf{w}_*$.

Bound $(E - E_*)\mathbf{w}_*$. We have:

$$E - E_* = \int_{\Omega_-} \mathbf{x} \mathbf{x}^\top \tilde{p}_l(\mathbf{x}) d\mathbf{x} - \int_{\Omega_+} \mathbf{x} \mathbf{x}^\top \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (73)$$

and thus

$$(E - E_*)\mathbf{w}_* = \int_{\Omega_-} \mathbf{x} \mathbf{x}^\top \mathbf{w}_* \tilde{p}_l(\mathbf{x}) d\mathbf{x} - \int_{\Omega_+} \mathbf{x} \mathbf{x}^\top \mathbf{w}_* \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (74)$$

For any $\mathbf{x} \in \Omega_+$, we have:

$$0 \leq \mathbf{w}_*^\top \mathbf{x} = \mathbf{w}^\top \mathbf{x} + (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \leq (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \leq C_0 \|\mathbf{w}_* - \mathbf{w}\|_2 \quad (75)$$

Therefore, $|\mathbf{w}_*^\top \mathbf{x}| \leq M \|\mathbf{w}_* - \mathbf{w}\|_2$ and we have

$$\left\| \int_{\Omega_+} \mathbf{x} \mathbf{x}^\top \mathbf{w}_* \tilde{p}_l(\mathbf{x}) d\mathbf{x} \right\|_2 \leq \int_{\Omega_+} |\mathbf{w}_*^\top \mathbf{x}| \|\mathbf{x}\|_2 \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (76)$$

$$\leq C_0^2 \|\mathbf{w}_* - \mathbf{w}\|_2 \max_{\mathbf{x} \in \Omega_+} \tilde{p}_l(\mathbf{x}) \int_{\Omega_+, \|\mathbf{x}\|_2 \leq C_0} d\mathbf{x} \quad (77)$$

$$= C_0^3 \|\mathbf{w}_* - \mathbf{w}\|_2 \max_{\mathbf{x} \in \Omega_+} \tilde{p}_l(\mathbf{x}) \frac{\text{vol}(C_0)}{2\pi} \arccos \mathbf{w}^\top \mathbf{w}_* \quad (78)$$

where $\text{vol}(C_0)$ is the volume of the d -dimensional ball of radius C_0 . Similarly for $\mathbf{x} \in \Omega_-$, we have

$$0 \geq \mathbf{w}_*^\top \mathbf{x} = \mathbf{w}^\top \mathbf{x} + (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \geq (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \geq -C_0 \|\mathbf{w}_* - \mathbf{w}\|_2 \quad (79)$$

hence $|\mathbf{w}_*^\top \mathbf{x}| \leq C_0 \|\mathbf{w}_* - \mathbf{w}\|_2$ and overall we have:

$$\|(E - E_*)\mathbf{w}_*\|_2 \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \|\mathbf{w}_* - \mathbf{w}\|_2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \arccos \mathbf{w}^\top \mathbf{w}_* \quad (80)$$

Since for $x \in (0, 1]$, $\arcsin \sqrt{1-x^2} \leq \frac{\sqrt{1-x^2}}{x}$, we have:

$$\arccos \mathbf{w}^\top \mathbf{w}_* = \arcsin \sqrt{1 - (\mathbf{w}^\top \mathbf{w}_*)^2} \leq \frac{\sqrt{1 - (\mathbf{w}^\top \mathbf{w}_*)^2}}{\mathbf{w}^\top \mathbf{w}_*} \quad (81)$$

$$= \frac{\sqrt{1 + \mathbf{w}^\top \mathbf{w}_*} \sqrt{1 - \mathbf{w}^\top \mathbf{w}_*}}{\mathbf{w}^\top \mathbf{w}_*} \leq \frac{\sqrt{2(1 - \mathbf{w}^\top \mathbf{w}_*)}}{\mathbf{w}^\top \mathbf{w}_*} \quad (82)$$

$$= \frac{1}{1 - \epsilon} \|\mathbf{w} - \mathbf{w}_*\|_2 \quad (83)$$

we have:

$$\|(E - E_*)\mathbf{w}_*\|_2 \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \frac{1}{1 - \epsilon} \|\mathbf{w}_* - \mathbf{w}\|_2^2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \quad (84)$$

Therefore, $\|(E - E_*)\mathbf{w}_*\|_2$ is a second-order term w.r.t. $\|\mathbf{w} - \mathbf{w}_*\|_2$.

Bound $(\mathbf{e}_* \mathbf{e}_*^\top - \mathbf{e} \mathbf{e}^\top) \mathbf{w}_*$. On the other hand:

$$\mathbf{e} \mathbf{e}^\top - \mathbf{e}_* \mathbf{e}_*^\top = \mathbf{e}(\mathbf{e} - \mathbf{e}_*)^\top + (\mathbf{e} - \mathbf{e}_*) \mathbf{e}_*^\top \quad (85)$$

We have $\|\mathbf{e}\|_2, \|\mathbf{e}_*\|_2$ bounded and

$$\mathbf{e} - \mathbf{e}_* = \int_{\Omega_-} \mathbf{x} \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} - \int_{\Omega_+} \mathbf{x} \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} \quad (86)$$

Using similar derivation, we conclude that $\|\mathbf{e}(\mathbf{e} - \mathbf{e}_*)^\top \mathbf{w}_*\|_2$ is also a second-order term. The only first-order term is $\|(\mathbf{e} - \mathbf{e}_*) \mathbf{e}_*^\top \mathbf{w}_*\|_2$:

$$\|(\mathbf{e} - \mathbf{e}_*) \mathbf{e}_*^\top \mathbf{w}_*\|_2 \leq \mathbb{E}_{\tilde{\rho}_l} [h(\mathbf{w}^\top \mathbf{x})] \int_{\Omega} \|\mathbf{x}\|_2 \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} \quad (87)$$

$$\leq C_0^2 \int_{\Omega} \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} \leq C_0^2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \int_{\Omega: \|\mathbf{x}\|_2 \leq C_0} d\mathbf{x} \quad (88)$$

$$\leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \arccos \mathbf{w}^\top \mathbf{w}_* \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \quad (89)$$

$$\leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \frac{1}{1 - \epsilon} \|\mathbf{w} - \mathbf{w}_*\|_2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \quad (90)$$

Overall we have:

$$\|(A(\mathbf{w}) - A(\mathbf{w}_*))\mathbf{w}_*\|_2 \leq \sum_l z_l^2 \|(\nabla_{\tilde{\rho}_l} [h(\mathbf{x}^\top \mathbf{w})] - \nabla_{\tilde{\rho}_l} [h(\mathbf{x}^\top \mathbf{w}_*)]) \mathbf{w}_*\|_2 \quad (91)$$

$$\leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \frac{1}{1 - \epsilon} \left(\sum_l z_l^2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \right) \|\mathbf{w} - \mathbf{w}_*\|_2 + \mathcal{O}(\|\mathbf{w} - \mathbf{w}_*\|_2^2) \quad (92)$$

Since $\rho(\mathbf{w}_*)$ is the smallest scalar that makes the local roughness metric hold and ϵ is arbitrarily small, we have:

$$\rho(\mathbf{w}_*) \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} r(\mathbf{w}_*, \alpha) \quad (93)$$

where $r(\mathbf{w}, \alpha) := \sum_l z_l^2 \max_{\mathbf{w}^\top \mathbf{x}=0} \tilde{\rho}_l(\mathbf{x}; \alpha)$. \square

Corollary 2 (Effect of different α). *For uniform α_u ($\alpha_{ij} := 1$) and 1-D Gaussian α_g ($\alpha_{ij} := \exp(-\|h(\mathbf{w}^\top \mathbf{x}_0[i]) - h(\mathbf{w}^\top \mathbf{x}_0[j])\|_2^2 / 2\tau)$), we have $r(\mathbf{w}_*, \alpha_g) = z_0(\alpha_g) r(\mathbf{w}_*, \alpha_u)$ with $z_0(\alpha_g) := \int \exp(-h^2(\mathbf{w}_*^\top \mathbf{x}) / 2\tau) p_D(\mathbf{x}) d\mathbf{x} \leq 1$. As a result, $z_0(\alpha_g) \ll 1$ leads to $r(\mathbf{w}_*, \alpha_g) \ll r(\mathbf{w}_*, \alpha_u)$.*

Proof. For uniform α_u , it is clear that the mapping $\phi_u(\mathbf{x}) \equiv 1$ is 1-dimensional. Therefore, $\tilde{p}_0(\mathbf{x}; \alpha_u) := \frac{1}{z_0(\alpha_u)} \phi_{u0}(\mathbf{x}) p_D(\mathbf{x}) = p_D(\mathbf{x})$ with $z_0(\alpha_u) = \int \phi_{u0}(\mathbf{x}) p_D(\mathbf{x}) d\mathbf{x} = 1$. This means that

$$r(\mathbf{w}_*, \alpha_u) := \sum_{l=0}^{+\infty} z_l^2(\alpha_u) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{p}_l(\mathbf{x}; \alpha_u) \quad (94)$$

$$= z_0^2(\alpha_u) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{p}_0(\mathbf{x}; \alpha_u) = \max_{\mathbf{w}_*^\top \mathbf{x}=0} p_D(\mathbf{x}) \quad (95)$$

For Gaussian α_g , from Lemma 1 we know that its infinite-dimensional mapping $\phi_g(\mathbf{x})$ has the following form for $\mathbf{w} = \mathbf{w}_*$:

$$\phi_g(\mathbf{x}) = e^{-\frac{h^2(\mathbf{w}_*^\top \mathbf{x})}{2\tau}} \begin{bmatrix} 1 \\ \tau^{-1/2} h(\mathbf{w}_*^\top \mathbf{x}) \\ \frac{1}{\tau^{2/2} \sqrt{2!}} h^2(\mathbf{w}_*^\top \mathbf{x}) \\ \dots \\ \frac{1}{\tau^{k/2} \sqrt{k!}} h^k(\mathbf{w}_*^\top \mathbf{x}) \\ \dots \end{bmatrix} \quad (96)$$

When $l \geq 1$, $z_l^2 \tilde{p}_l(\mathbf{x}; \alpha_g) = z_l \phi_{gl}(\mathbf{x}) p_D(\mathbf{x}) = 0$ for any \mathbf{x} on the plane $\mathbf{w}_*^\top \mathbf{x} = 0$, since $\phi_{gl}(\mathbf{x}) = 0$ on the plane. On the other hand, $\phi_{g0}(\mathbf{x}) = e^{-\frac{h^2(\mathbf{w}_*^\top \mathbf{x})}{2\tau}}$. On the plane, $\phi_{g0}(\mathbf{x}) = 1$ and is a constant. Therefore, we have:

$$r(\mathbf{w}_*, \alpha_g) := \sum_{l=0}^{+\infty} z_l^2 \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{p}_l(\mathbf{x}; \alpha_g) = z_0^2(\alpha_g) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{p}_0(\mathbf{x}; \alpha_g) \quad (97)$$

$$= z_0(\alpha_g) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \phi_{g0}(\mathbf{x}) p_D(\mathbf{x}) \quad (98)$$

$$= z_0(\alpha_g) \max_{\mathbf{w}_*^\top \mathbf{x}=0} p_D(\mathbf{x}) = z_0(\alpha_g) r(\mathbf{w}_*, \alpha_u) \quad (99)$$

Here

$$z_0(\alpha_g) := \int \phi_{g0}(\mathbf{x}) p_D(\mathbf{x}) d\mathbf{x} = \int e^{-\frac{h^2(\mathbf{w}_*^\top \mathbf{x})}{2\tau}} p_D(\mathbf{x}) d\mathbf{x} \leq 1 \quad (100)$$

□

B.4 FINDING CRITICAL POINTS WITH INITIAL GUESS (SEC. 3.3)

Notation. Let $\lambda_i(\mathbf{w})$ and $\phi_i(\mathbf{w})$ be the i -th eigenvalue and unit eigenvector of $A(\mathbf{w})$ where $\phi_1(\mathbf{w})$ is the largest. We first assume $A(\mathbf{w})$ is positive definite (PD) and then remove this assumption later. In this case, $\lambda_1(\mathbf{w}) \geq \lambda_2(\mathbf{w}) \geq \dots \geq \lambda_d(\mathbf{w}) > 0$. Let $c(\mathbf{w}) := \mathbf{w}^\top \phi_1(\mathbf{w})$ be the inner product between \mathbf{w} and the maximal eigenvector of $A(\mathbf{w})$.

Consider the following Power Iteration (PI) format:

$$\tilde{\mathbf{w}}(t+1) \leftarrow A(\mathbf{w}(t))\mathbf{w}(t), \quad \mathbf{w}(t+1) \leftarrow \frac{\tilde{\mathbf{w}}(t+1)}{\|\tilde{\mathbf{w}}(t+1)\|_2} \quad (101)$$

Along the trajectory, let $\phi_i(t) := \phi_i(A(\mathbf{w}(t)))$ be the i -th unit eigenvector of $A(\mathbf{w}(t))$ and $\lambda_i(t)$ to be the i -th eigenvalue. Define $\delta\mathbf{w}(t) := \mathbf{w}(t+1) - \mathbf{w}(t)$, $\delta A(t) := A(\mathbf{w}(t+1)) - A(\mathbf{w}(t))$, and

$$c_t := c(\mathbf{w}(t)) = \phi_1^\top(t)\mathbf{w}(t), \quad d_t := \phi_1^\top(t)\mathbf{w}(t+1) \quad (102)$$

Then $-1 \leq c_t, d_t \leq 1$ since they are inner product of two unit vectors.

Theorem 3 (Existence of critical points). *Let $c_0 := c(\mathbf{w}(0)) \neq 0$. If there exists $\gamma < 1$ so that:*

$$\sup_{\mathbf{w} \in B_\gamma} \omega(\mathbf{w}) \leq \gamma, \quad (11)$$

where $B_\gamma := \left\{ \mathbf{w} : \mathbf{w}^\top \mathbf{w}(0) \geq \frac{c_0 - c_\gamma}{1 - c_\gamma}, c_\gamma := \frac{2\sqrt{\gamma}}{1+\gamma} \right\}$ is the neighborhood of initial value $\mathbf{w}(0)$. Then Power Iteration (Eqn. PI) converges to a critical point $\mathbf{w}_* \in B_\gamma$ of Eqn. 7.

Proof. Note that if $c_0 < 0$, we can always use $-\phi_1(\mathbf{w})$ as the maximal eigenvector.

First we assume $A(\mathbf{w})$ is positive definite (PD) over the entire unit sphere $\|\mathbf{w}\|_2 = 1$, then follow Lemma 11, and notice that $\|\mathbf{w} - \mathbf{w}(0)\|_2 = \sqrt{2(1 - \mathbf{w}^\top \mathbf{w}(0))}$, so

$$\|\mathbf{w} - \mathbf{w}(0)\|_2 \leq \frac{\sqrt{2(1 + \gamma)(1 - c_0)}}{1 - \sqrt{\gamma}} \iff \mathbf{w}^\top \mathbf{w}(0) \geq \frac{c_0 - c_\gamma}{1 - c_\gamma} \quad (103)$$

When $A(\mathbf{w})$ is not PD, Theorem 3 still applies to the PD matrix $\hat{A}(\mathbf{w}) := A(\mathbf{w}) - \lambda_{\min}(\mathbf{w})I + \epsilon I$ with L and κ specified by $\hat{A}(\mathbf{w})$, where $\epsilon > 0$ is a small constant.

This transformation keeps c_0 since the eigenvectors of $\hat{A}(\mathbf{w})$ are the same as $A(\mathbf{w})$. The resulting fixed point $\hat{\mathbf{w}}_*$ is also the fixed point of the original problem with $A(\mathbf{w})$, due to the fact that

$$P_{\mathbf{w}}^\perp \hat{A}(\mathbf{w})\mathbf{w} = P_{\mathbf{w}}^\perp A(\mathbf{w})\mathbf{w} - (\lambda_{\min}(\mathbf{w}) - \epsilon)P_{\mathbf{w}}^\perp \mathbf{w} = P_{\mathbf{w}}^\perp A(\mathbf{w})\mathbf{w} \quad (104)$$

□

Remarks. Note that Lemma 11 assumes that along the trajectory $\{\mathbf{w}(t)\}$, $\mu_t + \nu_t \leq \gamma$ holds. In Theorem 3, this can not be assumed true until we prove that the entire trajectory is within B_γ .

B.5 THE EFFECT OF DATA AUGMENTATION ON LOCAL OPTIMA

While the majority of the analysis focuses on the cases where there are no data augmentation (i.e., using Corollary 1), the original formulation Lemma 2 can still handle contrastive learning in the presence of data augmentation.

In fact, data augmentation plays an important role by removing unnecessary local optima. First, Lemma 2 tells that the objective Eqn. 5, when $K = 1$, takes the following form:

$$2\mathcal{E}_\alpha(\mathbf{w}) := \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})] \rightarrow \sum_l z_l^2 \mathbb{V}_{\mathbf{x}_0 \sim \tilde{p}_l(\cdot; \alpha)} [b(\mathbf{w}|\mathbf{x}_0)] \quad (105)$$

where $b(\mathbf{w}|\mathbf{x}_0) := \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)} [h(\mathbf{w}^\top \mathbf{x})|\mathbf{x}_0]$.

Now let us consider the following simple data augmentation of \mathbf{x}_0 :

$$\mathbf{x} = R(t)\mathbf{x}_0, \quad t \sim \text{Uniform}(\mathcal{T}) \quad (106)$$

where $R(t) \in \mathbb{R}^{d \times d}$ is some rotation parameterized by t , which is drawn uniformly from a parameter family \mathcal{T} .

We assume $\{R(t)\}_{t \in \mathcal{T}}$ forms a 1-dimensional *Lie group* parameterized by \mathcal{T} . This means that

- **Closeness.** For any $t, t' \in \mathcal{T}$, there exists $t'' \in \mathcal{T}$ so that $R(t'') = R(t)R(t')$.
- **Existence of inverse element.** For each $R(t)$, there exists an inverse element $t' \in \mathcal{T}$ so that $R(t') = R^{-1}(t) = R^\top(t)$. The last equality is due to the fact that $R(t)$ is a rotation.
- **Existence of identity map.** $R(0) = I$.

Then for any small transformation $R(t')$ applied to the weights \mathbf{w} (here “small” means $\|R(t') - I\|_2$ is small), we can write down $b(R(t')\mathbf{w}|\mathbf{x}_0)$ using reparameterization trick:

$$\begin{aligned} b(R(t)\mathbf{w}|\mathbf{x}_0) &:= \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)} [h((R(t)\mathbf{w})^\top \mathbf{x})|\mathbf{x}_0] = \int h(\mathbf{w}^\top R^\top(t)R(t')\mathbf{x}_0) \mathbb{P}[t'] dt' \\ &= \int h(\mathbf{w}^\top R^{-1}(t)R(t')\mathbf{x}_0) \mathbb{P}[t'] dt' \end{aligned} \quad (107)$$

$$= \int h(\mathbf{w}^\top R(t'')\mathbf{x}_0) \mathbb{P}[t''] dt'' = b(\mathbf{w}|\mathbf{x}_0) \quad (108)$$

Note that the last equality is due to the fact that $\{R(t)\}_{t \in \mathcal{T}}$ is a Lie group, so that $R^{-1}(t)R(t')$ always maps to another group element $R(t'')$, and t'' as the resulting parameterization is still uniform.

Due to stop gradient, α and thus $\phi_l(\cdot; \alpha)$ is treated as a constant term when checking the local property of the current parameters \mathbf{w} . This means that in the local neighborhood of \mathbf{w} , $\mathcal{E}_\alpha(\mathbf{w}) = \mathcal{E}_\alpha(R(t')\mathbf{w})$.

Now notice an important observation: if $\mathbf{w}' := R(t')\mathbf{w} \neq \mathbf{w}$, then $\mathcal{E}_\alpha(\mathbf{w}) = \mathcal{E}_\alpha(R(t')\mathbf{w}) = \mathcal{E}_\alpha(\mathbf{w}')$ and therefore, \mathbf{w} *cannot* be a local optimal.

Intuitively, this means that the data augmentation can *remove* certain local optima of \mathbf{w} , if they are not *locally invariant* (i.e., $R(t')\mathbf{w} \neq \mathbf{w}$) to the transformation of the data augmentation. Therefore, augmentation removes certain patterns in the input data and their local optima in the training, to only keep patterns (local optima) that are most relevant to the tasks.

Here we only use 1-dimensional rotation group as one simple example. In practice, the augmentation may not globally form a Lie group, and there could be multiple different types of augmentations, yielding high-dimensional transformation space. Therefore, we may use Lie algebra instead to capture the local transformation structure, without making assumptions about the global structure. We will give a formal study in the future work.

C TWO LAYER CASE (SEC. 4)

C.1 LEARNING DYNAMICS

Lemma 4 (Dynamics of 2-layer nonlinear network with contrastive loss).

$$\dot{V} = VC_\alpha[\mathbf{f}_1], \quad \dot{\mathbf{w}} = P_{\mathbf{w}}^\perp [(S \otimes \mathbf{1}_d \mathbf{1}_d^\top) \circ \mathbb{C}_\alpha[\tilde{\mathbf{x}}]] \mathbf{w} \quad (12)$$

where $\mathbf{1}_d$ is d -dimensional all-one vector, \otimes is Kronecker product and \circ is Hadamard product.

Proof. The output of the 2-layer network can be written as the following:

$$f_{2l} = \sum_k v_{lk} h(\mathbf{w}_k^\top \mathbf{x}_k) \quad (109)$$

For convenience, we use $\mathbf{f}_1 := [h(\mathbf{w}_k^\top \mathbf{x}_k)]$ to represent the column vector that collects all the outputs of intermediate nodes, and \mathbf{v}_l^\top is the l -th row vector in V .

According to Theorem 1 in Tian (2022), the gradient descent direction of contrastive loss corresponds to the gradient ascent direction of the energy function $\mathcal{E}_\alpha(\boldsymbol{\theta})$. From Eqn. 25 of that theorem, we have:

$$\frac{\partial \mathcal{E}}{\partial \boldsymbol{\theta}} = \sum_l \mathbb{C}_\alpha \left[\frac{\partial f_{2l}}{\partial \boldsymbol{\theta}}, f_{2l} \right] \quad (110)$$

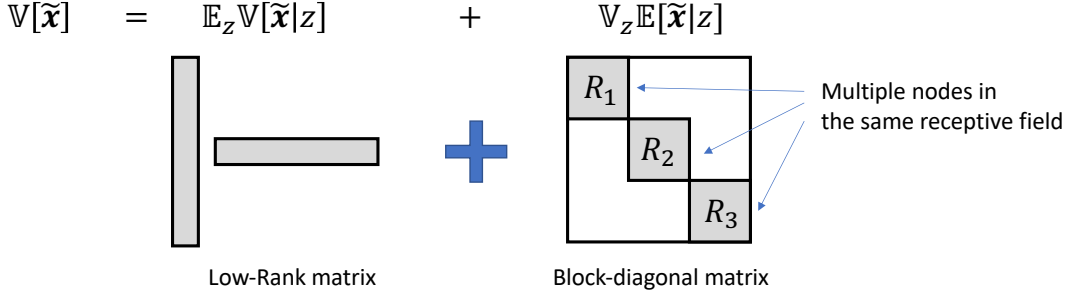
Therefore, for $V = [v_{ik}]$ we have:

$$\dot{\mathbf{v}}_i = \frac{\partial \mathcal{E}}{\partial \mathbf{v}_i} = \sum_l \mathbb{C}_\alpha \left[\frac{\partial f_{2l}}{\partial \mathbf{v}_i}, f_{2l} \right] \quad (111)$$

$$= \mathbb{C}_\alpha [\mathbf{f}_1, \mathbf{v}_i^\top \mathbf{f}_1] \quad (112)$$

$$= \mathbb{C}_\alpha [\mathbf{f}_1, \mathbf{f}_1] \mathbf{v}_i \quad (113)$$

So we have $\dot{\mathbf{v}}_i = \mathbb{C}_\alpha[\mathbf{f}_1] \mathbf{v}_i$, or $\dot{V} = VC_\alpha[\mathbf{f}_1]$.

Figure 6: Decomposition of the variance term $\mathbb{V}[\tilde{\mathbf{x}}]$.

Now we compute $\partial \mathcal{E} / \partial \mathbf{w}_k$:

$$\dot{\mathbf{w}}_k = \frac{\partial \mathcal{E}}{\partial \mathbf{w}_k} = \sum_l \mathbf{C}_\alpha \left[\frac{\partial f_{2l}}{\partial \mathbf{w}_k}, f_{2l} \right] \quad (114)$$

$$= \sum_l \mathbf{C}_\alpha [v_{lk} h'(\mathbf{w}_k^\top \mathbf{x}_k) \mathbf{x}_k, \mathbf{v}_l^\top \mathbf{f}_1] \quad (115)$$

$$= \sum_l v_{lk} \mathbf{C}_\alpha [\tilde{\mathbf{x}}_k, \mathbf{v}_l^\top \mathbf{f}_1] \quad (116)$$

$$= \sum_l v_{lk} \mathbf{C}_\alpha \left[\tilde{\mathbf{x}}_k, \sum_{k'} v_{lk'} h(\mathbf{w}_{k'}^\top \mathbf{x}_{k'}) \right] \quad (117)$$

$$= \sum_{k'} \left(\sum_l v_{lk} v_{lk'} \right) \mathbf{C}_\alpha [\tilde{\mathbf{x}}_k, \tilde{\mathbf{x}}_{k'}] \mathbf{w}_{k'} \quad (118)$$

$$= \sum_{k'} s_{kk'} \mathbf{C}_\alpha [\tilde{\mathbf{x}}_k, \tilde{\mathbf{x}}_{k'}] \mathbf{w}_{k'} \quad (119)$$

where $S = [s_{kk'}] = V^\top V = \sum_l \mathbf{v}_l \mathbf{v}_l^\top$. Let $\mathbf{w} := [\mathbf{w}_1; \dots; \mathbf{w}_K]$ and it leads to the conclusion. When $M > 1$, the proof is similar. \square

C.2 VARIANCE DECOMPOSITION

Let $p_c := \mathbb{P}[z = c]$ be the probability that the latent variable z takes categorical value c .

Lemma 5 (Close-form of variance under Assumption 2). *With Assumption 2, we have*

$$\mathbb{V}[\tilde{\mathbf{x}}] = \text{diag}_k [L_k] + \sum_{c=0}^{C-1} p_c (1 - p_c)^2 \Delta(c) \Delta^\top(c) \quad (120)$$

where $L_k := \mathbb{E}_z \mathbb{V}[\tilde{\mathbf{x}}_k | z] \in \mathbb{R}^{M^d}$ and $\Delta(c) := \mathbb{E}[\tilde{\mathbf{x}} | z = c] - \mathbb{E}[\tilde{\mathbf{x}} | z \neq c] \in \mathbb{R}^{MK^d}$. In particular when $C = 2$, the second term becomes $p_0 p_1 \Delta \Delta^\top$, a rank-1 matrix. Here $\Delta := \Delta(0)$ for brevity.

Proof. Use variance decomposition, we have:

$$\mathbb{V}[\tilde{\mathbf{x}}] = \mathbb{E}_z \mathbb{V}[\tilde{\mathbf{x}} | z] + \mathbb{V}_z \mathbb{E}[\tilde{\mathbf{x}} | z] \quad (121)$$

Remember that $\tilde{\mathbf{x}}_{km}$ is an abbreviation of gated input:

$$\tilde{\mathbf{x}}_{km} := \tilde{\mathbf{x}}_k^{\mathbf{w}_{km}} := \mathbf{x}_k \cdot h'(\mathbf{w}_{km}^\top \mathbf{x}_k) \quad (122)$$

By conditional independence, we have

$$\text{Cov}[\tilde{\mathbf{x}}_{km}, \tilde{\mathbf{x}}_{k'm'} | z] = 0 \quad \forall k \neq k' \quad (123)$$

This is because $\tilde{\mathbf{x}}_{km}$ and $\tilde{\mathbf{x}}_{k'm'}$ are deterministic functions of \mathbf{x}_k and $\mathbf{x}_{k'}$ and thus are also independent of each other.

Let

$$\tilde{\mathbf{x}}_k := \begin{bmatrix} \tilde{\mathbf{x}}_{k1} \\ \tilde{\mathbf{x}}_{k2} \\ \dots \\ \tilde{\mathbf{x}}_{kM} \end{bmatrix} \in \mathbb{R}^{Md} \quad (124)$$

and $L_k := \mathbb{E}_z \mathbb{V}[\tilde{\mathbf{x}}_k | z]$. Then we know that $\mathbb{E}_z \mathbb{V}[\tilde{\mathbf{x}} | z] = \text{diag}_k[L_k]$ is a block diagonal matrix (See Fig. 6).

On the other hand, $\mathbb{V}_z \mathbb{E}[\tilde{\mathbf{x}} | z]$ is a low-rank matrix:

$$\mathbb{V}_z \mathbb{E}[\tilde{\mathbf{x}} | z] = \mathbb{E}_z [(\mathbb{E}[\tilde{\mathbf{x}} | z] - \mathbb{E}[\tilde{\mathbf{x}}])(\mathbb{E}[\tilde{\mathbf{x}} | z] - \mathbb{E}[\tilde{\mathbf{x}}])^\top] \quad (125)$$

Let $\mathbf{q}_c := \mathbb{E}[\tilde{\mathbf{x}} | z = c]$ and $\mathbf{q}_{-c} := \mathbb{E}[\tilde{\mathbf{x}} | z \neq c]$, then we have:

$$\mathbb{E}[\tilde{\mathbf{x}} | z = c] - \mathbb{E}[\tilde{\mathbf{x}}] = \mathbf{q}_c - \sum_c p_c \mathbf{q}_c = (1 - p_c) \left(\mathbf{q}_c - \sum_{c' \neq c} \frac{p_{c'}}{1 - p_c} \mathbf{q}_{c'} \right) \quad (126)$$

$$= (1 - p_c) \left(\mathbf{q}_c - \sum_{c' \neq c} \mathbb{P}[z = c' | z \neq c] \mathbf{q}_{c'} \right) \quad (127)$$

$$= (1 - p_c)(\mathbf{q}_c - \mathbf{q}_{-c}) \quad (128)$$

Therefore, we have:

$$\mathbb{V}_z \mathbb{E}[\tilde{\mathbf{x}} | z] = \mathbb{E}_z [(\mathbb{E}[\tilde{\mathbf{x}} | z] - \mathbb{E}[\tilde{\mathbf{x}}])(\mathbb{E}[\tilde{\mathbf{x}} | z] - \mathbb{E}[\tilde{\mathbf{x}}])^\top] \quad (129)$$

$$= \sum_c p_c (1 - p_c)^2 (\mathbf{q}_c - \mathbf{q}_{-c})(\mathbf{q}_c - \mathbf{q}_{-c})^\top \quad (130)$$

$$= \sum_c p_c (1 - p_c)^2 \Delta(c) \Delta^\top(c) \quad (131)$$

where

$$\Delta(c) := \Delta(c; W) := \mathbf{q}_c - \mathbf{q}_{-c} = \begin{bmatrix} \Delta_{11}(c) \\ \dots \\ \Delta_{KM}(c) \end{bmatrix} \in \mathbb{R}^{KMd} \quad (132)$$

and

$$\Delta_{km}(c) := \Delta_{km}(c; \mathbf{w}_{km}) := \mathbb{E}[\tilde{\mathbf{x}}_{km} | z = c] - \mathbb{E}[\tilde{\mathbf{x}}_{km} | z \neq c] \quad (133)$$

We can see that $\mathbb{V}_z \mathbb{E}[\tilde{\mathbf{x}} | z]$ is at most rank- C , since it is a summation of C rank-1 matrix.

In particular, when $C = 2$, it is clear that $\Delta(0) = -\Delta(1)$ and thus $\Delta(0)\Delta^\top(0) = \Delta(1)\Delta^\top(1)$ and $\sum_c p_c (1 - p_c)^2 = p_0 p_1^2 + p_1 p_0^2 = p_0 p_1$. Hence the conclusion. \square

C.3 GLOBAL MODULATION WHEN $C = 2$ AND $M = 1$

Theorem 5 (Dynamics of \mathbf{w}_k under conditional independence). *When $C = 2$ and $M = 1$, the dynamics of \mathbf{w}_k is given by $(s_k^2$ and $\delta_k \geq 0$ are scalars defined in the proof):*

$$\dot{\mathbf{w}}_k = P_{\mathbf{w}_k}^\perp (s_k^2 A_k(\mathbf{w}_k) + \delta_k \Delta_k \Delta_k^\top) \mathbf{w}_k \quad (15)$$

Proof. Since $M = 1$, each receptive field (RF) R_k only output a single node with output f_k . Let:

$$L_k := \mathbb{E}_z \mathbb{V}[\tilde{\mathbf{x}}_k | z] \quad (134)$$

$$d_k := \mathbf{w}_k^\top L_k \mathbf{w}_k = \mathbb{E}_z \mathbb{V}[f_k | z] \geq 0 \quad (135)$$

$$D := \text{diag}_k[d_k] \quad (136)$$

$$\mathbf{b} := [b_k] := [\mathbf{w}_k^\top \Delta_k] \in \mathbb{R}^K \quad (137)$$

and λ be the maximal eigenvalue of $\mathbb{V}[f_1]$. Here L_k is a PSD matrix and D is a diagonal matrix. Then

$$\mathbb{V}[f_1] = D + p_0 p_1 \mathbf{b} \mathbf{b}^\top \quad (138)$$

is a diagonal matrix plus a rank-1 matrix. Since $p_0 p_1 \mathbf{b} \mathbf{b}^\top$ is always PSD, $\lambda = \lambda_{\max}(\mathbb{V}[\mathbf{f}_1]) \geq \lambda_{\max}(D) = \max_k d_k$. Then using Bunch–Nielsen–Sorensen formula (Bunch et al., 1978), for largest eigenvector \mathbf{s} , we have:

$$s_k = \frac{1}{Z} \frac{b_k}{d_k - \lambda} \quad (139)$$

where λ is the corresponding largest eigenvalue satisfying $1 + p_0 p_1 \sum_k \frac{b_k^2}{d_k - \lambda} = 0$, and $Z = \sqrt{\sum_k \left(\frac{b_k}{d_k - \lambda} \right)^2}$. Note that the above is well-defined, since if $k^* = \arg \max_k d_k$ and $b_{k^*} \neq 0$, then $\lambda > \max_k d_k = d_{k^*}$. So $b_k/(d_k - \lambda)$ won't be infinite.

So we have:

$$\dot{\mathbf{w}}_k = \sum_{k'} s_k s_{k'} \mathbb{C}_\alpha[\tilde{\mathbf{x}}_k, \tilde{\mathbf{x}}_{k'}] \mathbf{w}_{k'} \quad (140)$$

$$\begin{aligned} &= \sum_{k'} s_k s_{k'} (L_k \mathbb{I}(k = k') + p_0 p_1 \Delta_k \Delta_{k'}^\top) \mathbf{w}_{k'} \\ &= s_k^2 \mathbb{V}[\tilde{\mathbf{x}}_k] \mathbf{w}_k + p_0 p_1 s_k \Delta_k \sum_{k' \neq k} s_{k'} \Delta_{k'}^\top \mathbf{w}_{k'} \end{aligned} \quad (141)$$

$$\begin{aligned} &= s_k^2 \mathbb{V}[\tilde{\mathbf{x}}_k] \mathbf{w}_k + \frac{p_0 p_1 b_k}{Z^2 (d_k - \lambda)} \Delta_k \sum_{k' \neq k} \frac{b_{k'}^2}{d_{k'} - \lambda} \\ &= s_k^2 \mathbb{V}[\tilde{\mathbf{x}}_k] \mathbf{w}_k + \delta_k \Delta_k \Delta_k^\top \mathbf{w}_k \\ &= (s_k^2 \mathbb{V}[\tilde{\mathbf{x}}_k] + \delta_k \Delta_k \Delta_k^\top) \mathbf{w}_k \end{aligned} \quad (142)$$

where

$$\delta_k := \frac{p_0 p_1}{Z^2 (\lambda - d_k)} \sum_{k' \neq k} \frac{b_{k'}^2}{\lambda - d_{k'}} \quad (143)$$

Since $\lambda \geq \max_k d_k$, we have $\delta_k \geq 0$ and thus the modulation term is non-negative. Note that since $p_0 p_1 \sum_k \frac{b_k^2}{\lambda - d_k} = 1$, we can also write $\delta_k = 1 - \frac{p_0 p_1 b_k^2}{\lambda - d_k}$. \square

Theorem 6 (Global modulation of attractive basin). *If the structural assumption holds: $A_k(\mathbf{w}_k) = \sum_l g(\mathbf{u}_l^\top \mathbf{w}_k) \mathbf{u}_l \mathbf{u}_l^\top$ with $g(\cdot) > 0$ a linear increasing function and $\{\mathbf{u}_l\}$ orthonormal bases, then for $A_k + c \mathbf{u}_l \mathbf{u}_l^\top$, its attractive basin of $\mathbf{w}_k = \mathbf{u}_l$ is larger than A_k 's for $c > 0$.*

Proof. Since $A_k(\mathbf{w}) = \sum_l g(\mathbf{u}_l^\top \mathbf{w}) \mathbf{u}_l \mathbf{u}_l^\top$, we could write down its dynamics (we omit the projection P_w^\perp for now):

$$\dot{\mathbf{w}} = A_k(\mathbf{w}) \mathbf{w} = \sum_l g(\mathbf{u}_l^\top \mathbf{w}) \mathbf{u}_l \mathbf{u}_l^\top \mathbf{w} \quad (144)$$

Let $y_l(t) := \mathbf{u}_l^\top \mathbf{w}(t)$, i.e., $y_l(t)$ is the projected component of the weight $\mathbf{w}(t)$ onto the l -th direction, i.e., a change of bases to orthonormal bases $\{\mathbf{u}_l\}$, then the dynamics above can be written as

$$\dot{y}_l = g(y_l) y_l \quad (145)$$

which is the same for all l , so we just need to study $\dot{x} = g(x)x$. $g(x) > 0$ is a linear increasing function, so we can assume $g(x) = ax + b$ with $a > 0$. Without loss of generality, we could just set $a = 1$.

Then we just want to analyze the dynamics:

$$\dot{y}_l = (y_l + b_l) y_l, \quad b_l > 0 \quad (146)$$

which also includes the case of $A_k + c \mathbf{u}_l \mathbf{u}_l^\top$, that basically sets $b_l = b + c$. Solving the dynamics leads to the following close-form solution:

$$\frac{y_l(t)}{y_l(t) + b_l} = \frac{y_l(0)}{y_l(0) + b_l} e^{b_l t} \quad (147)$$

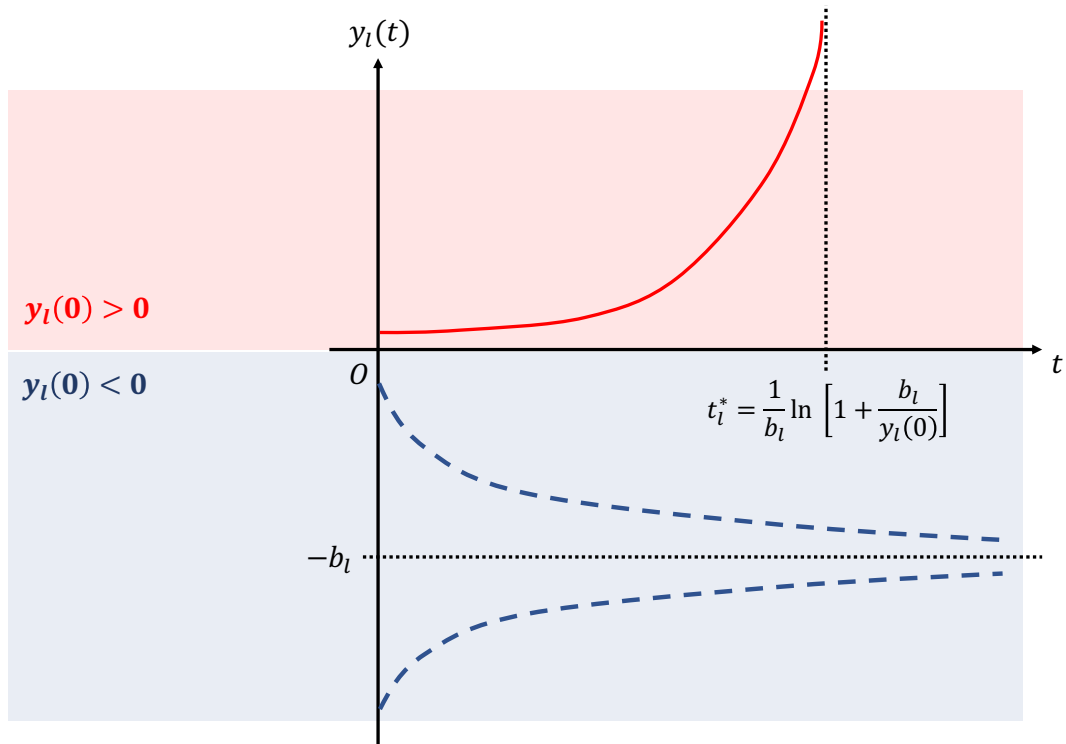


Figure 7: The one-dimensional dynamics (Eqn. 146) ($b_l > 0$). There exists a stable critical point $y_l = -b_l$ and one unstable critical point $y_l = 0$. When $y_l(0) > 0$, the dynamics blows up in finite time $t_l^* = \frac{1}{b_l} \ln \left(1 + \frac{b_l}{y_l(0)} \right)$.

The 1-d dynamics has an unstable fixed points $y_l = 0$ and a stable one $y_l = -b_l < 0$. Therefore, when the initial condition $y_l(0) < 0$, the dynamics will converge to $y_l(+\infty) = -b_l$, which is a finite number. On the other hand, when $y_l(0) > 0$, the dynamics has *finite-time blow-up* Thompson et al. (1990); Goriely & Hyde (1998), i.e., there exists a critical time $t_l^* < +\infty$ so that $y_l(t_l^*) = +\infty$. See Fig. 7.

Note that this finite time blow-up is not physical, since we don't take into consideration of normalization $Z(t)$, which depends on all $y_l(t)$. The real quality to be considered is $\hat{y}_l(t) = \frac{1}{Z(t)}y_l(t)$. Fortunately, we don't need to estimate $Z(t)$ since we are only interested in the ratio:

$$r_{l/l'}(t) := \frac{\hat{y}_l(t)}{\hat{y}_{l'}(t)} = \frac{y_l(t)}{y_{l'}(t)} \quad (148)$$

If for some l and any $l' \neq l$, $r_{l/l'}(t) \rightarrow \infty$, then $y_l(t)$ dominates and $\hat{y}_l(t) \rightarrow 1$, i.e., the dynamics converges to \mathbf{u}_l .

Now our task is to know which initial condition of y_l and b_l makes $r_{l/l'}(t) \rightarrow +\infty$. By comparing the critical time we know which component l shoots up the earliest and that $l^* = \arg \min_l t_l^*$ is the winner, without computing the normalization constant $Z(t)$.

The critical time satisfies

$$\frac{y_l(0)}{y_l(0) + b_l} e^{b_l t_l^*} = 1 \quad (149)$$

so

$$t_l^* = \frac{1}{b_l} \ln \left(1 + \frac{b_l}{y_l(0)} \right) \quad (150)$$

It is clear that when $y_l(0)$ is larger, the critical time t_l^* becomes smaller and the l -th component becomes more advantageous over other components.

For $b_l > 0$, we have:

$$\frac{\partial t_l^*}{\partial b_l} = \frac{1}{b_l^2} \left[\frac{b_l/y_l(0)}{1 + b_l/y_l(0)} - \ln(1 + b_l/y_l(0)) \right] < 0 \quad (151)$$

where the last inequality is due to the fact that $\frac{x}{1+x} < \ln(1+x)$ for $x > 0$. Therefore, larger b_l leads to smaller t_l^* . Since adding $c\mathbf{u}_l\mathbf{u}_l^\top$ with $c > 0$ to A_k increase b_l , it leads to smaller t_l^* and thus increases the advantage of the l -th component.

Therefore, larger b_l and larger $y_l(0)$ both leads to smaller t_l^* . For the same t_l^* , larger b_l can trade for smaller $y_l(0)$, i.e., larger attractive basin. \square

Remark. Special case. We start by assuming only one $\epsilon_l \neq 0$ and all other $\epsilon_{l'} = 0$ for $l' \neq l$, and then we generalize to the case when all $\{\epsilon_l\}$ are real numbers.

To quantify the probability that a random weight initialization leads to convergence of \mathbf{u}_l , we setup some notations. Let the event E_l be ‘‘a random weight initialization of \mathbf{y} leads to $\mathbf{y} \rightarrow \mathbf{e}_l$ ’’, or equivalently $\mathbf{w} \rightarrow \mathbf{u}_l$. Let Y_l be the random variable that instantiates the initial value of $y_l(0)$ due to random weight initialization. Then the convergence event E_l is equivalent to the following: (1) $Y_l > 0$ (so that the l -component has the opportunity to grow), and (2) $Y_l + \epsilon_l$ is the maximum over all $Y_{l'}$ for any $l' \neq l$, where ϵ_l is an advantage (> 0) or disadvantage (< 0) achieved by having larger/smaller b_l due to global modulation (e.g., c). Therefore, we also call ϵ_l the *modulation factor*.

Here we discuss about a simple case that $Y_l \sim U[-1, 1]$ and for $l' \neq l$, Y_l and $Y_{l'}$ are independent. In this case, for a given l , $\max_{l' \neq l} Y_{l'}$ is a random variable that is independent of Y_l , and has cumulative density function (CDF) $F_{\max}(x) := \mathbb{P}[\max_{l' \neq l} Y_{l'} \leq x] = F^{d-1}(x)$, where $F(x)$ is the CDF for Y_l .

Then we have:

$$\mathbb{P}[E_l] = \mathbb{P} \left[Y_l > 0, Y_l + \epsilon_l \geq \max_{l' \neq l} Y_{l'} \right] \quad (152)$$

$$= \int_0^{+\infty} \mathbb{P} \left[\max_{l' \neq l} Y_{l'} \leq Y_l + \epsilon_l \mid Y_l = y_l \right] \mathbb{P}[Y_l = y_l] dy_l \quad (153)$$

$$= \int_0^{+\infty} F^{d-1}(y_l + \epsilon_l) dF(y_l) \quad (154)$$

When $Y_l \sim U[-1, 1]$, $F(x) = \min\{\frac{1}{2}(x+1), 1\}$ has a close form and we can compute the integral:

$$\mathbb{P}[E_l] = \mathbb{P}\left[Y_l > 0, Y_l + \epsilon_l \geq \max_{l' \neq l} Y_{l'}\right] = \begin{cases} \epsilon_l > 1 \\ \frac{1}{d} \left[1 - \left(\frac{1+\epsilon_l}{2}\right)^d\right] + \frac{\epsilon_l}{2} & 0 \leq \epsilon_l \leq 1 \\ \frac{1}{d} \left[\left(1 + \frac{\epsilon_l}{2}\right)^d - \left(\frac{1+\epsilon_l}{2}\right)^d\right] & -1 < \epsilon_l < 0 \end{cases} \quad (155)$$

We can see that the modulation factor ϵ_l plays an important role in deciding the probability that $\mathbf{w} \rightarrow \mathbf{u}_l$:

- **No modulation.** If $\epsilon_l = 0$, then $\mathbb{P}[E_l] \sim \frac{1}{d}$. This means that each dimension of \mathbf{y} has equal probability to be the dominant component after training;
- **Positive modulation.** If $\epsilon_l > 0$, then $\mathbb{P}[E_l] \geq \frac{\epsilon_l}{2}$, and that particular l -th component has much higher probability to become the dominant component, independent of the dimensionality d . Furthermore, the stronger the modulation, the higher the probability becomes.
- **Negative modulation.** Finally, if $\epsilon_l < 0$, since $1 + \epsilon_l/2 < 1$, $\mathbb{P}[E_l] \leq \frac{1}{d}(1 + \frac{\epsilon_l}{2})^d$ decays exponentially w.r.t the dimensionality d .

General case. We then analyze cases if all ϵ_l are real numbers. Let $l^* = \arg \max_l \epsilon_l$ and $c(k)$ be the k -th index of ϵ_l in descending order, i.e., $c(1) = l^*$.

- For $l = c(1) = l^*$, ϵ_l is the largest over $\{\epsilon_l\}$. Since

$$\begin{aligned} \mathbb{P}[E_l] &= \mathbb{P}\left[Y_l \geq 0, Y_l + \epsilon_l \geq \max_{l' \neq l} Y_{l'} + \epsilon_{l'}\right] \\ &\geq \mathbb{P}\left[Y_l \geq 0, Y_l + \epsilon_{c(1)} - \epsilon_{c(2)} \geq \max_{l' \neq l} Y_{l'}\right] \end{aligned}$$

where $\epsilon_{c(1)} - \epsilon_{c(2)}$ is the gap between the largest ϵ_l and second largest ϵ_l . Then this case is similar to positive modulation and thus

$$\mathbb{P}[E_{c(1)}] \geq \frac{1}{2} (\epsilon_{c(1)} - \epsilon_{c(2)}) \quad (156)$$

- For l with rank r (i.e., $c(r) = l$), and any $r' < r$, we have:

$$\begin{aligned} \mathbb{P}[E_l] &= \mathbb{P}\left[Y_l \geq 0, Y_l + \epsilon_l \geq \max_{l' \neq l} Y_{l'} + \epsilon_{l'}\right] \\ &\leq \mathbb{P}\left[Y_l \geq 0, Y_l + \epsilon_l \geq \max_{l': c^{-1}(l') \leq r'} Y_{l'} + \epsilon_{l'}\right] \\ &= \mathbb{P}\left[Y_l \geq 0, Y_l + \epsilon_l - \epsilon_{c(r')} \geq \max_{l': c^{-1}(l') \leq r'} Y_{l'} + \epsilon_{l'} - \epsilon_{c(r')}\right] \\ &\leq \mathbb{P}\left[Y_l \geq 0, Y_l + \epsilon_l - \epsilon_{c(r')} \geq \max_{l': c^{-1}(l') \leq r'} Y_{l'}\right] \end{aligned}$$

Then it reduces to the case of negative modulation. Therefore, we have:

$$\mathbb{P}[E_{c(r)}] \leq \min_{r' < r} \frac{1}{r' + 1} \left(1 - \frac{\epsilon_{c(r')} - \epsilon_{c(r)}}{2}\right)^{r'+1} \quad (157)$$

and the probability is exponentially small if r is large, i.e., ϵ_l ranks low.

C.4 FUNDAMENTAL LIMITATION OF LINEAR MODELS

Theorem 4 (Gradient Colinearity in linear networks). *With linear activation, W follows the dynamics:*

$$\dot{\mathbf{w}}_{km} = s_{km} \mathbf{b}_k(W, V) \quad (14)$$

where $\mathbf{b}_k(W, V) := \mathcal{C}_\alpha \left[\mathbf{x}_k, \sum_{k', m'} s_{k'm'} \mathbf{w}_{k'm'}^\top \mathbf{x}_{k'} \right]$ is a linear function w.r.t. W . As a result, (1) $\dot{\mathbf{w}}_{km}$ are co-linear over m , and (2) If $s_{km} \neq 0$, from any critical point with distinct $\{\mathbf{w}_{km}\}$, there exists a path of critical points to identical weights ($\mathbf{w}_{km} = \mathbf{w}_k$).

Proof. In the linear case, we have $\tilde{\mathbf{x}}_{km} = \mathbf{x}_k$ since there is no gating and all M shares the same input \mathbf{x}_k . Therefore, we can write down the dynamics of \mathbf{w}_{km} as the following:

$$\dot{\mathbf{w}}_{km} = \sum_{k',m'} s_{km,k'm'} \mathbb{C}_\alpha[\tilde{\mathbf{x}}_{km}, \tilde{\mathbf{x}}_{k'm'}] \mathbf{w}_{k'm'} \quad (158)$$

$$= \sum_{k',m'} s_{km,k'm'} \mathbb{C}_\alpha[\mathbf{x}_k, \mathbf{x}_{k'}] \mathbf{w}_{k'm'} \quad (159)$$

Now we use the fact that the top-level learns fast so that $s_{km,k'm'} = s_{km}s_{k'm'}$, which gives:

$$\dot{\mathbf{w}}_{km} = s_{km} \sum_{k',m'} s_{k'm'} \mathbb{C}_\alpha[\mathbf{x}_k, \mathbf{x}_{k'}] \mathbf{w}_{k'm'} \quad (160)$$

$$= s_{km} \mathbb{C}_\alpha \left[\mathbf{x}_k, \sum_{k',m'} s_{k'm'} \mathbf{w}_{k'm'}^\top \mathbf{x}_{k'} \right] \quad (161)$$

Let $\mathbf{b}_k(W, V) := \mathbb{C}_\alpha \left[\mathbf{x}_k, \sum_{k',m'} s_{k'm'} \mathbf{w}_{k'm'}^\top \mathbf{x}_{k'} \right]$ be a linear function of W , and we have:

$$\dot{\mathbf{w}}_{km} = s_{km} \mathbf{b}_k(W, V) \quad (162)$$

Since \mathbf{b}_k is independent of m , all $\dot{\mathbf{w}}_{km}$ are co-linear.

For the second part, first all if W^* is a critical point, we have the following two facts:

- Since there exists m so that $s_{km} \neq 0$, we know that $\mathbf{b}_k(W^*) = 0$;
- If W^* contains two distinct filters $\mathbf{w}_{k1} = \boldsymbol{\mu}_1 \neq \mathbf{w}_{k2} = \boldsymbol{\mu}_2$ covering the same receptive field R_k , then by symmetry of the weights, W'^* in which $\mathbf{w}_{k1} = \boldsymbol{\mu}_2$ and $\mathbf{w}_{k2} = \boldsymbol{\mu}_1$, is also a critical point.

Then for any $c \in [0, 1]$, since $\mathbf{b}_k(W)$ is linear w.r.t. W , for the linear combination $W^c := cW^* + (1-c)W'^*$, we have:

$$\mathbf{b}_k(W^c) = \mathbf{b}_k(cW^* + (1-c)W'^*) = c\mathbf{b}_k(W^*) + (1-c)\mathbf{b}_k(W'^*) = 0 \quad (163)$$

Therefore, W^c is also a critical point, in which $\mathbf{w}_{k1} = c\boldsymbol{\mu}_1 + (1-c)\boldsymbol{\mu}_2$ and $\mathbf{w}_{k2} = (1-c)\boldsymbol{\mu}_1 + c\boldsymbol{\mu}_2$. In particular when $c = 1/2$, $\mathbf{w}_{k1} = \mathbf{w}_{k2}$. Repeating this process for different m , we could finally reach a critical point in which all $\mathbf{w}_{km} = \mathbf{w}_k$. \square

D ANALYSIS OF BATCH NORMALIZATION

From the previous analysis of global modulation, it is clear that the weight updating can be much slower for RF with small d_k , due to the factor $\frac{1}{\lambda - d_k}$ in both s_k^2 (Eqn. 139) and β_k (Eqn. 143) and the fact that $\lambda \geq \max_k d_k$. This happens when the variance of each receptive fields varies a lot (i.e., some d_k are large while others are small). In this case, adding BatchNorm at each node alleviates this issue, as shown below.

We consider BatchNorm right after \mathbf{f} : $f_k^{\text{bn}}[i] = (f_k[i] - \mu_k)/\sigma_k$, where μ_k and σ_k are the batch statistics computed from BatchNorm on all $2N$ samples in a batch:

$$\mu_k := \frac{1}{2N} \sum_i f_k[i] + f_k[i'] \quad (164)$$

$$\sigma_k^2 := \frac{1}{2N} \sum_i (f_k[i] - \mu_k)^2 + (f_k[i'] - \mu_k)^2 \quad (165)$$

When $N \rightarrow +\infty$, we have $\mu_k \rightarrow \mathbb{E}[f_k]$ and $\sigma_k^2 \rightarrow \mathbb{V}[f_k] = \mathbf{w}_k^\top \mathbb{V}[\tilde{\mathbf{x}}_k] \mathbf{w}_k$.

Let $\tilde{\mathbf{x}}_k^{\text{bn}} := \sigma_k^{-1} \tilde{\mathbf{x}}_k$ and $\tilde{\mathbf{x}}^{\text{bn}} := \begin{bmatrix} \tilde{\mathbf{x}}_1^{\text{bn}} \\ \tilde{\mathbf{x}}_2^{\text{bn}} \\ \dots \\ \tilde{\mathbf{x}}_K^{\text{bn}} \end{bmatrix}$. When computing gradient through BatchNorm layer, we consider the following variant:

Definition 5 (mean-backprop BatchNorm). *When computing backpropagated gradient through BatchNorm, we only backprop through μ_k .*

This leads to a model dynamics that has a very similar form as Lemma 4:

Lemma 6 (Dynamics with mean-backprop BatchNorm). *With mean-backprop BatchNorm (Def. 5), the dynamics is:*

$$\dot{V} = VC_\alpha[\mathbf{f}_1^{\text{bn}}], \quad \dot{\mathbf{w}} = [(S \otimes \mathbf{1}_d \mathbf{1}_d^\top) \circ C_\alpha[\tilde{\mathbf{x}}^{\text{bn}}]] \mathbf{w} \quad (166)$$

Proof. The proof is similar to Lemma 4. For \dot{V} it is the same by replacing \mathbf{f}_1 with \mathbf{f}_1^{bn} , which is the input to the top layer.

For $\dot{\mathbf{w}}$, similarly we have:

$$\dot{\mathbf{w}}_k = \frac{\partial \mathcal{E}}{\partial \mathbf{w}_k} = \sum_l C_\alpha \left[\frac{\partial f_{2l}}{\partial \mathbf{w}_k}, f_{2l} \right] \quad (167)$$

$$= \sum_l C_\alpha \left[v_{lk} \frac{\partial f_{1k}^{\text{bn}}}{\partial \mathbf{w}_k}, \sum_{k'} v_{lk'} f_{1k'}^{\text{bn}} \right] \quad (168)$$

$$= \sum_l C_\alpha \left[v_{lk} \frac{\partial f_{1k}^{\text{bn}}}{\partial \mathbf{w}_k}, \sum_{k'} v_{lk'} (f_{1k'} - \mu_{k'}) \sigma_{k'}^{-1} \right] \quad (169)$$

Note that $C_\alpha[\cdot, \mu_{k'} \sigma_{k'}^{-1}] = 0$ since $\mu_{k'}$ and $\sigma_{k'}$ are statistics of the batch and is constant. On the other hand, for $\partial f_{1k}^{\text{bn}} / \partial \mathbf{w}_k$, we have:

$$\frac{\partial f_{1k}^{\text{bn}}}{\partial \mathbf{w}_k} = \frac{1}{\sigma_k} \left(\frac{\partial f_{1k}}{\partial \mathbf{w}_k} - \frac{\partial \mu_k}{\partial \mathbf{w}_k} \right) - \frac{f_{1k}^{\text{bn}}}{\sigma_k} \frac{\partial \sigma_k}{\partial \mathbf{w}_k} \quad (170)$$

Note that

$$\frac{\partial \mu_k}{\partial \mathbf{w}_k} = \mathbb{E}_{\text{sample}}[\tilde{\mathbf{x}}_k] \quad (171)$$

where $\mathbb{E}_{\text{sample}}[\cdot]$ is the sample mean, which is a constant over the batch. Therefore $C_\alpha[\cdot, \partial \mu_k / \partial \mathbf{w}_k] = 0$. For mean-backprop BatchNorm, since the gradient didn't backpropagate through the variance, the second term is simply zero. Therefore, we have:

$$\dot{\mathbf{w}}_k = \sum_l C_\alpha \left[v_{lk} \sigma_k^{-1} \frac{\partial f_{1k}}{\partial \mathbf{w}_k}, \sum_{k'} v_{lk'} f_{1k'} \sigma_{k'}^{-1} \right] \quad (172)$$

$$= \sum_l C_\alpha \left[v_{lk} \sigma_k^{-1} \tilde{\mathbf{x}}_k, \sum_{k'} v_{lk'} \mathbf{w}_{k'}^\top \tilde{\mathbf{x}}_{k'} \sigma_{k'}^{-1} \right] \quad (173)$$

$$= \sum_{k'} s_{kk'} C_\alpha[\sigma_k^{-1} \tilde{\mathbf{x}}_k, \sigma_{k'}^{-1} \tilde{\mathbf{x}}_{k'}] \mathbf{w}_{k'} \quad (174)$$

Let $\tilde{\mathbf{x}}_k^{\text{bn}} := \sigma_k^{-1} \tilde{\mathbf{x}}_k$ and $\tilde{\mathbf{x}}^{\text{bn}} := \begin{bmatrix} \tilde{\mathbf{x}}_1^{\text{bn}} \\ \tilde{\mathbf{x}}_2^{\text{bn}} \\ \dots \\ \tilde{\mathbf{x}}_{K'}^{\text{bn}} \end{bmatrix} \in \mathbb{R}^{Kd}$. The conclusion follows. \square

Corollary 3 (Dynamics of \mathbf{w}_k under conditional independence and BatchNorm). *Let*

$$A_k^{\text{bn}} := \mathbb{V}[\tilde{\mathbf{x}}_k^{\text{bn}}] = \sigma_k^{-2} A_k \quad (175)$$

$$d_k^{\text{bn}} := \sigma_k^{-2} d_k \quad (176)$$

$$\Delta_k^{\text{bn}} := \sigma_k^{-1} \Delta_k = \mathbb{E}[\tilde{\mathbf{x}}_k^{\text{bn}} | z = 1] - \mathbb{E}[\tilde{\mathbf{x}}_k^{\text{bn}} | z = 0] \quad (177)$$

and λ^{bn} be the maximal eigenvalue of $\mathbb{V}[\mathbf{f}_1^{\text{bn}}]$. Then we have

- (I) $\lambda^{\text{bn}} \geq \max_k d_k^{\text{bn}}$;

- (2) For λ^{bn} , the associated unit eigenvector is

$$\mathbf{s}^{\text{bn}} := \frac{1}{Z^{\text{bn}}} \left[\frac{\mathbf{w}_k^\top \Delta_k^{\text{bn}}}{\lambda^{\text{bn}} - d_k^{\text{bn}}} \right] \in \mathbb{R}^K,$$

where Z^{bn} is the normalization constant;

- (3) the dynamics of \mathbf{w}_k is given by:

$$\dot{\mathbf{w}}_k = [(s_k^{\text{bn}})^2 A_k^{\text{bn}} + \delta_k^{\text{bn}} (\Delta_k^{\text{bn}})(\Delta_k^{\text{bn}})^\top] \mathbf{w}_k \quad (178)$$

where

$$\delta_k^{\text{bn}} := \frac{p_0 p_1}{(Z^{\text{bn}})^2 (\lambda^{\text{bn}} - d_k^{\text{bn}})} \sum_{k' \neq k} \frac{(\mathbf{w}_{k'}^\top \Delta_{k'}^{\text{bn}})^2}{\lambda^{\text{bn}} - d_{k'}^{\text{bn}}} \geq 0 \quad (179)$$

Proof. Similar to Theorem 5. □

Remarks In the presence of BatchNorm, Lemma 5 still holds, since it only depends on the generative structure of the data. Therefore, we have

$$\sigma_k^2 \rightarrow \mathbb{V}[f_k] = \mathbf{w}_k^\top \mathbb{V}[\tilde{\mathbf{x}}_k] \mathbf{w}_k = d_k + p_0 p_1 (\mathbf{w}_k^\top \Delta_k)^2$$

and thus

$$d_k^{\text{bn}} = \sigma_k^{-2} d_k \rightarrow \frac{d_k}{d_k + p_0 p_1 (\mathbf{w}_k^\top \Delta_k)^2} = \frac{1}{1 + p_0 p_1 (\mathbf{w}_k^\top \Delta_k / \sqrt{d_k})^2}$$

becomes more uniform. This is because $d_k := \mathbf{w}_k^\top L_k \mathbf{w}_k = \mathbb{E}_z \mathbb{V}[f_k | z] \geq 0$ (Eqn. 135) is approximately the variance of f_k , and thus $\mathbf{w}_k^\top \Delta_k / \sqrt{d_k}$ is normalized across different receptive field k , reducing the effect of magnitude of the input.

Since there is no much variation within $\{d_k^{\text{bn}}\}$, $\lambda^{\text{bn}} - d_k^{\text{bn}}$ becomes almost constant across different receptive field R_k and won't lead to slowness of feature learning.

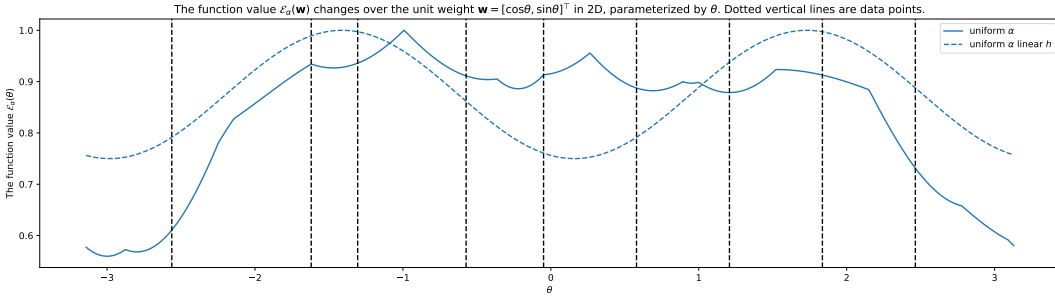


Figure 8: Local optima of objective $\mathcal{E}_\alpha(\mathbf{w})$ under uniform $\alpha := 1$. $\mathbf{w} = [\cos \theta, \sin \theta]^\top$ is parameterized by $\theta \in [-\pi, \pi]$, and each vertical dotted line is a data point. Dotted line is \mathcal{E}_α with linear activation $h(x) = x$, and solid line is using ReLU activation $h(x) = \max(x, 0)$. Both lines are scaled so that their maximal value is 1.

E ADDITIONAL EXPERIMENTS

E.1 VISUALIZATION OF LOCAL OPTIMA IN 1-LAYER SETTING

We added a simple experiment to visualize the local maxima of $\mathcal{E}_\alpha(\mathbf{w}) := \frac{1}{2} \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})]$, when $\mathbf{w} = [\cos \theta, \sin \theta]^\top$ is a 2D unit vector parameterized by θ , and h is ReLU activation. For simplicity, here we use a uniform $\alpha := 1$.

We put a few data points $\{\mathbf{x}_i\}$ on the unit circle, which are also parameterized by θ . The data points are located at $\{-\frac{4\pi}{5}, -\frac{\pi}{2}, -\frac{2\pi}{5}, -\frac{\pi}{6}, 0, \frac{\pi}{5}, \frac{2\pi}{5}, \frac{3\pi}{5}, \frac{4\pi}{5}\}$ and no data augmentation is used. The objective function $\mathcal{E}_\alpha(\theta)$ is plotted in Fig. 8.

From the figure, we can see many local maxima (≥ 8) caused by nonlinearity (solid line), much more than $2 \times 2 = 4$, the maximal possible number of local maxima counting all PCA components in 2D case (i.e., $\pm\phi_1$ and $\pm\phi_2$, where ϕ_1 and ϕ_2 are orthogonal PCA directions in this 2D example). Moreover, unlike PCA directions, these local optima are not orthogonal to each other.

On the other hand, in the linear case (dotted line), the curve is much smoother. There are only two local maxima corresponding to $\pm\phi_1$, where ϕ_1 is the largest PCA eigenvector.

E.2 2-LAYER SETTING

We also do more experiments on the 2-layer setting, to further verify our theoretical findings.

Overall matching score $\bar{\chi}_+$ and overall irrelevant-matching score $\bar{\chi}_-$. As defined in the main text (Eqn. 16), the matching score $\chi_+(R_k)$ is the degree of matching between learned weights and the embeddings of the subset R_k^g of tokens that are allowed in the global patterns at each receptive field R_k . And the overall matching score $\bar{\chi}_+$ is $\bar{\chi}_+$ averaged over all receptive fields:

$$\chi_+(R_k) = \frac{1}{P} \sum_{a \in R_k^g} \max_m \frac{\mathbf{w}_{km}^\top \mathbf{u}_a}{\|\mathbf{w}_{km}\|_2 \|\mathbf{u}_a\|_2}, \quad \bar{\chi}_+ = \frac{1}{K} \sum_k \chi_+(R_k) \quad (180)$$

Similarly, we can also define *irrelevant-matching score* $\chi_-(R_k)$ which is the degree of matching between learned weights and the embeddings of the tokens that are NOT in the subset R_k^g at each receptive field R_k . And the overall *irrelevant-matching score* $\bar{\chi}_-$ is defined similarly.

$$\chi_-(R_k) = \frac{1}{P} \sum_{a \notin R_k^g} \max_m \frac{\mathbf{w}_{km}^\top \mathbf{u}_a}{\|\mathbf{w}_{km}\|_2 \|\mathbf{u}_a\|_2}, \quad \bar{\chi}_- = \frac{1}{K} \sum_k \chi_-(R_k) \quad (181)$$

Ideally, we want to see high overall matching score $\bar{\chi}_+$ and low overall irrelevant-matching score $\bar{\chi}_-$, which means that the important patterns in R_k^g (i.e., the patterns that are allowed in the global generators) are learned, but noisy patterns that are not part of the global patterns (i.e., the generators) are not learned. Fig. 9 shows that this indeed is the case.

Non-uniformity ζ and how BatchNorm interacts with it. When the scale of input data varies a lot, BatchNorm starts to matter in discovering features with low magnitude (Sec. D). To model the

scale non-uniformity, we set $\|\mathbf{u}_a\|_2 = \zeta$ for $\lfloor d/2 \rfloor$ tokens and $\|\mathbf{u}_a\|_2 = 1/\zeta$ for the remaining tokens. Larger ζ corresponds to higher non-uniformity across inputs.

Fig. 11 shows that BN with ReLU activations handles large non-uniformity (large ζ) very well, compared to the case without BN. Specifically, BN yields higher $\bar{\chi}_+$ in the presence of high non-uniformity (e.g., $\zeta = 10$) when the network is over-parameterized ($\beta > 1$) and there are multiple candidates per R_k ($P > 1$), a setting that is likely to hold in real-world scenarios.

Note that in the real-world scenario, features from different channels/modalities indeed will have very different scales, and some local features that turn out to be super important to global features, can have very small scale. In such cases, normalization techniques (e.g., BatchNorm) can be very useful and our formulation justifies it in a mathematically consistent way.

Selectively Backpropagating μ_k and σ_k^2 in BatchNorm. In our analysis of BatchNorm, we assume that gradient backpropagating the mean statistics μ_k , but not variance σ_k^2 (see Def. 5). Note that this is different from regular BatchNorm, in which both μ_k and σ_k^2 get backpropagated gradients. Therefore, we test how this modified BN affects the matching score $\bar{\chi}_+$: we change whether μ_k and σ_k^2 gets backpropagated gradients, while the forward pass remains the same, yielding the four following variants:

$$\begin{aligned}
 f_k^{\text{bn}}[i] &:= \frac{f_k^{\text{bn}}[i] - \mu_k}{\sigma_k} && \text{(Vanilla BatchNorm)} \\
 f_k^{\text{bn}}[i] &:= \frac{f_k^{\text{bn}}[i] - \text{stop-gradient}(\mu_k)}{\sigma_k} && \text{(BatchNorm with backpropated } \sigma_k) \\
 f_k^{\text{bn}}[i] &:= \frac{f_k^{\text{bn}}[i] - \mu_k}{\text{stop-gradient}(\sigma_k)} && \text{(BatchNorm with backpropated } \mu_k) \\
 f_k^{\text{bn}}[i] &:= \frac{f_k^{\text{bn}}[i] - \text{stop-gradient}(\mu_k)}{\text{stop-gradient}(\sigma_k)} && \text{(BatchNorm without backpropating statistics)}
 \end{aligned}$$

As shown in Tbl. 1, it is interesting to see that if σ_k^2 is not backpropagated, then the matching score $\bar{\chi}_+$ is actually better. This justifies our BN variant.

Quadratic versus InfoNCE loss. Fig. 10 shows that quadratic loss (constant pairwise importance α) shows worse matching score than InfoNCE. A high-level intuition is that InfoNCE dynamically adjusts the pairwise importance α (i.e., the focus of different sample pairs) during training to focus on the most important sample pairs, which makes learning patterns more efficient. We leave a comprehensive study for future work.

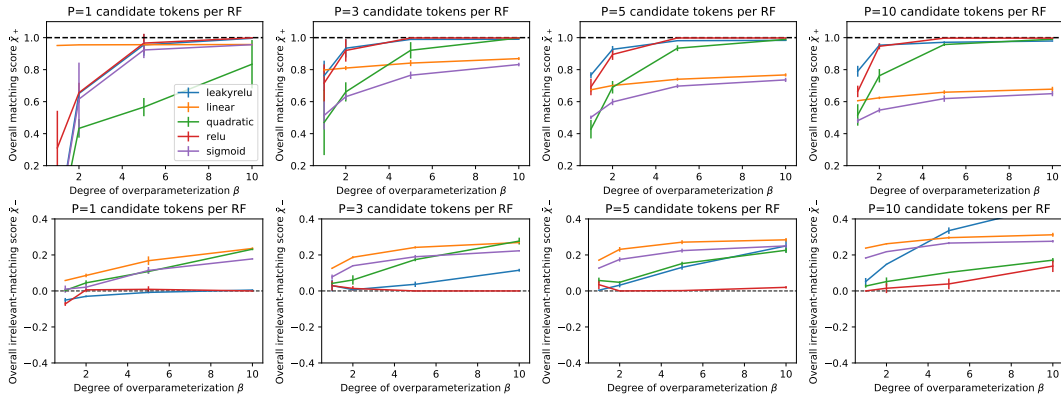


Figure 9: Overall matching score $\bar{\chi}_+$ (Eqn. 180, the top row) and irrelevant-matching score $\bar{\chi}_-$ (Eqn. 181, the bottom row). This is an extended version of Fig. 4. (a) When $P = 1$, linear model works well regardless of the degree of over-parameterization β , while ReLU model requires large over-parameterization to perform well; (b) When each R_k has multiple local patterns that are related to the global patterns ($P > 1$) related to generators, ReLU models can capture diverse patterns better than linear ones in the over-parameterization region $\beta > 1$ and stay focus on relevant local patterns that are related to the global patterns (i.e., low $\bar{\chi}_-$). Among all activations (homogeneous or non-homogeneous), ReLU shows its strength by achieving the lowest irrelevant-matching score $\bar{\chi}_-$. In contrast, linear models are much less affected by over-parameterization. Each setting is repeated 3 times and mean/standard derivations are reported.

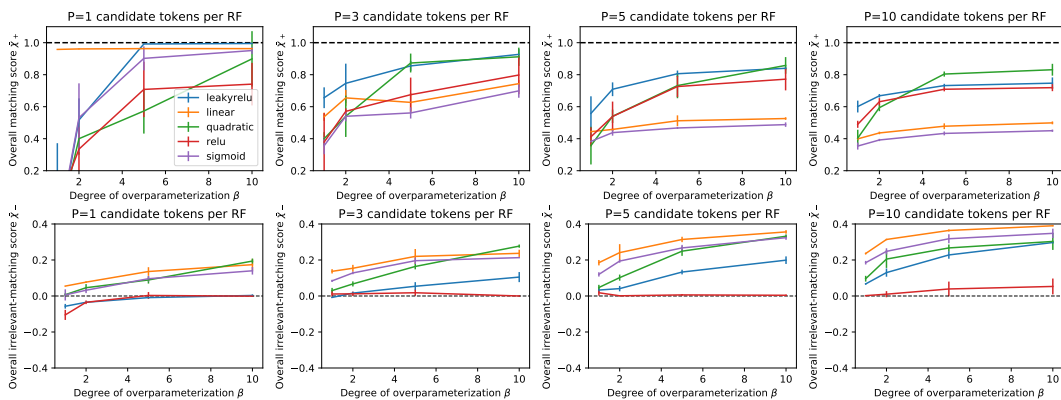


Figure 10: Overall matching score $\bar{\chi}_+$ (top row) and overall irrelevant-matching score $\bar{\chi}_-$ (Eqn. 16, bottom row) using **quadratic** loss function rather than InfoNCE. The result using InfoNCE is shown in Fig. 9, with all experiments setting being the same, except for the loss function. While we see similar trends as in Sec. 5.1, quadratic loss is not as effective as InfoNCE in feature learning.

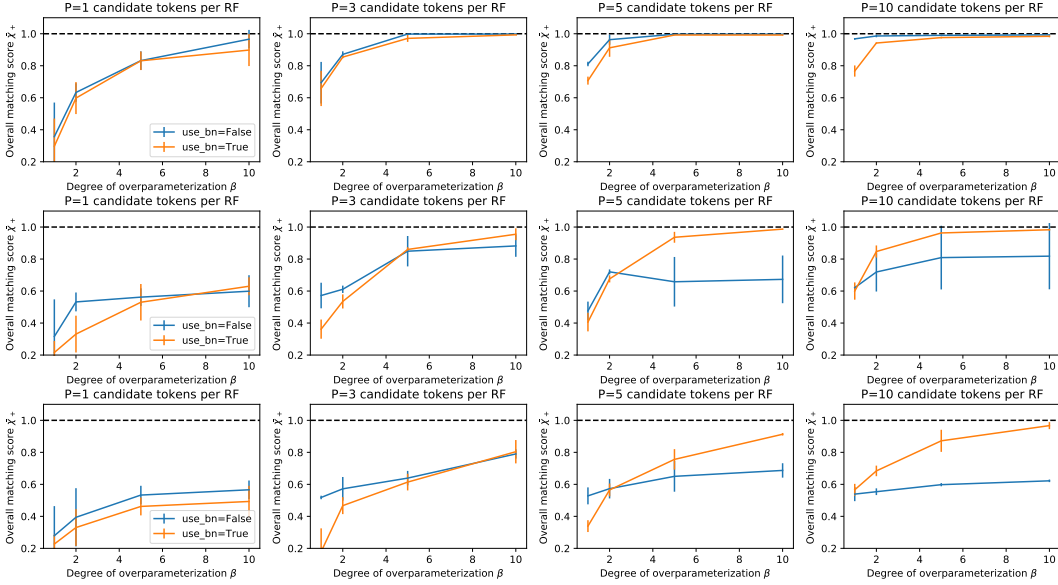


Figure 11: The effect of BatchNorm (BN) with ReLU activation in the presence of non-uniformity ζ of the input data. The non-uniformity is set to be $\zeta = 2$ (top), $\zeta = 5$ (middle) and $\zeta = 10$ (bottom). For small non-uniformity, BN doesn't help much. For larger nonuniformity, BN yields better matching score $\bar{\chi}_+$ in the over-parameterization region (large β) and multiple tokens per RF (large P).

β	P	μ_k no backprop		μ_k backprop	
		σ_k^2 no backprop	σ_k^2 backprop	σ_k^2 no backprop	σ_k^2 backprop
1	1	0.31 ± 0.24	0.23 ± 0.06	0.30 ± 0.24	0.23 ± 0.06
	3	0.65 ± 0.09	-0.03 ± 0.01	0.64 ± 0.07	0.24 ± 0.11
	5	0.61 ± 0.06	-0.00 ± 0.00	0.62 ± 0.02	0.38 ± 0.04
	10	0.66 ± 0.06	0.53 ± 0.05	0.70 ± 0.03	0.56 ± 0.04
2	1	0.36 ± 0.13	0.27 ± 0.06	0.63 ± 0.11	0.33 ± 0.12
	3	0.78 ± 0.01	0.00 ± 0.01	0.80 ± 0.02	0.41 ± 0.07
	5	0.78 ± 0.02	0.22 ± 0.21	0.77 ± 0.07	0.56 ± 0.02
	10	0.83 ± 0.08	0.72 ± 0.06	0.80 ± 0.04	0.71 ± 0.05
5	1	0.63 ± 0.06	0.43 ± 0.12	0.67 ± 0.06	0.46 ± 0.06
	3	0.90 ± 0.06	0.45 ± 0.07	0.90 ± 0.03	0.60 ± 0.08
	5	0.90 ± 0.01	0.72 ± 0.01	0.88 ± 0.06	0.74 ± 0.02
	10	0.88 ± 0.01	0.88 ± 0.04	0.92 ± 0.03	0.90 ± 0.03
10	1	0.63 ± 0.12	0.37 ± 0.15	0.70 ± 0.17	0.46 ± 0.15
	3	0.95 ± 0.05	0.72 ± 0.06	0.94 ± 0.05	0.80 ± 0.03
	5	0.98 ± 0.02	0.84 ± 0.05	0.96 ± 0.06	0.92 ± 0.01
	10	0.90 ± 0.01	0.97 ± 0.02	0.89 ± 0.03	0.97 ± 0.02

Table 1: The effect of backpropagating different BN statistics under nonuniformity $\zeta = 10$. Backpropagating the gradient through the sample mean μ_k but not the sample variance σ_k^2 gives overall good matching score $\bar{\chi}_+$, justifying our setting of mean-backprop BatchNorm (Def. 5).

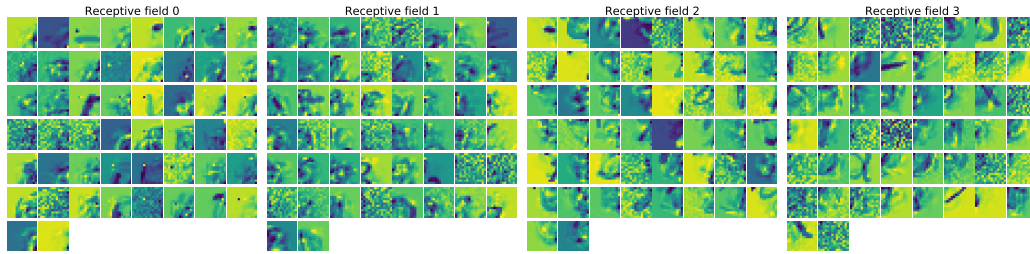


Figure 12: Learned filters (of the 4 disjoint receptive field) in MNIST dataset without using augmentation. The 4 receptive fields corresponds to upper left (0), upper right (1), bottom left (2) and bottom right (3) part of the input image.

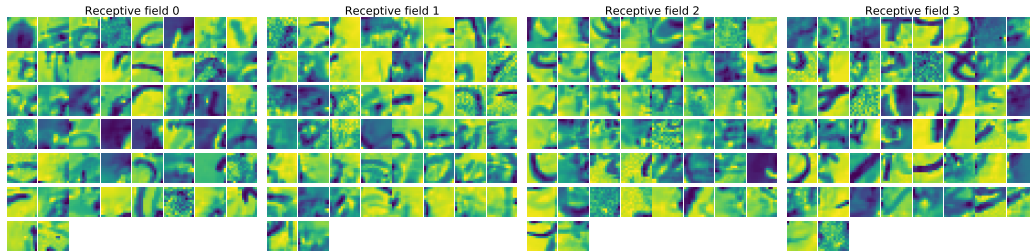


Figure 13: Same as Fig. 12 but with data augmentation during contrastive learning. The learned filters are smoother. According to the tentative theory in Sec. B.5, data augmentation removes some of the local optima.

E.3 MNIST EXPERIMENTS WITH 2-LAYER SETTING

We also run the same 2-layer network (as in Sec. E.2) on MNIST Deng (2012) dataset. In its training, the MNIST dataset consists of 50,000 images, each with the size of 28 by 28. We split the 28-by-28 images into 4 disjoint receptive fields, each with a size of 14 by 14, just like Fig. 2(b). In each region, we vectorize the receptive field into $14 * 14 = 196$ dimensional vector. We use smaller batchsize (8), since there are only 10 true classes in MNIST, and the probability that two samples from the same class are incorrectly treated as negative pair increases with large batchsizes. We train for 50000 minibatches and start new pass (epoch) of the dataset if needed.

Fig 12 shows the initial results. We could see that indeed filters in each receptive field capture a diverse set of pattern of the input data points (e.g., part of the digits). Furthermore, with additional data augmentation (i.e., random cropping and resizing with `transforms.RandomResizedCrop((28,28), scale=(0.5, 1.0), ratio=(0.5, 1.5))` in PyTorch), the resulting learned patterns becomes smoother with much weaker high-frequency components. This is because the augmentation implicitly removes some of the local optima (Sec. B.5), leaving more informative local optima.

F CONTEXT AND MOTIVATIONS OF THEOREMS

Here are the motivations and intuitions behind each major theoretical results:

- First of all, Lemma 2 and Corollary 1 try to make connection between a relatively new (and somehow abstract) concept (i.e., contrastive covariance $\mathbb{C}_\alpha[\cdot]$) and a well-known concept (i.e., regular covariance $\mathbb{V}[\cdot]$). This will also enable us to leverage existing properties of covariance, in order to deepen our understanding of this concept, which seems to play an important role in contrastive learning (CL).
- After that, we mainly focus on studying the CL energy function \mathcal{E}_α , which can be represented in terms of the contrastive covariance $\mathbb{C}_\alpha[\cdot]$. One important property of the energy function is whether there exists any local optima and what are the properties of these local optima, since these local optima are the final destinations that network weights will converge into. Previous works in landscape analysis often talk about local optima in neural networks as abstract objects in high-dimensional space, but here, we would like to make them as concrete as possible.
- For such analysis, we always start from the simplest case (e.g., one-layer). Therefore, naturally we have Lemma 3 that characterizes critical points (as a superset of local optima), a few examples in Sec. 3.2, properties of these critical points and when they become local optima in Sec. 3. Finally, Appendix B.5 further gives a preliminary study on how the data augmentation affects the distribution of the local optima.
- Then we extend our analysis to 2-layer setting. The key question is to study the additional benefits of 2-layer network compared to K independent 1-layer cases. Here the assumption of *disjoint* receptive fields is to make sure there is an apple-to-apple comparison, otherwise additional complexity would be involved, e.g., overlapping receptive fields. As demonstrated in Theorem 5 and Theorem 6, we find the effect of *global modulation* in 2-layer case, which clearly tells that the interactions across different receptive fields lead to additional terms in the dynamics that favors patterns related to latent variable z that leads to conditional independence across the disjointed receptive fields.
- As a side track, in 2-layer case, we also have Theorem 4 that shows linear activation does not learn distinct features, which is consistent with 1-layer case that linear activation $h(x) = x$ only gives to maximal PCA directions (Sec. 6).

G OTHER LEMMAS

Lemma 7 (Bound of $1 - d_t$). *Define*

$$\mu(\mathbf{w}) := \frac{1 + c(\mathbf{w})}{2c^2(\mathbf{w})} \left(\frac{\lambda_2(A(\mathbf{w}(t)))}{\lambda_1(A(\mathbf{w}(t)))} \right)^2 = \frac{1 + c(\mathbf{w})}{2c^2(\mathbf{w})} \left[1 - \frac{\lambda_{\text{gap}}(A(t))}{\lambda_1(A(t))} \right]^2 \geq 0 \quad (182)$$

and $\mu_t := \mu(\mathbf{w}(t))$. If $c_t > 0$ and $\lambda_1(t) > 0$, then $1 - d_t \leq \mu_t(1 - c_t)$.

Proof. We could write d_t :

$$d_t = \frac{\phi_1^\top(t) \tilde{\mathbf{w}}(t+1)}{\|\tilde{\mathbf{w}}(t+1)\|_2} = \frac{\lambda_1(t) \phi_1^\top(t) \mathbf{w}(t)}{\sqrt{\sum_i \lambda_i^2(t) (\phi_i^\top(t) \mathbf{w}(t))^2}} \quad (183)$$

$$\geq \frac{\lambda_1(t) c_t}{\sqrt{\lambda_1^2(t) c_t^2 + \lambda_2^2(t) (1 - c_t^2)}} = \frac{1}{\sqrt{1 + \left(\frac{\lambda_2(t)}{\lambda_1(t)} \right)^2 \left(\frac{1}{c_t^2} - 1 \right)}} \quad (184)$$

$$= \left[1 + \left(\frac{\lambda_2(t)}{\lambda_1(t)} \right)^2 \left(\frac{1}{c_t^2} - 1 \right) \right]^{-1/2} \quad (185)$$

$$\geq 1 - \frac{1}{2} \left(\frac{\lambda_2(t)}{\lambda_1(t)} \right)^2 \left(\frac{1}{c_t^2} - 1 \right) =: 1 - \mu_t(1 - c_t) \quad (186)$$

The first inequality is due to the fact that $\sum_{i>1} \lambda_i^2(t) (\phi_i^\top(t) \mathbf{w}(t))^2 = 1 - c_t^2$ (Parseval's identity). The last inequality is due to the fact that for $x > -1$, $(1+x)^\alpha \geq 1 + \alpha x$ when $\alpha \geq 1$ or $\alpha < 0$ (Bernoulli's inequality). Therefore the conclusion holds. \square

Lemma 8 (Bound of weight difference). *If $c_t > 0$ and $\lambda_i(t) > 0$ for all i , then $\|\delta \mathbf{w}(t)\|_2 \leq \sqrt{2(1 + \mu_t c_t)(1 - c_t)}$*

Proof. First, for $\mathbf{w}^\top(t+1)\mathbf{w}(t)$, we have (notice that $\lambda_i(t) \geq 0$):

$$\mathbf{w}^\top(t+1)\mathbf{w}(t) = \frac{\sum_i \lambda_i(t) (\phi_i^\top(t) \mathbf{w}(t))^2}{\sqrt{\sum_i \lambda_i^2(t) (\phi_i^\top(t) \mathbf{w}(t))^2}} \quad (187)$$

$$\geq \frac{\lambda_1(t) c_t^2}{\sqrt{\lambda_1^2(t) c_t^2 + \lambda_2^2(t) (1 - c_t^2)}} \geq [1 - \mu_t (1 - c_t)] c_t \quad (188)$$

Therefore,

$$\|\mathbf{w}(t+1) - \mathbf{w}(t)\|_2 = \sqrt{2} \sqrt{1 - \mathbf{w}^\top(t) \mathbf{w}(t+1)} \leq \sqrt{2(1 + \mu_t c_t)(1 - c_t)} \quad (189)$$

\square

Lemma 9. *Let $\delta A = A' - A$, then the maximal eigenvector $\phi_1 := \phi_1(A)$ and $\phi'_1 := \phi_1(A')$ has the following Taylor expansion:*

$$\phi'_1 = \phi_1 + \Delta \phi_1 + \mathcal{O}(\|\delta A\|_2^2) \quad (190)$$

where λ_i is the i -th eigenvalue of A , $\Delta \phi_1 := \sum_{j>1} \frac{\phi_j^\top \delta A \phi_1}{\lambda_1 - \lambda_j} \phi_j$ is the first-order term of eigenvector perturbation. In terms of inequality, there exist $\kappa > 0$ so that:

$$\|\phi'_1 - (\phi_1 + \Delta \phi_1)\|_2 \leq \kappa \|\delta A\|_2^2 \quad (191)$$

Proof. See time-independent perturbation theory in Quantum Mechanics (Fernández, 2000). \square

Lemma 10. *Let L be the minimal Lipschitz constant of A so that $\|A(\mathbf{w}') - A(\mathbf{w})\|_2 \leq L\|\mathbf{w} - \mathbf{w}'\|_2$ holds. If $c_t > 0$ and $\lambda_i(t) > 0$ for all i , then we have:*

$$|d_t - c_{t+1}| = \left| (\phi_1(t) - \phi_1(t+1))^\top \mathbf{w}(t+1) \right| \leq \nu_t (1 - c_t) \quad (192)$$

where

$$\nu(\mathbf{w}) := 2\kappa L^2 (1 + \mu(\mathbf{w})c(\mathbf{w})) + 2L\lambda_{\text{gap}}^{-1}(A\mathbf{w}(t)) \sqrt{\mu(\mathbf{w})(1 + \mu(\mathbf{w})c(\mathbf{w}))} \geq 0 \quad (193)$$

and $\nu_t := \nu(\mathbf{w}(t))$.

Proof. Using Lemma 9 and the fact that $\|\mathbf{w}(t+1)\|_2 = 1$, we have:

$$|d_t - c_{t+1}| = \left| (\phi_1(t) - \phi_1(t+1))^\top \mathbf{w}(t+1) \right| \leq |\Delta \phi_1^\top(t) \mathbf{w}(t+1)| + \kappa L^2 \|\delta \mathbf{w}(t)\|_2^2 \quad (194)$$

where

$$\Delta \phi_1(t) := \sum_{j>1} \frac{\phi_j^\top(t) \delta A(t) \phi_1(t)}{\lambda_1(t) - \lambda_j(t)} \phi_j(t) \quad (195)$$

and $\delta A(t) := A(t+1) - A(t)$. For brevity, we omit all temporal notation if the quantity is evaluated at iteration t . E.g., $\delta \mathbf{w}$ means $\delta \mathbf{w}(t)$ and ϕ_1 means $\phi_1(t)$.

Now we bound $|\Delta \phi_1^\top \mathbf{w}(t+1)|$. Using Cauchy-Schwarz inequality:

$$|\Delta \phi_1^\top \mathbf{w}(t+1)| = \left| \sum_{j>1} \left(\frac{\phi_j^\top \delta A \phi_1}{\lambda_1 - \lambda_j} \right) (\phi_j^\top \mathbf{w}(t+1)) \right| \quad (196)$$

$$\leq \sqrt{\sum_{j>1} \left(\frac{\phi_j^\top \delta A \phi_1}{\lambda_1 - \lambda_j} \right)^2} \sqrt{\sum_{j>1} (\phi_j^\top \mathbf{w}(t+1))^2} \quad (197)$$

$$\leq \frac{1}{\lambda_{\text{gap}}(A)} \sqrt{\sum_{j>1} (\phi_j^\top \delta A \phi_1)^2} \sqrt{\sum_{j>1} (\phi_j^\top \mathbf{w}(t+1))^2} \quad (198)$$

Since $\{\phi_j\}$ is a set of orthonormal bases, Parseval's identity tells that for any vector \mathbf{v} , its energy under any orthonormal bases are preserved: $\sum_j (\phi_j^\top \mathbf{v})^2 = \|\mathbf{v}\|_2^2$. Therefore, we have:

$$|\Delta \phi_1^\top \mathbf{w}(t+1)| \leq \frac{1}{\lambda_{\text{gap}}(A)} \|\delta A \phi_1\|_2 \sqrt{1-d_t^2} \quad (199)$$

$$\leq \frac{L}{\lambda_{\text{gap}}(A)} \|\delta \mathbf{w}(t)\|_2 \sqrt{1-d_t^2} \quad (200)$$

Note that using $-1 \leq d_t \leq 1$ and Lemma 7, we have:

$$\sqrt{1-d_t^2} = \sqrt{1+d_t} \sqrt{1-d_t} \leq \sqrt{2(1-d_t)} \leq \sqrt{2\mu_t(1-c_t)} \quad (201)$$

Finally using bound of weight difference (Lemma 8), we have:

$$|d_t - c_{t+1}| \leq 2\kappa L^2(1+\mu_t c_t)(1-c_t) + L\lambda_{\text{gap}}^{-1} \sqrt{2(1+\mu_t c_t)(1-c_t)} \sqrt{1-d_t^2} \quad (202)$$

$$\leq \nu_t(1-c_t) \quad (203)$$

Here $\nu_t := 2\kappa L^2(1+\mu_t c_t) + 2L\lambda_{\text{gap}}^{-1}(A(t))\sqrt{\mu_t(1+\mu_t c_t)}$. \square

Lemma 11. Let $c_0 := c(\mathbf{w}(0)) = \mathbf{w}^\top(0)\phi_1(A(\mathbf{w}(0))) > 0$. Define local region B_γ :

$$B_\gamma := \left\{ \mathbf{w} : \|\mathbf{w} - \mathbf{w}(0)\|_2 \leq \frac{\sqrt{2(1+\gamma)(1-c_0)}}{1-\sqrt{\gamma}} \right\} \quad (204)$$

Define $\omega(\mathbf{w}) := \mu(\mathbf{w}) + \nu(\mathbf{w})$ to be the irregularity (also defined in Def. 4). If there exists $\gamma < 1$ so that

$$\sup_{\mathbf{w} \in B_\gamma} \omega(\mathbf{w}) \leq \gamma, \quad (205)$$

then

- The sequence $\{c_t\}$ increases monotonously and converges to 1;
- There exists \mathbf{w}_* so that $\lim_{t \rightarrow +\infty} \mathbf{w}(t) = \mathbf{w}_*$.
- \mathbf{w}_* is the maximal eigenvector of $A(\mathbf{w}_*)$ and thus a fixed point of gradient update (Eqn. 7);
- For any t , $\|\mathbf{w}(t) - \mathbf{w}(0)\|_2 \leq \frac{\sqrt{2(1+\gamma)(1-c_0)}}{1-\sqrt{\gamma}}$.
- $\|\mathbf{w}_* - \mathbf{w}(0)\|_2 \leq \frac{\sqrt{2(1+\gamma)(1-c_0)}}{1-\sqrt{\gamma}}$. That is, \mathbf{w}_* is in the vicinity of the initial weight $\mathbf{w}(0)$.

Proof. We first prove by induction that the following induction arguments are true for any t :

- $c_{t+1} \geq c_t > 0$;
- $1 - c_t \leq \gamma^t(1 - c_0)$;
- $\mathbf{w}(t)$ is not far away from its initial value $\mathbf{w}(0)$:

$$\|\mathbf{w}(t) - \mathbf{w}(0)\|_2 \leq \sqrt{2(1+\gamma)(1-c_0)} \sum_{t'=0}^{t-1} \gamma^{t'/2} \quad (206)$$

which suggests that $\mathbf{w}(t) \in B_\gamma$.

Base case ($t = 1$). Since $1 \geq c_0 > 0$, $\mu(\mathbf{w}) \geq 0$, and $A(\mathbf{w})$ is PD, applying Lemma 8 to $\|\mathbf{w}(1) - \mathbf{w}(0)\|_2$, it is clear that

$$\|\mathbf{w}(1) - \mathbf{w}(0)\|_2 = \|\delta \mathbf{w}(0)\|_2 \leq \sqrt{2} \sqrt{(1+\mu_0 c_0)(1-c_0)} \leq \sqrt{2(1+\gamma)(1-c_0)} \quad (207)$$

Note that the last inequality is due to $\mu_0 \leq \gamma$. Note that

$$1 - c_1 = 1 - d_0 + d_0 - c_1 \leq 1 - d_t + |d_0 - c_1| \leq (\mu_0 + \nu_0)(1 - c_0) \leq \gamma(1 - c_0) \quad (208)$$

and finally we have $c_1 \geq 1 - \gamma(1 - c_0) \geq c_0 > 0$. So the base case is satisfied.

Inductive step. Assume for t , the induction argument is true and thus $\mathbf{w}(t) \in B_\gamma$. Therefore, by the condition, we know $\mu_t + \nu_t \leq \gamma$.

By Lemma 8, we know that

$$\|\mathbf{w}(t+1) - \mathbf{w}(t)\|_2 = \|\delta\mathbf{w}(t)\|_2 \leq \sqrt{2(1 + \mu_t c_t)(1 - c_t)} \leq \sqrt{2(1 + \gamma)(1 - c_0)} \gamma^{t/2} \quad (209)$$

Therefore, we know that $\mathbf{w}(t+1)$ also satisfies Eqn. 206:

$$\|\mathbf{w}(t+1) - \mathbf{w}(0)\|_2 \leq \|\mathbf{w}(t) - \mathbf{w}(0)\|_2 + \|\delta\mathbf{w}(t)\|_2 \quad (210)$$

$$\leq \sqrt{2(1 + \gamma)(1 - c_0)} \left[\sum_{t'=0}^{t-1} \gamma^{t'/2} + \gamma^{t/2} \right] \quad (211)$$

$$= \sqrt{2(1 + \gamma)(1 - c_0)} \sum_{t'=0}^t \gamma^{t'/2} \quad (212)$$

Also we have:

$$1 - c_{t+1} = 1 - d_t + d_t - c_{t+1} \leq 1 - d_t + |d_t - c_{t+1}| \quad (213)$$

$$\leq (\mu_t + \nu_t)(1 - c_t) \leq \gamma(1 - c_t) \quad (214)$$

$$\leq \gamma^{t+1}(1 - c_0) \quad (215)$$

and thus we have $c_{t+1} \geq 1 - \gamma(1 - c_t) \geq c_t > 0$.

Therefore, we have

$$1 - c_t \leq \gamma^t(1 - c_0) \rightarrow 0 \quad (216)$$

thus c_t is monotonously increasing to 1. This means that:

$$\lim_{t \rightarrow +\infty} c_t = \lim_{t \rightarrow +\infty} \phi_1^\top(t) \mathbf{w}(t) \rightarrow 1 \quad (217)$$

Therefore, we can show that $\mathbf{w}(t)$ is also convergent, by checking how fast $\|\delta\mathbf{w}(t)\|_2$ decays:

$$\|\delta\mathbf{w}(t)\|_2 \leq \sqrt{2(1 + \mu_t c_t)(1 - c_t)} \leq \sqrt{2(1 + \gamma)(1 - c_0)} \gamma^{t/2} \quad (218)$$

By Cauchy's convergence test, $\mathbf{w}(t) = \mathbf{w}(0) + \sum_{t'=0}^{t-1} \delta\mathbf{w}(t')$ also converges. Let

$$\lim_{t \rightarrow +\infty} \mathbf{w}(t) = \mathbf{w}_* \quad (219)$$

This means that $A(\mathbf{w}_*)\mathbf{w}_* = \lambda_* \mathbf{w}_*$ and thus $P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*)\mathbf{w}_* = 0$, i.e., \mathbf{w}_* is a fixed point of gradient update (Eqn. 7). Finally, we have:

$$\|\mathbf{w}(t) - \mathbf{w}(0)\|_2 \leq \sqrt{2(1 + \gamma)(1 - c_0)} \sum_{t'=0}^{t-1} \gamma^{t'/2} \leq \frac{\sqrt{2(1 + \gamma)(1 - c_0)}}{1 - \sqrt{\gamma}} \quad (220)$$

Since $\|\cdot\|_2$ is continuous, we have the conclusion. \square