
On the Role of Nonlinearity in Training Dynamics of Contrastive Learning on One-layer Network

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 1-layer nonlinear networks with homogeneous activation $h(x) = h'(x)x$, by extending recent α -CL framework [29] and linking it to kernels [26]. We find that 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.

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., [8; 17; 6; 4; 16; 15; 7]) without labels.

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 [19; 2; 22; 21; 32; 33]. 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 [29] and linking it to kernels [26] in large batchsize limit, 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, CL with linear network is PCA under certain conditions [29], and only the most salient pattern (i.e., the maximal eigenvector of the data covariance matrix) is learned, while other less salient ones are discarded, regardless of the number of hidden nodes.

Related works. Previous works [34; 23; 30; 28; 1] that analyze training dynamics mostly focus on supervised learning. Different from [27; 20] 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 [23] that assumes symmetric weight structure and data generated sparse linear models.

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 θ 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 [25] 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 [29] that proposes a general CL framework that covers a broad family of existing CL losses. α -CL maximizes an energy function $\mathcal{E}_\alpha(\theta)$ using gradient ascent:

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

where η is the learning rate, $\text{sg}(\cdot)$ is the stop gradient operator, the *energy* function $\mathcal{E}_\alpha(\theta) := \text{tr} \mathbb{C}_\alpha[\mathbf{f}, \mathbf{f}]$ and $\mathbb{C}_\alpha[\cdot, \cdot]$ is the *contrastive covariance* [29; 21]¹:

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

One important quantity is the *pairwise importance* $\alpha(\theta) = [\alpha_{ij}(\theta)]_{i,j=1}^N$, 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 [29] by choosing different $\alpha(\theta)$, e.g., quadratic loss corresponds to $\alpha_{ij} := \text{const}$ and InfoNCE (with $\epsilon = 0$) corresponds to $\alpha_{ij} := \exp(-d_{ij}^2/\tau) / \sum_{j \neq i} \exp(-d_{ij}^2/\tau)$.

In this paper, we mainly focus on how the nonlinearity affects representation learning when α is *fixed*, i.e., the behavior of the gradient ascent update (Eqn. 2) of the objective $\max_{\theta} \text{tr} \mathbb{C}_{\text{sg}(\alpha)}[\mathbf{f}(\mathbf{x}; \theta)]$, where \mathbf{f} is a nonlinear neural network with parameters θ . Note that in Eqn. 2, the subscript $\text{sg}(\alpha)$ means that while the value of α may depend on current network parameters θ , when running optimization we treat α as an independent variable (no backpropagate gradient through α).

For brevity $\mathbb{C}_\alpha[\mathbf{x}] := \mathbb{C}_\alpha[\mathbf{x}, \mathbf{x}]$. 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* [12; 26]:

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 $\mathcal{K}(\mathbf{a}, \mathbf{b}) = \phi^\top(\mathbf{a})\phi(\mathbf{b}) = \sum_{l=0}^{+\infty} \phi_l(\mathbf{a})\phi_l(\mathbf{b})$ with non-negative high-dimensional mapping $\phi(\cdot) = [\phi_l(\cdot)] \geq 0$.*

Definition 2 (Adjusted PDF $\tilde{p}_l(\mathbf{x})$). *For l -th component ϕ_l of the mapping ϕ , 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 $\phi \equiv 1$. 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 $\mathbb{C}_\alpha[\mathbf{g}] \rightarrow \sum_l z_l^2 \mathbb{V}_{\tilde{p}_l}[\mathbf{g}]$.*

¹Compared to [29], our \mathbb{C}_α definition has an additional constant term $1/2N^2$ to simply the notation.

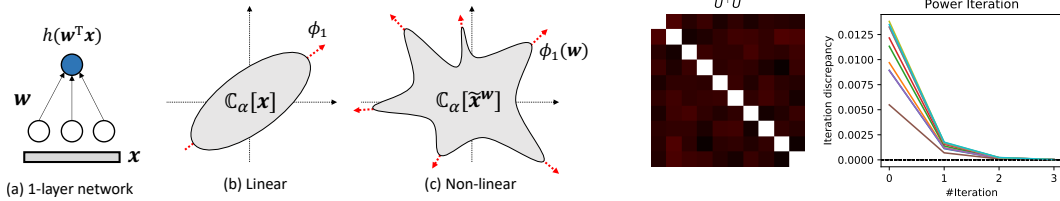


Figure 1: **Left:** Summary of Sec. 3. (a) We analyze the dynamics of one-layer network $h(\mathbf{w}^\top \mathbf{x})$ under CL loss (Eqn. 2). (b) With linear activation $h(x) = x$, then there is only one fixed point (PCA direction). (c) Non-linear activation $h(x)$ creates many critical points and 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. 94) 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 .

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 $W = [\mathbf{w}_1, \dots, \mathbf{w}_K]^\top \in \mathbb{R}^{K \times d}$ 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, $\text{tr} \mathbb{C}_\alpha[\mathbf{f}] = \sum_{k=1}^K \mathbb{C}_\alpha[f_k]$. For convenience, we consider per-filter normalization $\|\mathbf{w}_k\|_2 = 1$, which can be achieved by imposing BatchNorm [18] at each node k [29]. In this case, optimization with each filter \mathbf{w}_k can be considered separately:

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

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. All other patterns will be neglected due to the (overly-smooth) landscape of the objective function. This is the winner-take-all effect and will miss many patterns in the data.

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 [9]/Reversibility [31]). *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 fixed 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.3 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: $P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*) \mathbf{w}_* = \lambda P_{\mathbf{w}_*}^\perp \mathbf{w}_* = 0$. Appendix B.1 gives two concrete examples that show Eqn. 7 has multiple critical points with ReLU activations.

3.1 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 [29].

Possible relationship 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) [11; 24; 30; 35], recently also verified in CL [5], 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 probability of finding a specific local minimum, with randomly initialized weights. For this, we need to study the volume of *attractive basin* defined as $\text{Basin}(\mathbf{w}_*) := \{\mathbf{w} : \mathbf{w}(0) = \mathbf{w}, \lim_{t \rightarrow +\infty} \mathbf{w}(t) = \mathbf{w}_*\}$ for each local optimum \mathbf{w}_* of Eqn. 7. While Sec. B.5 in the Appendix gives hints, we leave this very challenging problem for future work.

References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- [2] 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*, 2019. URL <https://openreview.net/forum?id=SkMQg3C5K7>.
- [3] Mary L Boas and Philip Peters. *Mathematical methods in the physical sciences*, 1984.

- [4] 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.
- [5] 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.
- [6] 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.
- [7] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2020.
- [8] 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.
- [9] 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.
- [10] Francisco M Fernández. *Introduction to perturbation theory in quantum mechanics*. CRC press, 2000.
- [11] 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>.
- [12] 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.
- [13] Mike B Giles. Collected matrix derivative results for forward and reverse mode algorithmic differentiation. In *Advances in Automatic Differentiation*, pp. 35–44. Springer, 2008.
- [14] Gene H Golub and Charles F Van Loan. *Matrix computations*. JHU press, 2013.
- [15] Jean-Bastien Grill, Florian Strub, Florent Altché, 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.
- [16] 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.
- [17] 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.
- [18] 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.
- [19] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [20] 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.
- [21] Li Jing, Pascal Vincent, Yann LeCun, and Yuandong Tian. Understanding dimensional collapse in contrastive self-supervised learning. *ICLR*, 2022.

- [22] Kenji Kawaguchi. Deep learning without poor local minima. *NeurIPS*, 2016.
- [23] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. *Advances in neural information processing systems*, 30, 2017.
- [24] 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.
- [25] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [26] Vern I Paulsen and Mrinal Raghupathi. *An introduction to the theory of reproducing kernel Hilbert spaces*, volume 152. Cambridge university press, 2016.
- [27] 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.
- [28] 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.
- [29] Yuandong Tian. Understanding deep contrastive learning via coordinate-wise optimization. *NeurIPS*, 2022.
- [30] Yuandong Tian, Tina Jiang, Qucheng Gong, and Ari Morcos. Luck matters: Understanding training dynamics of deep relu networks. *arXiv preprint arXiv:1905.13405*, 2019.
- [31] Yuandong Tian, Lantao Yu, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning with dual deep networks. *arXiv preprint arXiv:2010.00578*, 2020.
- [32] 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.
- [33] 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.
- [34] Charles L Wilson, James L Blue, and Omid M Omidvar. Training dynamics and neural network performance. *Neural Networks*, 10(5):907–923, 1997.
- [35] 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:

$$\alpha_{ij} = \exp\left(-\frac{\|\mathbf{y}[i] - \mathbf{y}[j]\|_2^2}{2\tau}\right) \quad (8)$$

$$= \exp\left(-\frac{\|\mathbf{y}[i]\|_2^2}{2\tau}\right) \exp\left(-\frac{\|\mathbf{y}[j]\|_2^2}{2\tau}\right) \exp\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right) \quad (9)$$

And using Taylor expansion, we have

$$\exp\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right) = 1 + \frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau} + \frac{1}{2}\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right)^2 + \dots + \frac{1}{k!}\left(\frac{\mathbf{y}^\top[i]\mathbf{y}[j]}{\tau}\right)^k + \dots \quad (10)$$

Let

$$\tilde{\phi}(\mathbf{y}) := \begin{bmatrix} 1 \\ \tau^{-1/2}\mathbf{y} \\ \frac{1}{\sqrt{2!}}\text{AllChoose}(\tau^{-1/2}\mathbf{y}, 2) \\ \dots \\ \frac{1}{\sqrt{k!}}\text{AllChoose}(\tau^{-1/2}\mathbf{y}, k) \\ \dots \end{bmatrix} \geq 0 \quad (11)$$

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

$$\alpha_{ij} = \phi^\top(\mathbf{x}_0[i])\phi(\mathbf{x}_0[j]) = \sum_{l=0}^{+\infty} \phi_l(\mathbf{x}_0[i])\phi_l(\mathbf{x}_0[j]) \quad (12)$$

which satisfies Def. 1. Here

$$\phi(\mathbf{x}) := \exp\left(-\frac{\|\mathbf{y}\|_2^2}{2\tau}\right) \tilde{\phi}(\mathbf{y}) = \exp\left(-\frac{\|\mathbf{g}(\mathbf{x}) - \mathbf{v}\|_2^2}{2\tau}\right) \tilde{\phi}(\mathbf{g}(\mathbf{x}) - \mathbf{v}) \quad (13)$$

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

$$\mathbb{C}_\alpha^{\text{inter}}[\mathbf{a}, \mathbf{b}] := \frac{1}{2N^2} \sum_{i=1}^N \sum_{j \neq i}^N \alpha_{ij} (\mathbf{a}[i] - \mathbf{a}[j])(\mathbf{b}[i] - \mathbf{b}[j])^\top \quad (14)$$

$$\mathbb{C}_\alpha^{\text{intra}}[\mathbf{a}, \mathbf{b}] := \frac{1}{2N} \sum_{i=1}^N \left(\frac{1}{N} \sum_{j \neq i}^N \alpha_{ij} \right) (\mathbf{a}[i] - \mathbf{a}[i'])(\mathbf{b}[i] - \mathbf{b}[i'])^\top \quad (15)$$

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

$$\mathbb{C}_\alpha[\mathbf{g}] = \mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}] - \mathbb{C}_\alpha^{\text{intra}}[\mathbf{g}]. \quad (16)$$

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

$$\mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}] = \frac{1}{2N^2} \sum_{ij} \mathcal{K}(\mathbf{x}_0[i], \mathbf{x}_0[j]) (\mathbf{g}(\mathbf{x}[i]) - \mathbf{g}(\mathbf{x}[j])) (\mathbf{g}(\mathbf{x}[i]) - \mathbf{g}(\mathbf{x}[j]))^\top \quad (17)$$

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:

$$\int (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y}))^\top \mathbb{P}(\mathbf{x}|\mathbf{x}_0) \mathbb{P}(\mathbf{y}|\mathbf{y}_0) d\mathbf{x} d\mathbf{y} \quad (18)$$

$$= \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] + \mathbb{E}_{\cdot|\mathbf{y}_0}[\mathbf{g}\mathbf{g}^\top] - \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}]\mathbb{E}_{\cdot|\mathbf{y}_0}[\mathbf{g}^\top] - \mathbb{E}_{\cdot|\mathbf{y}_0}[\mathbf{g}]\mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}^\top] \quad (19)$$

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:

$$\begin{aligned} & \int \phi_l(\mathbf{x}_0) \phi_l(\mathbf{y}_0) \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] \mathbb{P}(\mathbf{x}_0) \mathbb{P}(\mathbf{y}_0) d\mathbf{x}_0 d\mathbf{y}_0 \\ &= z_l^2 \int \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] \tilde{p}_l(\mathbf{x}_0) d\mathbf{x}_0 \quad (20) \\ &= z_l^2 \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] \quad (21) \end{aligned}$$

So we have:

$$\mathbb{C}_\alpha^{\text{inter}}[\mathbf{g}] \rightarrow \sum_l z_l^2 (\mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}\mathbf{g}^\top] - \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}] \mathbb{E}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\cdot|\mathbf{x}_0}[\mathbf{g}^\top]) \quad (22)$$

$$= \sum_l z_l^2 \mathbb{V}_{\mathbf{x}_0 \sim \tilde{p}_l, \mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \quad (23)$$

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

$$\frac{1}{N} \sum_{j \neq i} \alpha_{ij} = \frac{1}{N} \sum_{j \neq i} \mathcal{K}(\mathbf{x}_0[i], \mathbf{x}_0[j]) \rightarrow \int \mathcal{K}(\mathbf{x}_0, \mathbf{y}_0) \mathbb{P}(\mathbf{y}_0) d\mathbf{y}_0 \quad (24)$$

$$= \sum_l \phi_l(\mathbf{x}_0) \int \phi_l(\mathbf{y}_0) \mathbb{P}(\mathbf{y}_0) d\mathbf{y}_0 = \sum_l z_l \phi_l(\mathbf{x}_0) \quad (25)$$

Therefore, we have:

$$\mathbb{C}_\alpha^{\text{intra}}[\mathbf{g}] \rightarrow \frac{1}{2} \sum_l z_l \int \phi_l(\mathbf{x}_0) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}')) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}'))^\top \mathbb{P}(\mathbf{x}, \mathbf{x}'|\mathbf{x}_0) \mathbb{P}(\mathbf{x}_0) d\mathbf{x} d\mathbf{x}' d\mathbf{x}_0 \quad (26)$$

Similarly,

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

$$= 2 \int \mathbf{g}(\mathbf{x}) \mathbf{g}^\top(\mathbf{x}) \mathbb{P}(\mathbf{x}|\mathbf{x}_0) d\mathbf{x} - 2 \int \mathbf{g}(\mathbf{x}) \mathbb{P}(\mathbf{x}|\mathbf{x}_0) d\mathbf{x} \int \mathbf{g}^\top(\mathbf{x}') \mathbb{P}(\mathbf{x}'|\mathbf{x}_0) d\mathbf{x}' \quad (28)$$

$$= 2 \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}\mathbf{g}^\top] - 2 \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}^\top] \quad (29)$$

$$= 2 \mathbb{V}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \quad (30)$$

So we have:

$$\mathbb{C}_\alpha^{\text{intra}}[\mathbf{g}] \rightarrow \frac{1}{2} \sum_l z_l \int \phi_l(\mathbf{x}_0) 2 \mathbb{V}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \mathbb{P}(\mathbf{x}_0) d\mathbf{x}_0 \quad (31)$$

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

Using the law of total variation, finally we have:

$$\mathbb{C}_\alpha[\mathbf{g}] \rightarrow \sum_l z_l^2 \mathbb{V}_{\mathbf{x}_0 \sim \tilde{p}_l} \mathbb{E}_{\mathbf{x} \sim p_{\text{aug}}(\cdot|\mathbf{x}_0)}[\mathbf{g}] \quad (33)$$

□

B One-layer model (Sec. 3)

B.1 Computation of the two example models

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 use ReLU activation $h(x) = \max(x, 0)$. 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 $U = [\mathbf{u}_1, \dots, \mathbf{u}_M]$ 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 Sec. B.2 for detailed steps):

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{P}[y = m](1 - \mathbb{P}[y = m])\mathbf{u}_m\mathbf{u}_m^\top \quad (34)$$

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:

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}] = (1 - q_m)^2\mathbf{u}_m\mathbf{u}_m^\top + q_m(I - \mathbf{u}_m\mathbf{u}_m^\top) \quad (35)$$

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 .

B.2 Detailed computation of the two examples

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

$$\mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}_z[\mathbb{E}[\tilde{\mathbf{x}}^{\mathbf{w}}|z]] + \mathbb{E}_z[\mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}|z]] = p_0p_1\Delta(\mathbf{w})\Delta^\top(\mathbf{w}) + p_0\Sigma_0(\mathbf{w}) + p_1\Sigma_1(\mathbf{w}) \quad (36)$$

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:

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{P}[y = m](1 - \mathbb{P}[y = m])\mathbf{u}_m\mathbf{u}_m^\top \quad (37)$$

Latent summation model. If $\mathbf{w} = \mathbf{u}_m$, first notice that due to orthogonal constraints we have $\mathbf{w}^\top \mathbf{x} = \sum_{m'} y_{m'} \mathbf{u}_{m'}^\top \mathbf{w} = y_m$. 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:

$$A(\mathbf{w})|_{\mathbf{w}=\mathbf{u}_m} := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \mathbb{V}[\tilde{\mathbf{x}}] = (1 - q_m)^2\mathbf{u}_m\mathbf{u}_m^\top + q_m(I - \mathbf{u}_m\mathbf{u}_m^\top) \quad (38)$$

B.3 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 fixed 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 w_k evolves independently. Therefore, we omit the subscript k and derive the dynamics of one node 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 [13] to make the derivation easier to understand.

Let $J(\mathbf{w}) := \frac{1}{2}\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x})] = \frac{1}{2}\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h(\mathbf{w}^\top \mathbf{x})]$ 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:

$$\begin{aligned} dJ &= \frac{1}{2}\mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h'(\mathbf{w}^\top \mathbf{x})d\mathbf{w}^\top \mathbf{x}] + \frac{1}{2}\mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x})d\mathbf{w}^\top \mathbf{x}, h(\mathbf{w}^\top \mathbf{x})] \\ &= \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h'(\mathbf{w}^\top \mathbf{x})\mathbf{x}]d\mathbf{w} \end{aligned} \quad (39)$$

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:

$$dJ = \mathbf{w}^\top \mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x})\mathbf{x}, h'(\mathbf{w}^\top \mathbf{x})\mathbf{x}]d\mathbf{w} = \mathbf{w}^\top A(\mathbf{w})d\mathbf{w} \quad (40)$$

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 $\frac{\partial J}{\partial \mathbf{w}} = A(\mathbf{w})\mathbf{w}$. 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 :

$$\begin{aligned} dA &= d\mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x})] \\ &= \mathbb{C}_\alpha[h''(\mathbf{w}^\top \mathbf{x})(d\mathbf{w}^\top \mathbf{x})\mathbf{x}, h'(\mathbf{w}^\top \mathbf{x})\mathbf{x}] + \mathbb{C}_\alpha[h'(\mathbf{w}^\top \mathbf{x})\mathbf{x}, h''(\mathbf{w}^\top \mathbf{x})(d\mathbf{w}^\top \mathbf{x})\mathbf{x}] \end{aligned} \quad (41)$$

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:

$$\begin{aligned} \mathbf{w}^\top (dA)\mathbf{w} &= \mathbb{C}_\alpha[(d\mathbf{w}^\top \mathbf{x})h''(\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x}, h(\mathbf{w}^\top \mathbf{x})] + \mathbb{C}_\alpha[h(\mathbf{w}^\top \mathbf{x}), h''(\mathbf{w}^\top \mathbf{x})(d\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x}] \\ &= 2\mathbb{C}_\alpha[(d\mathbf{w}^\top \mathbf{x})h''(\mathbf{w}^\top \mathbf{x})\mathbf{w}^\top \mathbf{x}, h(\mathbf{w}^\top \mathbf{x})] \end{aligned} \quad (43)$$

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 [3]. 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.4 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 $P_{\mathbf{v}}^\perp A(\mathbf{v})\mathbf{v}$. First, we have:

$$P_{\mathbf{v}}^\perp = I - \mathbf{v}\mathbf{v}^\top = I - \left(\sqrt{1 - \epsilon^2}\mathbf{w}_* + \epsilon\mathbf{u}\right) \left(\sqrt{1 - \epsilon^2}\mathbf{w}_* + \epsilon\mathbf{u}\right)^\top \quad (44)$$

$$= I - \mathbf{w}_*\mathbf{w}_*^\top - \epsilon(\mathbf{u}\mathbf{w}_*^\top + \mathbf{w}_*\mathbf{u}^\top) + \mathcal{O}(\epsilon^2) \quad (45)$$

$$= P_{\mathbf{w}_*}^\perp - \epsilon(\mathbf{u}\mathbf{w}_*^\top + \mathbf{w}_*\mathbf{u}^\top) + \mathcal{O}(\epsilon^2) \quad (46)$$

So we have:

$$P_v^\perp A(\mathbf{w}_*)\mathbf{v} = P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*)\mathbf{v} - \epsilon(\mathbf{u}\mathbf{w}_*^\top + \mathbf{w}_*\mathbf{u}^\top)A(\mathbf{w}_*)\mathbf{v} + \mathcal{O}(\epsilon^2) \quad (47)$$

$$= P_{\mathbf{w}_*}^\perp A(\mathbf{w}_*)\epsilon\mathbf{u} - \epsilon\lambda_*\mathbf{u} + \mathcal{O}(\epsilon^2) \quad (48)$$

$$= P_{\mathbf{w}_*}^\perp (A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} + \mathcal{O}(\epsilon^2) \quad (49)$$

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:

$$P_v^\perp A(\mathbf{v})\mathbf{v} = P_v^\perp A(\mathbf{w}_*)\mathbf{v} + P_v^\perp (A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v} \quad (50)$$

$$= P_{\mathbf{w}_*}^\perp (A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} + P_v^\perp (A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v} + \mathcal{O}(\epsilon^2) \quad (51)$$

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

$$\mathbf{u}^\top P_{\mathbf{w}_*}^\perp (A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} = \mathbf{u}^\top (I - \mathbf{w}_*\mathbf{w}_*^\top)(A(\mathbf{w}_*) - \lambda_*I)\epsilon\mathbf{u} \quad (52)$$

$$= \epsilon\mathbf{u}^\top (A(\mathbf{w}_*) - \lambda_*I)\mathbf{u} \leq -\lambda_{\text{gap}}(\mathbf{w}_*)\epsilon + \mathcal{O}(\epsilon^2) \quad (53)$$

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

$$|\mathbf{u}^\top P_v^\perp (A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v}| \leq \|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v}\|_2 \quad (54)$$

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

$$\|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{w}_*\|_2 \leq \rho(\mathbf{w}_*)\|\mathbf{v} - \mathbf{w}_*\|_2 + \mathcal{O}(\|\mathbf{v} - \mathbf{w}_*\|_2^2) = \rho(\mathbf{w}_*)\epsilon + \mathcal{O}(\epsilon^2) \quad (55)$$

This leads to

$$\|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{v}\|_2 \leq \|(A(\mathbf{v}) - A(\mathbf{w}_*))\mathbf{w}_*\|_2 + \|(A(\mathbf{v}) - A(\mathbf{w}_*))(\mathbf{v} - \mathbf{w}_*)\|_2 \quad (56)$$

$$\leq \rho(\mathbf{w}_*)\epsilon + \mathcal{O}(\epsilon^2) \quad (57)$$

Therefore, we have:

$$\mathbf{u}^\top P_v^\perp A(\mathbf{v})\mathbf{v} \leq -(\lambda_{\text{gap}}(\mathbf{w}_*) - \rho(\mathbf{w}_*))\epsilon + \mathcal{O}(\epsilon^2) \quad (58)$$

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

$$A(\mathbf{w}) := \mathbb{C}_\alpha[\tilde{\mathbf{x}}^{\mathbf{w}}] = \sum_l z_l^2 \mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}]. \quad (59)$$

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

$$\mathbf{e} := \int_{\mathbf{w}^\top \mathbf{x} \geq 0} \mathbf{x}\tilde{p}_l(\mathbf{x})d\mathbf{x}, \quad \mathbf{e}_* := \int_{\mathbf{w}_*^\top \mathbf{x} \geq 0} \mathbf{x}\tilde{p}_l(\mathbf{x})d\mathbf{x} \quad (60)$$

$$\mathbf{E} := \int_{\mathbf{w}^\top \mathbf{x} \geq 0} \mathbf{x}\mathbf{x}^\top \tilde{p}_l(\mathbf{x})d\mathbf{x}, \quad \mathbf{E}_* := \int_{\mathbf{w}_*^\top \mathbf{x} \geq 0} \mathbf{x}\mathbf{x}^\top \tilde{p}_l(\mathbf{x})d\mathbf{x} \quad (61)$$

So we can write

$$\mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}] = E - \mathbf{e}\mathbf{e}^\top, \quad \mathbb{V}_{\tilde{p}_l}[\tilde{\mathbf{x}}^{\mathbf{w}_*}] = E_* - \mathbf{e}_*\mathbf{e}_*^\top \quad (62)$$

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

$$\Omega_+ := \{\mathbf{x} : \mathbf{w}_*^\top \mathbf{x} \geq 0, \mathbf{w}^\top \mathbf{x} \leq 0\} \quad (63)$$

$$\Omega_- := \{\mathbf{x} : \mathbf{w}_*^\top \mathbf{x} \leq 0, \mathbf{w}^\top \mathbf{x} \geq 0\} \quad (64)$$

$$\Omega := \Omega_+ \cup \Omega_- \quad (65)$$

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:

$$E - E_* = \int_{\Omega_-} \mathbf{x}\mathbf{x}^\top \tilde{p}_l(\mathbf{x}) d\mathbf{x} - \int_{\Omega_+} \mathbf{x}\mathbf{x}^\top \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (66)$$

and thus

$$(E - E_*)\mathbf{w}_* = \int_{\Omega_-} \mathbf{x}\mathbf{x}^\top \mathbf{w}_* \tilde{p}_l(\mathbf{x}) d\mathbf{x} - \int_{\Omega_+} \mathbf{x}\mathbf{x}^\top \mathbf{w}_* \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (67)$$

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

$$0 \leq \mathbf{w}_*^\top \mathbf{x} = \mathbf{w}^\top \mathbf{x} + (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \leq (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \leq C_0 \|\mathbf{w}_* - \mathbf{w}\|_2 \quad (68)$$

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

$$\left\| \int_{\Omega_+} \mathbf{x}\mathbf{x}^\top \mathbf{w}_* \tilde{p}_l(\mathbf{x}) d\mathbf{x} \right\|_2 \leq \int_{\Omega_+} |\mathbf{w}_*^\top \mathbf{x}| \|\mathbf{x}\|_2 \tilde{p}_l(\mathbf{x}) d\mathbf{x} \quad (69)$$

$$\leq C_0^2 \|\mathbf{w}_* - \mathbf{w}\|_2 \max_{\mathbf{x} \in \Omega_+} \tilde{p}_l(\mathbf{x}) \int_{\Omega_+, \|\mathbf{x}\|_2 \leq C_0} d\mathbf{x} \quad (70)$$

$$= C_0^3 \|\mathbf{w}_* - \mathbf{w}\|_2 \max_{\mathbf{x} \in \Omega_+} \tilde{p}_l(\mathbf{x}) \frac{\text{vol}(C_0)}{2\pi} \arccos \mathbf{w}^\top \mathbf{w}_* \quad (71)$$

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

$$0 \geq \mathbf{w}_*^\top \mathbf{x} = \mathbf{w}^\top \mathbf{x} + (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \geq (\mathbf{w}_* - \mathbf{w})^\top \mathbf{x} \geq -C_0 \|\mathbf{w}_* - \mathbf{w}\|_2 \quad (72)$$

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

$$\|(E - E_*)\mathbf{w}_*\|_2 \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \|\mathbf{w}_* - \mathbf{w}\|_2 \max_{\mathbf{x} \in \Omega} \tilde{p}_l(\mathbf{x}) \arccos \mathbf{w}^\top \mathbf{w}_* \quad (73)$$

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

$$\arccos \mathbf{w}^\top \mathbf{w}_* = \arcsin \sqrt{1 - (\mathbf{w}^\top \mathbf{w}_*)^2} \leq \frac{\sqrt{1 - (\mathbf{w}^\top \mathbf{w}_*)^2}}{\mathbf{w}^\top \mathbf{w}_*} \quad (74)$$

$$= \frac{\sqrt{1 + \mathbf{w}^\top \mathbf{w}_*} \sqrt{1 - \mathbf{w}^\top \mathbf{w}_*}}{\mathbf{w}^\top \mathbf{w}_*} \leq \frac{\sqrt{2(1 - \mathbf{w}^\top \mathbf{w}_*)}}{\mathbf{w}^\top \mathbf{w}_*} \quad (75)$$

$$= \frac{1}{1 - \epsilon} \|\mathbf{w} - \mathbf{w}_*\|_2 \quad (76)$$

we have:

$$\|(E - E_*)\mathbf{w}_*\|_2 \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \frac{1}{1 - \epsilon} \|\mathbf{w}_* - \mathbf{w}\|_2^2 \max_{\mathbf{x} \in \Omega} \tilde{p}_l(\mathbf{x}) \quad (77)$$

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:

$$\mathbf{e}\mathbf{e}^\top - \mathbf{e}_*\mathbf{e}_*^\top = \mathbf{e}(\mathbf{e} - \mathbf{e}_*)^\top + (\mathbf{e} - \mathbf{e}_*)\mathbf{e}_*^\top \quad (78)$$

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

$$\mathbf{e} - \mathbf{e}_* = \int_{\Omega_-} \mathbf{x} \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} - \int_{\Omega_+} \mathbf{x} \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} \quad (79)$$

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 $\|(\mathbf{e} - \mathbf{e}_*)\mathbf{e}_*^\top \mathbf{w}_*\|_2$:

$$\|(\mathbf{e} - \mathbf{e}_*)\mathbf{e}_*^\top \mathbf{w}_*\|_2 \leq \mathbb{E}_{\tilde{\rho}_l}[h(\mathbf{w}^\top \mathbf{x})] \int_{\Omega} \|\mathbf{x}\|_2 \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} \quad (80)$$

$$\leq C_0^2 \int_{\Omega} \tilde{\rho}_l(\mathbf{x}) d\mathbf{x} \leq C_0^2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \int_{\Omega: \|\mathbf{x}\|_2 \leq C_0} d\mathbf{x} \quad (81)$$

$$\leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \arccos \mathbf{w}^\top \mathbf{w}_* \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \quad (82)$$

$$\leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \frac{1}{1 - \epsilon} \|\mathbf{w} - \mathbf{w}_*\|_2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \quad (83)$$

Overall we have:

$$\|(A(\mathbf{w}) - A(\mathbf{w}_*))\mathbf{w}_*\|_2 \leq \sum_l z_l^2 \|(\nabla_{\tilde{\rho}_l}[\tilde{\mathbf{x}}^{\mathbf{w}}] - \nabla_{\tilde{\rho}_l}[\tilde{\mathbf{x}}^{\mathbf{w}_*}])\mathbf{w}_*\|_2 \quad (84)$$

$$\leq \frac{C_0^3 \text{vol}(C_0)}{\pi} \frac{1}{1 - \epsilon} \left(\sum_l z_l^2 \max_{\mathbf{x} \in \Omega} \tilde{\rho}_l(\mathbf{x}) \right) \|\mathbf{w} - \mathbf{w}_*\|_2 + \mathcal{O}(\|\mathbf{w} - \mathbf{w}_*\|_2^2) \quad (85)$$

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

$$\rho(\mathbf{w}_*) \leq \frac{C_0^3 \text{vol}(C_0)}{\pi} r(\mathbf{w}_*, \alpha) \quad (86)$$

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, $\tilde{\rho}_0(\mathbf{x}; \alpha_u) := \frac{1}{z_0(\alpha_u)} \phi_{u0}(\mathbf{x}) p_D(\mathbf{x}) = p_D(\mathbf{x})$ with $z_0(\alpha_u) = \int \phi_{u0}(\mathbf{x}) p_D(\mathbf{x}) d\mathbf{x} = 1$. This means that

$$r(\mathbf{w}_*, \alpha_u) := \sum_{l=0}^{+\infty} z_l^2(\alpha_u) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{\rho}_l(\mathbf{x}; \alpha_u) \quad (87)$$

$$= z_0^2(\alpha_u) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{\rho}_0(\mathbf{x}; \alpha_u) = \max_{\mathbf{w}_*^\top \mathbf{x}=0} p_D(\mathbf{x}) \quad (88)$$

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

$$\phi_g(\mathbf{x}) = e^{-\frac{h^2(\mathbf{w}_*^\top \mathbf{x})}{2\tau}} \begin{bmatrix} 1 \\ \tau^{-1/2} h(\mathbf{w}_*^\top \mathbf{x}) \\ \frac{1}{\tau^{2/2} \sqrt{2!}} h^2(\mathbf{w}_*^\top \mathbf{x}) \\ \dots \\ \frac{1}{\tau^{k/2} \sqrt{k!}} h^k(\mathbf{w}_*^\top \mathbf{x}) \\ \dots \end{bmatrix} \quad (89)$$

When $l \geq 1$, $z_l^2 \tilde{\rho}_l(\mathbf{x}; \alpha_g) = z_l \phi_{gl}(\mathbf{x}) p_D(\mathbf{x}) = 0$ 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:

$$r(\mathbf{w}_*, \alpha_g) := \sum_{l=0}^{+\infty} z_l^2 \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{p}_l(\mathbf{x}; \alpha_g) = z_0^2(\alpha_g) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \tilde{p}_0(\mathbf{x}; \alpha_g) \quad (90)$$

$$= z_0(\alpha_g) \max_{\mathbf{w}_*^\top \mathbf{x}=0} \phi_{g0}(\mathbf{x}) p_D(\mathbf{x}) \quad (91)$$

$$= z_0(\alpha_g) \max_{\mathbf{w}_*^\top \mathbf{x}=0} p_D(\mathbf{x}) = z_0(\alpha_g) r(\mathbf{w}_*, \alpha_u) \quad (92)$$

Here

$$z_0(\alpha_g) := \int \phi_{g0}(\mathbf{x}) p_D(\mathbf{x}) d\mathbf{x} = \int e^{-\frac{h^2(\mathbf{w}_*^\top \mathbf{x})}{2\tau}} p_D(\mathbf{x}) d\mathbf{x} \leq 1 \quad (93)$$

□

B.5 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) [14], starting from initial value $\mathbf{w} = \mathbf{w}(0)$:

$$\tilde{\mathbf{w}}(t+1) \leftarrow A(\mathbf{w}(t))\mathbf{w}(t), \quad \mathbf{w}(t+1) \leftarrow \frac{\tilde{\mathbf{w}}(t+1)}{\|\tilde{\mathbf{w}}(t+1)\|_2} \quad (94)$$

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*

$$\omega(\mathbf{w}) := \mu(\mathbf{w}) + 2\kappa L^2(1 + \mu(\mathbf{w})c(\mathbf{w})) + 2L\lambda_{\text{gap}}^{-1}(\mathbf{w})\sqrt{\mu(\mathbf{w})(1 + \mu(\mathbf{w})c(\mathbf{w}))}, \quad (95)$$

here κ is the high-order eigenvector bound defined in Appendix (Lemma 6).

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. 94).

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

$$\sup_{\mathbf{w} \in B_\gamma} \omega(\mathbf{w}) \leq \gamma, \quad (96)$$

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. 94) 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. Along the trajectory, let $\phi_i(t) := \phi_1(A(\mathbf{w}(t)))$ 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

$$c_t := c(\mathbf{w}(t)) = \phi_1^\top(t)\mathbf{w}(t), \quad d_t := \phi_1^\top(t)\mathbf{w}(t+1) \quad (97)$$

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

First we assume $A(\mathbf{w})$ is positive definite (PD) over the entire unit sphere $\|\mathbf{w}\|_2 = 1$, then follow Lemma 8, and notice that $\|\mathbf{w} - \mathbf{w}(0)\|_2 = \sqrt{2(1 - \mathbf{w}^\top \mathbf{w}(0))}$, so

$$\|\mathbf{w} - \mathbf{w}(0)\|_2 \leq \frac{\sqrt{2(1 + \gamma)(1 - c_0)}}{1 - \sqrt{\gamma}} \iff \mathbf{w}^\top \mathbf{w}(0) \geq \frac{c_0 - c_\gamma}{1 - c_\gamma} \quad (98)$$

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

$$P_{\mathbf{w}}^\perp \hat{A}(\mathbf{w})\mathbf{w} = P_{\mathbf{w}}^\perp A(\mathbf{w})\mathbf{w} - (\lambda_{\min}(\mathbf{w}) - \epsilon)P_{\mathbf{w}}^\perp \mathbf{w} = P_{\mathbf{w}}^\perp A(\mathbf{w})\mathbf{w} \quad (99)$$

□

Remarks. Note that Lemma 8 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 .

Intuitively, with L and κ small, c_0 close to 1, and λ_{gap} large, Eqn. 96 can always hold with $\gamma < 1$ and the fixed point exists. For example, for the two cases in Sec. B.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.

C Other Lemmas

Lemma 4 (Bound of $1 - d_t$). *Define*

$$\mu(\mathbf{w}) := \frac{1 + c(\mathbf{w})}{2c^2(\mathbf{w})} \left(\frac{\lambda_2(A(\mathbf{w}(t)))}{\lambda_1(A(\mathbf{w}(t)))} \right)^2 = \frac{1 + c(\mathbf{w})}{2c^2(\mathbf{w})} \left[1 - \frac{\lambda_{\text{gap}}(A(t))}{\lambda_1(A(t))} \right]^2 \geq 0 \quad (100)$$

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 :

$$d_t = \frac{\phi_1^\top(t)\tilde{\mathbf{w}}(t+1)}{\|\tilde{\mathbf{w}}(t+1)\|_2} = \frac{\lambda_1(t)\phi_1^\top(t)\mathbf{w}(t)}{\sqrt{\sum_i \lambda_i^2(t) (\phi_i^\top(t)\mathbf{w}(t))^2}} \quad (101)$$

$$\geq \frac{\lambda_1(t)c_t}{\sqrt{\lambda_1^2(t)c_t^2 + \lambda_2^2(t)(1 - c_t^2)}} = \frac{1}{\sqrt{1 + \left(\frac{\lambda_2(t)}{\lambda_1(t)}\right)^2 \left(\frac{1}{c_t^2} - 1\right)}} \quad (102)$$

$$= \left[1 + \left(\frac{\lambda_2(t)}{\lambda_1(t)}\right)^2 \left(\frac{1}{c_t^2} - 1\right) \right]^{-1/2} \quad (103)$$

$$\geq 1 - \frac{1}{2} \left(\frac{\lambda_2(t)}{\lambda_1(t)}\right)^2 \left(\frac{1}{c_t^2} - 1\right) =: 1 - \mu_t(1 - c_t) \quad (104)$$

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. □

Lemma 5 (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 $\lambda_i(t) \geq 0$):

$$\mathbf{w}^\top(t+1)\mathbf{w}(t) = \frac{\sum_i \lambda_i(t) (\phi_i^\top(t)\mathbf{w}(t))^2}{\sqrt{\sum_i \lambda_i^2(t) (\phi_i^\top(t)\mathbf{w}(t))^2}} \quad (105)$$

$$\geq \frac{\lambda_1(t)c_t^2}{\sqrt{\lambda_1^2(t)c_t^2 + \lambda_2^2(t)(1 - c_t^2)}} \geq [1 - \mu_t(1 - c_t)] c_t \quad (106)$$

Therefore,

$$\|\mathbf{w}(t+1) - \mathbf{w}(t)\|_2 = \sqrt{2} \sqrt{1 - \mathbf{w}^\top(t)\mathbf{w}(t+1)} \leq \sqrt{2(1 + \mu_t c_t)(1 - c_t)} \quad (107)$$

□

Lemma 6. 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:

$$\phi'_1 = \phi_1 + \Delta\phi_1 + \mathcal{O}(\|\delta A\|_2^2) \quad (108)$$

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:

$$\|\phi'_1 - (\phi_1 + \Delta\phi_1)\|_2 \leq \kappa \|\delta A\|_2^2 \quad (109)$$

Proof. See time-independent perturbation theory in Quantum Mechanics [10]. \square

Lemma 7. 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:

$$|d_t - c_{t+1}| = \left| (\phi_1(t) - \phi_1(t+1))^\top \mathbf{w}(t+1) \right| \leq \nu_t(1 - c_t) \quad (110)$$

where

$$\nu(\mathbf{w}) := 2\kappa L^2(1 + \mu(\mathbf{w})c(\mathbf{w})) + 2L\lambda_{\text{gap}}^{-1}(A\mathbf{w}(t))\sqrt{\mu(\mathbf{w})(1 + \mu(\mathbf{w})c(\mathbf{w}))} \geq 0 \quad (111)$$

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

Proof. Using Lemma 6 and the fact that $\|\mathbf{w}(t+1)\|_2 = 1$, we have:

$$|d_t - c_{t+1}| = \left| (\phi_1(t) - \phi_1(t+1))^\top \mathbf{w}(t+1) \right| \leq |\Delta\phi_1^\top(t)\mathbf{w}(t+1)| + \kappa L^2 \|\delta\mathbf{w}(t)\|_2^2 \quad (112)$$

where

$$\Delta\phi_1(t) := \sum_{j>1} \frac{\phi_j^\top(t)\delta A(t)\phi_1(t)}{\lambda_1(t) - \lambda_j(t)} \phi_j(t) \quad (113)$$

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:

$$|\Delta\phi_1^\top \mathbf{w}(t+1)| = \left| \sum_{j>1} \left(\frac{\phi_j^\top \delta A \phi_1}{\lambda_1 - \lambda_j} \right) (\phi_j^\top \mathbf{w}(t+1)) \right| \quad (114)$$

$$\leq \sqrt{\sum_{j>1} \left(\frac{\phi_j^\top \delta A \phi_1}{\lambda_1 - \lambda_j} \right)^2} \sqrt{\sum_{j>1} (\phi_j^\top \mathbf{w}(t+1))^2} \quad (115)$$

$$\leq \frac{1}{\lambda_{\text{gap}}(A)} \sqrt{\sum_{j>1} (\phi_j^\top \delta A \phi_1)^2} \sqrt{\sum_{j>1} (\phi_j^\top \mathbf{w}(t+1))^2} \quad (116)$$

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:

$$|\Delta\phi_1^\top \mathbf{w}(t+1)| \leq \frac{1}{\lambda_{\text{gap}}(A)} \|\delta A \phi_1\|_2 \sqrt{1 - d_t^2} \quad (117)$$

$$\leq \frac{L}{\lambda_{\text{gap}}(A)} \|\delta\mathbf{w}(t)\|_2 \sqrt{1 - d_t^2} \quad (118)$$

Note that using $-1 \leq d_t \leq 1$ and Lemma 4, we have:

$$\sqrt{1 - d_t^2} = \sqrt{1 + d_t} \sqrt{1 - d_t} \leq \sqrt{2(1 - d_t)} \leq \sqrt{2\mu_t(1 - c_t)} \quad (119)$$

Finally using bound of weight difference (Lemma 5), we have:

$$|d_t - c_{t+1}| \leq 2\kappa L^2(1 + \mu_t c_t)(1 - c_t) + L\lambda_{\text{gap}}^{-1} \sqrt{2(1 + \mu_t c_t)(1 - c_t)} \sqrt{1 - d_t^2} \quad (120)$$

$$\leq \nu_t(1 - c_t) \quad (121)$$

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

$$B_\gamma := \left\{ \mathbf{w} : \|\mathbf{w} - \mathbf{w}(0)\|_2 \leq \frac{\sqrt{2(1+\gamma)(1-c_0)}}{1-\sqrt{\gamma}} \right\} \quad (122)$$

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

$$\sup_{\mathbf{w} \in B_\gamma} \omega(\mathbf{w}) \leq \gamma, \quad (123)$$

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 $\mathbf{w}(0)$:

$$\|\mathbf{w}(t) - \mathbf{w}(0)\|_2 \leq \sqrt{2(1+\gamma)(1-c_0)} \sum_{t'=0}^{t-1} \gamma^{t'/2} \quad (124)$$

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 5 to $\|\mathbf{w}(1) - \mathbf{w}(0)\|_2$, it is clear that

$$\|\mathbf{w}(1) - \mathbf{w}(0)\|_2 = \|\delta\mathbf{w}(0)\|_2 \leq \sqrt{2}\sqrt{(1+\mu_0c_0)(1-c_0)} \leq \sqrt{2(1+\gamma)(1-c_0)} \quad (125)$$

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

$$1 - c_1 = 1 - d_0 + d_0 - c_1 \leq 1 - d_t + |d_0 - c_1| \leq (\mu_0 + \nu_0)(1 - c_0) \leq \gamma(1 - c_0) \quad (126)$$

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 5, we know that

$$\|\mathbf{w}(t+1) - \mathbf{w}(t)\|_2 = \|\delta\mathbf{w}(t)\|_2 \leq \sqrt{2(1+\mu_t c_t)(1-c_t)} \leq \sqrt{2(1+\gamma)(1-c_0)} \gamma^{t/2} \quad (127)$$

Therefore, we know that $\mathbf{w}(t+1)$ also satisfies Eqn. 124:

$$\|\mathbf{w}(t+1) - \mathbf{w}(0)\|_2 \leq \|\mathbf{w}(t) - \mathbf{w}(0)\|_2 + \|\delta\mathbf{w}(t)\|_2 \quad (128)$$

$$\leq \sqrt{2(1+\gamma)(1-c_0)} \left[\sum_{t'=0}^{t-1} \gamma^{t'/2} + \gamma^{t/2} \right] \quad (129)$$

$$= \sqrt{2(1+\gamma)(1-c_0)} \sum_{t'=0}^t \gamma^{t'/2} \quad (130)$$

Also we have:

$$1 - c_{t+1} = 1 - d_t + d_t - c_{t+1} \leq 1 - d_t + |d_t - c_{t+1}| \quad (131)$$

$$\leq (\mu_t + \nu_t)(1 - c_t) \leq \gamma(1 - c_t) \quad (132)$$

$$\leq \gamma^{t+1}(1 - c_0) \quad (133)$$

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

Therefore, we have

$$1 - c_t \leq \gamma^t(1 - c_0) \rightarrow 0 \quad (134)$$

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

$$\lim_{t \rightarrow +\infty} c_t = \lim_{t \rightarrow +\infty} \phi_1^\top(t) \mathbf{w}(t) \rightarrow 1 \quad (135)$$

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

$$\|\delta \mathbf{w}(t)\|_2 \leq \sqrt{2(1 + \mu_t c_t)(1 - c_t)} \leq \sqrt{2(1 + \gamma)(1 - c_0)} \gamma^{t/2} \quad (136)$$

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

$$\lim_{t \rightarrow +\infty} \mathbf{w}(t) = \mathbf{w}_* \quad (137)$$

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:

$$\|\mathbf{w}(t) - \mathbf{w}(0)\|_2 \leq \sqrt{2(1 + \gamma)(1 - c_0)} \sum_{t'=0}^{t-1} \gamma^{t'/2} \leq \frac{\sqrt{2(1 + \gamma)(1 - c_0)}}{1 - \sqrt{\gamma}} \quad (138)$$

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