EMERGENCE OF ALIGNMENT AND LOCAL ELASTIC ITY IN TWO-LAYER NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

Abstract

Investigating phenomena such as Alignment and Local Elasticity is essential for understanding feature space of Neural Networks and enhancing performance across a wide range of tasks. In this context, we investigate the emergence of these phenomena in two-layer neural networks performing a classification task. This paper reveals Alignment and Local Elasticity emergence condition after one step of training are identical. In particular, we demonstrate that intra-class features are more aligned when the inner product of their mean and the covariance of the training data-label i.e. *train-unseen similarity* is large, with stronger Local Elasticity occurring under this condition. We validate our theory through experiments with a two-layer network showing that both Alignment and Local Elasticity improve as the train-unseen similarity increases. Furthermore, we claim that our analysis provides both theoretical and practical insights into the relationship between train-unseen similarity, alignment, and the improvement of clustering performance on unseen data for neural networks trained on similar domain data. This is supported by experiments, including a multi-layer CNN setup and detailed discussions. Specifically, we show that higher train-unseen similarity improves Recall@1 in two-layer networks and that Alignment and Recall@1 exhibit a positive correlation in metric learning. We also present novel techniques for deriving operator norm bounds of non-centered Sub-Gaussian matrices, extending conventional regression analysis with standard Gaussian assumptions to the binary classification setting.

033

043

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

1 INTRODUCTION

034 Representation learning has been advanced thanks to the introduction of deep learning (Goodfellow et al., 2016; Bengio et al., 2014), surpassing the generalization performance of the conventional 035 machine learning techniques (Bach, 2016; Sánchez & Perronnin, 2011). However, the underlying feature training dynamics that enable deep network to learn more generalizable features An et al. 037 (2023); Radford et al. (2021) remain unclear, prompting studies aimed at theoretically resolving this issue (Damian et al., 2022; Abbe et al., 2021). To understand the learning dynamics, we argue that the following *three challenges* must be addressed: First, under what conditions does learning 040 occur (He & Su, 2019)? Second, to what extent does learning take place under those conditions e.g. 041 Local Elasticity (Dan et al., 2023)? Third, how are the resulting features structured after learning 042 e.g. Alignment (Wang & Isola, 2022; Beaglehole et al., 2024)?

One approach to addressing this challenge is the Neural Tangent Kernel (NTK) (Jacot et al., 2020). 044 NTK studies have explored the alignment structure of features and the concept of Local Elasticity in NTK (Seleznova et al., 2023; Chen et al., 2020; Atanasov et al., 2021; Shan & Bordelon, 2022). 046 However, the NTK operates under a lazy training regime, and its empirical variants exhibit signifi-047 cant discrepancies in modeling neural networks (Chizat et al., 2020; Vyas et al., 2022; Yang & Hu, 048 2022). This makes it challenging to conduct a theoretical analysis of feature learning without additional assumptions, such as whitened data, feature block structure, or label awareness. On the other hand, Conjugate Kernel (CK) approaches have been studied (Pennington & Worah, 2017; Fan & 051 Wang, 2020; Benigni & Péché, 2022), with a key distinction from NTK in their ability to facilitate the analysis of feature learning (Ba et al., 2022; Dandi et al., 2023; Moniri et al., 2024), thereby 052 offering a framework for explaining generalization performance. Building on these properties, we claim that the CK feature learning model not only explains the generalization performance on test data from the same distribution as the training data but also offers a structural analysis of features
 derived from data sampled from unseen distributions that differ from the training distribution.

Deep representations are used in problems where the distributions are "unseen" or "almost similar but different downstream task data" such as in transfer learning (Yosinski et al., 2014; Weiss et al., 2016; Bozinovski, 2020; Galanti et al., 2022), linear probing (Kumar et al., 2022; He et al., 2020; Kornblith et al., 2019), and metric learning (Huang et al., 2024). In these applications, learned features remain effective for data outside the training distributions, even though statistical theories suggest that perfect extrapolation is not attainable (Balestriero et al., 2021; Kang et al., 2024; Armstrong, 1984). Therefore, it is essential to investigate such a problem to advance the deep learning theory. Specifically, this paper investigates the emergence conditions of *Alignment Structure* and *Local Elasticity* for data from unseen distributions to address the *three challenges* mentioned above.

065 066

067

1.1 RELATED WORKS

Conjugate Kernel Many works (Benigni & Péché, 2021; Louart et al., 2017; Hu & Lu, 2022; 068 Goldt et al., 2020) study the CK, which models neural networks and enables the analysis of the 069 structure of the first layer in two-layer networks after the Gradient Descent. Ba et al. (2022) analyze regression tasks in the teacher-student setups to study feature learning in the proportional regime. 071 They demonstrate that neural networks exhibit superior performance compared to linear models, 072 particularly at higher learning rates since the feature learning reflects the structure of the teacher's 073 weights. Moniri et al. (2024) utilize Hermite decomposition to analyze how nonlinear features are 074 learned based on the polynomials. Ba et al. (2023) theoretically compute the condition when neural 075 networks learn the low-dimensional structure of the dataset with spiked covariance Gaussian distri-076 bution data. Bietti et al. (2022) analyze the loss landscape and sample complexity which enables us 077 to learn a single-index model. Ba et al. (2022); Moniri et al. (2024); Ba et al. (2023); Bietti et al. (2022) argue that in teacher-student settings for solving regression problems with centered Gaussian distributions, neural network features can learn the structure of the teacher, thereby improving 079 generalization performance. Unlike these studies, we extend the two-layer network setting to classification with non-centered Sub-Gaussian distributions and examine the phenomena that arise when 081 the network is exposed to input drawn from a distribution different from train distributions. To the best of our knowledge, our work provides the first analysis of non-centered training distributions. 083 We believe this contributes a framework that can be further utilized in analyzing classification.

084 085

Alignment Structure Alignment has been used with various definitions in the study of neural network structure and applications. For instance, there are studies on the following: intra-class fea-087 ture alignment (Deng et al., 2022; Wang & Isola, 2022), feature-weight alignment (Papyan et al., 088 2020), feature-label alignment (Shan & Bordelon, 2022; Atanasov et al., 2021), feature-gradient 089 alignment (Ziyin et al., 2024). We are interested in intra-class feature alignment. Therefore, in the following, "Alignment" will refer to intra-class feature Alignment. As training progresses, 091 the Alignment where features of a given class align towards a single point has been observed (Papyan et al., 2020). It is linked to generalization performance on unseen distributions (Liu et al., 092 2018). For example, some works claim that increasing intra-class alignment of train distributions 093 with inductive bias improves task performance on unseen distributions, particularly in metric learn-094 ing. (Wang et al., 2018; Liu et al., 2017). However, to the best of our knowledge, the conditions 095 under which alignment strongly emerges have not been established. In this work, we demonstrate 096 that the emergence of higher alignment is governed by the relationship between the training data and the input data distribution. Specifically, we show that in a binary classification problem, where 098 β represents the covariance vector of the training data and labels, and μ is the mean of the unseen 099 class conditioned distributions, a larger inner product $|\beta^{\top}\mu|$ i.e. train-unseen similarity leads to 100 higher alignment, thereby providing a theoretical basis for alignment.

101

Neural Collapse (NC) and Unconstrained Layer-Peeled Model (ULPM) research is related to intra-class feature and feature-weight alignment. NC (Papyan et al., 2020) addresses the phenomena that occur with the features and the weights of the classifier head at the final stages of classifier training. At this stage, phenomena related to alignment occur: First, Variability Collapse, i.e. intra-class feature alignment, and Second, self-duality, i.e. feature-weight alignment. Several studies propose the ULPM to analyze NC treating features and weights as unconstrained free variables (Ji et al., 2022; Tirer & Bruna, 2022; Zhu et al., 2021; Fang et al., 2021). However, ULPM, unlike



Figure 1: Emergence of Alignment and Local Elasticity: The Neural Networks feature F(x), $F_0(x)$ from data points surrounding the training data (i.e. unseen data) are influenced by the gradient step, leading to both phenomena. We denote unseen data as α, β, c, d and two classes of training data are represented as a sphere. Notably, distribution α , being closer to the training data, undergoes stronger Alignment and Local Elasticity i.e. the intra-class inner product is enlarged, and the features undergo substantial movement during this single step, compared to other distant distributions. β, c, d .

132

133

134 135

123

124

125

126

127

128

the CK model we use, assumes the features as free variables, which limits its ability analysis about input the data distribution and, consequently, prevents studying the structure of the features. This motivates and provides the need to explore internal features using CK.

136 **Studies on the concept of Local Elasticity (LE)** have been established after observing that data 137 points closer to the training samples are updated more significantly than those farther away (He & Su, 2019). Thus, Local elasticity has been informally described using terms such as "similar-138 ity/closeness". In other words, it is argued that the greater the "similarity" between the training data 139 and the input data, the higher the elasticity of the feature. Subsequently, in He & Su (2019), the 140 elasticity score was formalized as a metric to quantify this informally defined notion of "similarity". 141 Meanwhile, there have been attempts to theoretically understand LE. Zhang et al. (2021) model the 142 learning process of neural networks using SDE to verify its occurrence, but they have a limitation 143 that actual neural networks are not utilized as our CK modeling. Dan et al. (2023) sort training steps 144 into two phases by whether LE occurs or not using Gradient Flow, but they only empirically ob-145 served the basic condition of LE i.e. feature of "similar" sample is updated more, without engaging 146 in theoretical exploration. However, with theoretical assumption and analysis, we establish that this 147 similarity can be measured and expressed as train-unseen similarity.

- Additional related works are discussed in Appendix C.
- 150 151

152

1.2 OVERVIEW

This section provides basic definitions and informal Theorems of the results of Alignment and Local elasticity, which will be detailed in section 4. The phenomenon described here is also illustrated in Figure 1. Let θ be the set of every randomly initialized parameter of a neural network, let d, N denote the dimensions of the data space and feature space, respectively, for $x \in \mathbb{R}^d$ denote $F(x) \in$ \mathbb{R}^N is trained feature and $F_0(x)$ is initialized feature and c is given class conditioned distribution. Feature represents a network output obtained by peeling off the last task layer. The Alignment score and Elasticity score are defined as follows:

160

Definition 1.1 (Alignment score). The **Alignment score** is defined as the expected inner product between the features F(x) for two i.i.d. samples of $c : \mathbb{E}_{x,x' \sim c,\theta}[F(x)^\top F(x')]$.

162 **Definition 1.2** (Elasticity score similar¹ to He & Su (2019)). The **Elasticity score** is defined as the expected L2 distance between F(x), $F_0(x)$ for sample of c: $\mathbb{E}_{x \sim c, \theta}[||F(x) - F_0(x)||^2]$.

These two definitions are informally expressed as the Theorem below, which is an approximation with high probability in the proportional regime for a two-layer neural network after one step training and a Gaussian assumption of given class conditional distribution c.

Theorem 1.3 (Alignment and Elasticity score (Informal of Theorem 4.2, 4.3)). Let **n** be the number of data points. Assume $x, x' \sim \mathcal{N}(\mu, \Sigma)$ be i.i.d random vectors drawn from the arbitrary class conditional distribution given mean $\mu \in \mathbb{R}^d$ and Covariance $\Sigma \in \mathbb{R}^{d \times d}$ and the network allows for Hermite expansion. Let $\beta \triangleq \frac{1}{n\sqrt{N}} X^{\top} y$ from given training datasaet (X, y). Then train (β) unseen (μ) similarity $|\beta^{\top}\mu|$ and $\beta^{\top}\Sigma\beta$ is governing the Alignment and the Elasticity score.

173

Following Theorem, Alignment and Elasticity score approximately increase as Train-unseen similarity $|\beta^{\top}\mu|$ and $\beta^{\top}\Sigma\beta$ grow i.e., with fixed covariance Σ , both scores are approximately polynomial to $|\beta^{\top}\mu|$, which is the similarity between the training sample distributions and the arbitrary class data distributions.

It can be interpreted that the closer the unseen distribution is to the training data (i.e. , higher
Train-unseen similarity), the stronger the effect of Local Elasticity (LE) becomes, and leading to
a stronger Alignment of features. These implications can be observed in section 4, where comparable formulae for the two phenomena are derived, demonstrating their simultaneous occurrence and
correlation.

For theoretical analysis, we define two-layer networks with elementwise activation function that allows Hermite decomposition to decompose a one step trained feature function into initialized features and polynomial functions. This decomposition is explained thoroughly in subsection 3.1 and subsection 3.2. The decomposed feature is analyzed using unseen data distributions assumed to follow Gaussian distributions.

This paper also verifies the following supplementary contributions during our theoretical analysis.
We expand the previous two-layer network analysis method, which is based on regression tasks with standard Gaussian train distribution into binary classification with non-centered Sub-Gaussian distribution. This assumption makes two-layer model available to analyze classification problems in further works or any non-centered Sub-Gaussian training data, which is more discussed in section 3.

Finally, we conduct experiments that empirically verify our analyses using a synthetic dataset where classes of evaluation set are consecutively distant from the training set in section 5.

195 196 197

2 PROBLEM STATEMENTS

Notations Let $\|\cdot\|$ be L^2 or the operator norm. Let \odot be the Hadamad product. Let $A^{\circ k}$ be 199 the Hadamad power. Let C, c > 0 be absolute constants, and let $\kappa \in \mathbb{R}$ be a constant that may 200 change from line to line. Define $[d] \triangleq \{1, 2, \dots, d\}$. Let $\mathbf{1}_{\text{condition}}$ be 1 if the condition is true 201 and 0 otherwise. The operator $diag(\cdot)$ creates a matrix with the elements of the input vector placed 202 along the diagonal. Let $n!! \triangleq \prod_{k=0}^{\lfloor \frac{n}{2} \rfloor - 1} (n - 2k)$ be double factorial. For simplicity, we define (-1)!! = 0!! = 1. For two positive sequences A_n and B_n , we write $A_n = \Theta(B_n)$ if there exist 203 204 constants $c_1, c_2 > 0$ s.t. $c_1 B_n \le A_n \le c_2 B_n$ for sufficiently large n. Similarly, $A_n = \Theta_{\mathbb{P}}(B_n)$ indicates that the relationship holds with high probability as $n \to \infty$. We say $A_n = o(B_n)$ if, for 205 206 every $\epsilon > 0$, there exists $N \in \mathbb{N}$ such that $A_n \leq \epsilon B_n$ for all $n \geq N$. For a vector $x \in \mathbb{R}^n$, the expression x[i] denotes the *i*-th element of x. For a matrix $A \in \mathbb{R}^{n \times m}$, A[i] denotes the *i*-th column 207 208 of A, and A[i : j] denotes the columns from i to j. Additionally, A[:] refers to all elements of A. 209

Hermite Polynomials We employ the probabilist's Hermite polynomials (Moniri et al., 2024; Szegő, 1975; Grad, 1949; Bienstman, 2023). The *n*-th Hermite polynomials, $H_n(\cdot)$, are defined by the recurrence relation: $H_{n+1}(x) = xH_n(x) - nH_{n-1}(x)$, for $n \ge 1$, with the initial conditions $H_0(x) = 1, H_1(x) = x$. Using this recurrence, we have $H_2(x) = x^2 - 1, H_3(x) = x^3 - 3x, \cdots$.

¹Unlike the definition in the original paper, which uses network's predictions, this paper examines Elasticity in the feature level.

216 2.1 PROBLEM SETTINGS

242

243

244

249

250 251

Proportional Regime We consider a two-class classification problem with classes c_1 and c_2 , using two-layer neural networks in the proportional regime. Here, **n**, **d**, and **N** are sample size, data dimension, and feature dimension, respectively. We perform our analysis under the following regime: $n/d \rightarrow \psi_1, N/d \rightarrow \psi_2$ as $\mathbf{n}, \mathbf{d}, \mathbf{N} \rightarrow \infty$, where $\psi_1, \psi_2 \in (0, \infty)$. This setup reflects a scenario where the network width scales proportionally to the data size, aligning with common scaling practices in modern machine learning models.

Training Data Let $\mathcal{D} = (X, Y)$, where $X \in \mathbb{R}^{n \times d}$, $Y \in \{-1, 1\}^{n \times 2} \subseteq \mathbb{R}^{n \times 2}$, represent the training dataset. For any data point (x, y), $y = (1, -1)^{\top}$ if $x \sim c_1$ and $y = (-1, 1)^{\top}$ if $x \sim c_2$, where $x \sim c_i$ indicates that x belongs to class c_i . We denote the *i*-th column of Y, Y[i], as $y_i \in \mathbb{R}^n$. It follows that $y_1 = -y_2$. For every *i*-th row of $X[:][1 : \lfloor n/2 \rfloor]$, we have $X[:][i] \sim c_1$, for every *i*-th row of $X[:][\lfloor n/2 \rfloor + 1 : n]$, we have $X[:][i] \sim c_2$. Let $\tilde{\mathcal{D}} = (\tilde{X}, \tilde{Y})$ an i.i.d. copy of \mathcal{D} .

Evaluation Data In this paper, we employ the "Unseen" dataset as the Evaluation dataset, which is drawn from a distribution different from the one used to generate the training dataset. We assume that "Unseen" samples follow a Gaussian distribution $x \sim \mathcal{N}(\mu, \Sigma)$, where $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathbb{R}^{d \times d}$.

Network Structure We consider two-layer networks to be fully connected. The initial weight of the first layer, $W_0 \in \mathbb{R}^{d \times N}$, is initialized as $W_0[i] \sim Unif(\mathbb{S}^{d-1})$ for $i \in [d]$. We denote W as the one-step trained weight. The initial weights of the second layer, $a_c \in \mathbb{R}^N$ for $c \in \{1, 2\}$, are initialized as $a_c \sim N(0, \frac{1}{N}I)$. For an input x, we define the initialized feature as $F_0(x) \triangleq \sigma(W_0^\top x)$ and the one-step trained feature as $F(x) \triangleq \sigma(W^\top x)$. The network output is defined as the following two-dimensional vector: $(\frac{1}{\sqrt{N}}F(x)^\top a_1, \frac{1}{\sqrt{N}}F(x)^\top a_2)^\top$. The network is designed to output y =(1, -1) for c_1 and y = (-1, 1) for c_2 .

Optimization Problem Denote $\theta = \{W, a_1, a_2\}$ as the set of all network parameters. However, for feature analysis, we only train W and use a_1, a_2 for calculating gradient. To classify the given data, we introduce the Mean Squared Error (MSE) loss

$$L(X, y; \theta) = \frac{1}{2n} \sum_{c \in \{1, 2\}} ||y_c - \frac{1}{\sqrt{N}} \sigma(XW) a_c||^2.$$
(1)

The weight update formula for the first layer is given by $W' = W + \eta \sqrt{\mathbf{N}}G$, where η is the learning rate and G is the negative gradient of $L(X, y; \theta)$ with respect to W expressed as

$$G = -\frac{\partial L}{\partial W} = -\frac{1}{\mathbf{n}} \sum_{c=\{1,2\}} \left[X^{\top} \left[\left(\frac{1}{\sqrt{\mathbf{N}}} (\frac{1}{\sqrt{\mathbf{N}}} \sigma(XW) a_c - y_c) a_c^{\top} \right) \odot \sigma'(XW) \right] \right].$$
(2)

Now, we introduce the assumptions for our theoretical analysis.

Assumption 2.1 (Activation Function). Let $\sigma(x)$ be an element-wise activation s.t. $\sigma, \sigma', \sigma''$ is bounded by λ_{σ} almost surely (a.s.). For $z \sim \mathcal{N}(0, 1)$, it admits a Hermite decomposition i.e. $\sigma(z) = \sum_{k=0}^{\infty} c_k H_k(z)$, where $c_k = \frac{1}{k!} \mathbb{E}_z[\sigma(z)H_k(z)]$. Note that $\mathbb{E}[\sigma(z)] = c_0$ and $\mathbb{E}[z\sigma(z)] = c_1$. We denote $c_{\perp_{0,1}} \triangleq \sqrt{\mathbb{E}[\sigma^2(z) - c_1^2]}$. We assume $c_0 = 0, c_1 \neq 0$ and $c_k^2 k! \leq Ck^{-3/2-w}$, for some constants C, w > 0.

Assumption 2.2 (Learning Rate). $\eta = \Theta(\mathbf{n}^{\alpha}), \ \frac{l-1}{2l} < \alpha < \frac{l}{2l+2}, \ l \in \mathbb{N}.$

Assumption 2.3 (Training Data Structure). Let the class-conditional training data distributions c_1 and c_2 be Sub-Gaussian (Vershynin, 2018; Cole & Lu, 2024; Cao et al., 2021; Jambulapati et al., 2020; Sivakumar et al., 2015; Bombari et al., 2022; Bazinet et al., 2024).

Remark 2.4 (MSE for Classification). Note that utilizing MSE in classification is as well-established
 as using softmax-cross entropy, especially in theoretical analyses of classification problems (Han
 et al., 2022; Zhou et al., 2022).

267 Note 2.5 (Sub-Gaussian Training Data Distribution). The data structure described in Assumption 2.3 allows us to transform the analysis of CK solving linear regression under Gaussian assumptions (Ba et al., 2022; 2023; Moniri et al., 2024) to classification problems. This extension can open new avenues for theoretical analyses of deep representations in classification tasks.

Analysis of Feature in the Proportional Regime with Mse Classification Setting and Sub-Gaussian Data

In this section, we analyze the learning dynamics of a neural network in a single training step, assuming the training data \mathcal{D} originates from two distinct Sub-Gaussian distributions with non-zero means. To achieve this, we decompose the gradient (equation 2) using Hermite decomposition, which allows us to extract the essential rank-one matrix structure. As a result, we approximate the one-step trained feature function $F(x) = \sigma((W_0 + \eta \sqrt{\mathbf{N}}G^{\top}x))$ as F_l by deriving its Hermite expansion, which serves as a key step in deriving our main theorem. The entire process is carried out asymptotically in the proportional regime.

280 281

282

3.1 RANK-ONE APPROXIMATION OF THE FIRST GRADIENT

In this section, we follow the proof structure of Ba et al. (2022) to decompose gradient in our classification learning setting. Unlike their assumption of centered Gaussian training data, we consider non-centered Sub-Gaussian data distributions. In this process, we apply a novel approach involving the concentration of the operator norm on a random matrix. Also, since our framework is not in a teacher-student setting, we use class labels instead of a teacher function.

Starting from equation 2, by performing an orthogonal decomposition of the first Hermite expansion term and the remainder of $\sigma(x)$, we express $\sigma(x) = c_1 x + \sigma_{\perp}(x)$. The gradient G is then decomposed as follows $G_0 = \mathbb{A} + \mathbb{B} + \mathbb{C}$ i.e.

291 292

293

295 296

297

298

299 300

305

306 307

314

287

$$G_{0} = \underbrace{\frac{c_{1}}{\mathbf{n}\sqrt{\mathbf{N}}} X^{\top}(y_{1}a_{1}^{\top} + y_{2}a_{2}^{\top})}_{-\frac{1}{\mathbf{n}\mathbf{N}}} + \underbrace{\frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} X^{\top}(y_{1}a_{1}^{\top} + y_{2}a_{2}^{\top}) \odot \sigma'_{\perp}(XW_{0})}_{-\frac{1}{\mathbf{n}\mathbf{N}}} X^{\top}\sigma(XW_{0})(a_{1}a_{1}^{\top} + a_{2}a_{2}^{\top}) \odot \sigma'(XW_{0}) \dots c.$$

$$(3)$$

We derive the norm bound for the terms \mathbb{A} , \mathbb{B} , and \mathbb{C} in Lemma F.1. Using these bounds, we establish the following Proposition 3.1. For the proof, please refer to Appendix F.

Proposition 3.1. Under the assumptions in subsection 2.1, and when **n** satisfy $1 - \kappa' \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} > 1/2$, the following holds:

$$||G_0 - \mathbb{A}|| \le \kappa \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} ||G_0|| \quad w.p.1 - C(\mathbf{n}e^{-c\log^2 \mathbf{n}} + e^{-c\mathbf{n}}).$$
(4)

Now we utilize \mathbb{A} as the approximate gradient for training the CK model given the training set \mathcal{D} .

3.2 ANALYSIS OF FEATURES AFTER ONE-STEP GD

We now study the feature space induced by the conjugate kernel after one step of gradient descent (GD). We first analyze $\sigma(\tilde{X}W) = F(\tilde{X}W_1) = F(\tilde{X}W_0 + \eta\tilde{X}G)$ by an approximation using Hermite polynomials. Denote $\beta \triangleq \frac{1}{n\sqrt{N}}X^{\top}y_1$ and let $\alpha = a_1 - a_2$. By Proposition 3.1 and results D, E in Lemma G.1, we generalize this to Lemma 3.2. For the proof, see Appendix G.

Lemma 3.2 (Monomial Approximation of Data-Gradient). *For any* $k \in \mathbb{N}$, *sufficiently large* **n**, *and w.p.* 1 - o(1),

$$||(\tilde{X}G^{\top})^{\circ k} - c_1^k (\tilde{X}\beta)^{\circ k} (\alpha^{\circ k})^{\top}|| \le C^k \mathbf{n}^{-\frac{k}{2} \log^{2k} \mathbf{n}}.$$
(5)

Finally, we constructed the Data-Gradient form in our classification setup, satisfying the assumptions same to those in Theorem 3.2 of Moniri et al. (2024). We now decompose F into a feasible form.

Lemma 3.3 (Decomposition of Trained Features). Let $F_0 = \sigma(\tilde{X}W_0^{\top})$. With probability 1 - o(1), $F = F_l + \Delta$,

where $F_l = F_0 + \sum_{k=1}^l c_1^k c_k \eta^k (\tilde{X}\beta)^{\circ k} (\alpha^{\circ k})^\top$ and l is defined in section 2. Moreover, $||\Delta|| = o(\sqrt{\mathbf{n}}), ||F_0|| = \Theta_{\mathbb{P}}(\sqrt{\mathbf{n}}), and ||c_1^k c_k \eta^k (\tilde{X}\beta)^{\circ k} (\alpha^{\circ k})^\top ||$ has an order larger than $o(\sqrt{\mathbf{n}})$.

Based on these results, we analyze the feature representation using the approximation F_l , which dominates the residual term $||\Delta|| = o(\sqrt{n})$ with probability 1 - o(1).

³²⁴ 4 EMERGENCE OF ALIGNMENT AND LOCAL ELASTICITY

In this section we provide theorems indicating $train(\beta)$ -unseen(μ) similarity $|\beta^{\top}\mu|$ and $\beta^{\top}\Sigma\beta$ is governing the Alignment and the Elasticity score. Given separable Sub-Gaussian training data \mathcal{D} , we compute the approximate feature F_l after a single gradient step in the above section.

Condition 4.1 (Condition statement for Theorem 4.2 and Theorem 4.3). Let $x, x' \sim \mathcal{N}(\mu, \Sigma)$ be *i.i.d* random vectors drawn from the arbitrary class conditional distribution given μ, Σ . With assumption in section 2 and following from Lemma 3.3, remark the approximated initialized/trained neural network feature extractor as $F_0(x) = \sigma(W_0^{\top}x)$, $F_l(x) = F_0(x) + \sum_{k=1}^{l} c_1^k c_k \eta^k (x^{\top}\beta)^{\circ k} (\alpha^{\circ k})^{\top}$ where $\beta \triangleq \frac{1}{n\sqrt{N}} X^{\top} y_1$ with training datasaet (X, y).

Theorem 4.2 (Alignment). Following condition 4.1, denote $\mathcal{T}_k \triangleq c_1^k c_k \eta^k \mathbb{E}_x[(\beta^\top x)^k]$. Then, the average Inner Product between two approximated one-step trained features is as follows:

$$\mathbb{E}_{x,x',\theta}[F_l(x)^{\top}F_l(x')] = \mathbb{E}_{x,\theta}||F_0(x)||^2 + 2\langle \mathbb{E}_{x,\theta}F_0(x), \sum_{k=1}^l \mathcal{T}_k \mathbb{E}_{\theta}[\alpha^{\circ k}]\rangle] + \sum_{k=1,j=1}^l \mathcal{T}_k \mathcal{T}_j \mathbb{E}_{\theta} \langle \alpha^{\circ j}, \alpha^{\circ k} \rangle$$
(6)

The first term $\mathbb{E}_{x,\theta} \|F_0(x)\|^2$ only depends on unseen distribution parameter μ, Σ without train distribution. The second term $2\langle \mathbb{E}_{x,\theta}F_0(x), \sum_{k=1}^l \mathcal{T}_k \mathbb{E}_{\theta}[\alpha^{\circ k}] \rangle$] depends on $\mu, \Sigma, |\beta^{\top}\mu|$ and $\beta^{\top}\Sigma\beta$. The last term $\sum_{k=1,j=1}^l \mathcal{T}_k \mathcal{T}_j \mathbb{E}_{\theta} \langle \alpha^{\circ j}, \alpha^{\circ k} \rangle$ depends on $|\beta^{\top}\mu|$ and $\beta^{\top}\Sigma\beta$. Therefore, the alignment measure grows as $|\beta^{\top}\mu|, \beta^{\top}\Sigma\beta$ increases.

Proof. Proof is in Appendix I

Theorem 4.3 (Local Elasticity). Following Condition 4.1, Then, the average L^2 distance between the initialized features $F_0(x)$ and the approximated one step trained features $F_l(x)$ is as follows:

$$\mathbb{E}_{x,\theta} \|F_l(x) - F_0(x)\|^2 = \sum_{k=1}^l \sum_{m=1}^{l} \sum_{i=0}^{k+m} \kappa_{LE} \ |\beta^\top \mu|^{k+m-i} (\beta^\top \Sigma \beta)^{\frac{i}{2}} \ \mathbf{1}_{k+m \text{ and } i \text{ is even}}.$$
 (7)

 κ_{LE} depends only on $k, m, i, \mathbf{N}, c_1, \eta$, and is independent of the data distribution parameters. The local elasticity measure grows as $|\beta^{\top}\mu|, \beta^{\top}\Sigma\beta$ increases.

Proof. Proof is included in Appendix J

Note 4.4 (Interpretation of sign of β). If the two given classes have a zero-centered symmetric structure, a symmetric representation should be learned regardless of the sign of β . This can be observed in our results and observations as well. We defined $\beta = \frac{1}{n\sqrt{N}}X^{\top}y_1$ in subsection 3.2. When the sign of $\alpha = a_1 - a_2$ is flipped, β can also be defined as $\frac{1}{n\sqrt{N}}X^{\top}y_2$. With alternative definition, the same result is obtained for Theorem 4.2 and Theorem 4.3, where the scores are represented as polynomials of $|\beta^{\top}\mu|$ and non-negative $\beta^{\top}\Sigma\beta$.

Note 4.5 (Relationship between l and learning rate η). Our learning rate assumption is that it is determined by the parameter $l \in \mathbb{N}$, which determines the maximum Hermite expansion degree of the Alignment and LE scores as polynomials of $|\beta^{\top}\mu|$. This behavior aligns with the intuition that larger learning rates correspond to more aggressive updates of the features, causing them to shift and align more during the optimization process.

369 370

326

327

328

329

330

331

332

333

334

335

336 337 338

339 340

341

342 343 344

345 346

347

348

354

355

357

5 EXPERIMENTS

371 372 373

Remark 5.1. Recall@1= $\mathbb{E}\mathbf{1}_{y_i=\hat{y}_{i,1-NN}} \hat{y}_{i,1-NN}$ is the class of closest feature to x_i .

In our experiments, we examine the relationships between **train-unseen similarity** (i.e. $|\beta^{\top}\mu|$, **Alignment**), **Elasticity**, and Recall@1. The experimental setups range from synthetic datasets trained with two-layer networks (*Setup 1, 2*) to real-world datasets, including CARS196 (Krause et al., 2013) and CUB200 (Wah et al., 2011), trained with multi-layer networks such as ResNet18 (*Setup 3*) and ResNet50 (He et al., 2015) (*Setup 4*).

396 397

399 400

401 402

415

378 **Setup 1,2** To evaluate the theory, we follow the configurations described in section 2. We use three 379 different non-centered Sub-Gaussian distributions as training datasets: (i) a uniform distribution over 380 a radius- \sqrt{d} ball (Data 1); (ii) a multi-dimensional element-wise truncated Gaussian distribution 381 (Data 2); and (iii) a uniform distribution over a radius- \sqrt{d} sphere (Data 3).² We set d = n = N =382 2^{11} , and $\eta = \mathbf{n}^{0.25}$ in accordance with the assumptions. The means of Data 1 and 3 are v and -v, 383 respectively, where $v \triangleq 5r^2 \cdot \mathbf{u}$, with $\mathbf{u} \sim \text{Unif}(\mathbb{S}^{d-1})$. For Data 2, one class has support on $[1,\infty)$ 384 across all dimensions, while the other class has support on $(-\infty, -1]$. We define $v \triangleq (1, 1, \dots, 1)^{\top}$ 385 for Data 2 used in Evaluation data generation. 386

For the evaluation data, we introduce unseen samples x_{unseen} , which are projected Gaussian distributed and defined as $x_{unseen} \triangleq z - (z^{\top}\nu\nu)/||\nu||^4 + \nu$, where $z \sim \mathcal{N}(0, I)$, and $\nu \triangleq ev$ for **Setup** *I* and $\nu \triangleq Rv$ for **Setup** *2*, with $e \in (-1, 1)$, $R \in SO(d)$. We use this data for measuring Alignment, Elasticity, Recall@1. By adjusting *e* and *R*, one can control the **train-unseen similarity** $\beta^{\top}\mu$, where $\mu \triangleq \mathbb{E}[x_{unseen}]$. Please refer to Figure 2 for illustrations of these setups.



Train Data 1, 2, 3

Evaluation Data for Setup 1, 2

Figure 2: Examples of training datasets (Data 1, 2, 3) and evaluation data used in Setup 1, 2.

Setup 3. 4 We also conduct the experiment with practical settings i.e. the multi-layer networks 403 and the real-world data. In Setup 3, we designate either the CAR or CUB dataset and randomly 404 select two classes as the training set. Then we sample five classes from each evaluation set of CAR 405 and CUB as our new evaluation set. We set $d = N = 2^{11}$, n = 96. The whole model consists of 406 ResNet18 whose output dimension is d, a single nonlinear layer $F(x) = \sigma(W^{\top}x)$, and classifier 407 a_1, a_2 . We measure β and μ from the representations after ResNet18 architecture. Then they are 408 passed through F(x) and final classifier a_1, a_2 . Note that we randomly initialize ResNet18 and do 409 not freeze its layers during training. The Setup 4, conducted on the CARS196 and CUB200 datasets, 410 and its every configuration follows the approach outlined in Zhai & Wu (2019), which represents 411 a baseline in metric learning. We employ the normsoft metric learning loss function Zhai & Wu 412 (2019). This setup is particularly relevant to our focus on unseen distribution, as it conducts the 413 metric learning task with use of unseen data. The detailed configuration of two experiments is in 414 Table 2.

Alignment and Elasticity Observations With *Setup 1,2*, we analyze the behavior of Alignment, Elasticity as $|\beta^{\top}\mu|$ varies with *e* and *R*. Following Thm. 4.2 and 4.3, as $|\beta^{\top}\mu|$ increases, we 416 417 expect to observe a positive tendency in Alignment and Elasticity score defined in Definition 1.1, 418 1.2. In this experiment, the variable e span from -0.9 to 0.9, and the 300 random rotation matrix is 419 generated using a process in subsection K.4 for R. We repeat the experiment 30 times with different 420 initializations of the neural network parameters and include the results along with the mean and 421 standard deviation as in Figure 3, K.1, K.2. It demonstrates that Alignment and Elasticity score 422 occur strongly as e or $\beta^{\perp}\mu$ increase. This phenomenon corresponds to the results of our theoretical 423 findings, which suggest that the features from distributions closer to the training data emerge the 424 stronger Alignment and Local Elasticity.

In *Setup 3*, we validate the theoretical results by adapting the network and data to a practical setting for a binary classification problem. After each training epoch, we evaluate $|\beta^{\top}\mu|$, Alignment, and Elasticity across evaluation datasets CAR and CUB. We calculate the ranks for each of these metrics— $|\beta^{\top}\mu|$, Alignment, and Elasticity—using the values measured across five classes in each dataset. These rankings are then compared across all metrics to see if they maintained consistent ranking orders using Kendall's W (Kendall & Smith, 1939) ranking correlation. A W value of 1

²The Sub-Gaussian property is proven for Data 1 and 3 in Vershynin (2018), and for Data 2 in Lemma E.1.

432 indicates complete agreement in rank order, while a value of 0 indicates no agreement. As a result, 433 we found that, as the theory suggests, there is a rank correlation between Elasticity, Alignment, and 434 $|\beta^{+}\mu|$ on average across four different seeds. Numerically, during the middle stages of training, 435 before the model converges, we observed that the model trained on CAR showed a rank correlation 436 of at least 0.7 across all datasets, while the model trained on CUB exhibited a rank correlation of at least 0.5 across all datasets. See Figure 4. Additionally, on top of the strict order requirement of Kendall's W statistic, we directly observe that $|\beta^{\top}\mu|$, Alignment, and Elasticity simultaneously 438 increased or decreased without aggregation, as shown in Appendix N. 439



Figure 3: Observation of Alignment (a, d) and Elasticity (b, e) Recall@1 (c, f). Figure (a, b) are plotted across different e (lower x-axis, exactly overlapped) and $\beta^{\top}\mu$ (upper x-axis) values. Figure (c, d, e, f) are plotted across different $\beta^{\perp}\mu$ (x-axis) values. For figure (c, f) the blue line represents the clustering performance measured using the features in their initialized state, the orange line reflects the performance after one step of training, and the green line indicates the improvement, i.e., the difference between the two.

1 0.95 0.05 0.05 0.05 0.05		CAR CUB	R@1 v. Align 0.24±0.09 0.29±0.05	p-value (two-sided) 0.00 0.00
0 5 10 15 20	0 5 10 15 20			
(a) Model trained with CAR	(b) Model trained with CUB		Recall@1	Avg. Align
(1)	(-)	CAR	93.65±0.34	843.27±16.83
Figure 4: In Satur 3 the average Kandall's W value (v			68.13 ± 0.40	1072.49 ± 20.46

468 Figure 4: In Setup 3, the average Kendall's W value (y-469 axis) over step (x-axis) for a model trained on the (a)

CAR and (b) CUB dataset with 4 different seeds. The Table 1: In Setup 4, (top) The average cormagenta line represents the Kendall's W value for the relation of Recall@1 and Alignment with p test. (bottom) The final R@1 and Align. CAR dataset, the blue line for the CUB dataset.

472 473

470

471

456

457

458

459

460

437

474 **Connections between** $|\beta^{\top}\mu|$, Alignment and Recall@1 In this section, we analyze the rela-475 tionship between the train-unseen similarity, i.e., $|\beta^{\top}\mu|$ and Recall@1 performance as well as 476 Alignment score. Based on the theoretical finding that neural networks produce features with high alignment and elasticity for unseen classes close to the training data, we hypothesize that data from 477 unseen data distribution similar to train classes undergoes greater shifts during learning, resulting 478 in better alignment and cluster formation with superior Recall@1 performance in the feature space. 479 To validate this hypothesis, we observe whether the higher train-unseen similarity leads to improved 480 Recall@1 performance in Setup 1, 2. 481

482 Through Setup 1, 2, we measure Recall@1 with two classes using cosine similarity. See (c, f) in 483 Figure 3 and K.4. One class is instantiated according to the original definition as *Setup*, while the other is constructed by inverting signs across all axes in data space. After a single learning step, 484 we observe Recall@1 performance increases when the $|\beta^{\top}\mu|$ is higher across all neural networks 485 (orange line in (c, f) at Figure 3).

486 Additionally, in *Setup 1*, we confirm that as $e(\beta^{\top}\mu)$ increases, the recall@1 measure at initialization 487 also increases (blue line in (c) at Figure 3). This is a natural phenomenon because, by the definition 488 of the dataset, as e increases, the L2 distance between the mean of two evaluation classes increases. 489 Further discussion of the observations of Setup 1 is given in subsection K.5. In the setting of Setup 2, the value of $\beta^{\top}\mu$ changes due to variations in rotation; however, the distance between the two 490 classes remains unchanged. Consequently, we observe that the initial Recall@1 does not vary (blue 491 line in (f) at Figure 3). Moreover, we observed that when $\beta^{\top}\mu$ is too small, recall performance 492 decreases after a single step of training (orange line in (f) at Figure 3). This suggests that unseen 493 datasets, which are not too related to the domain of the train dataset, fail to generate meaningful 494 representations. 495

In Setup 4, we extend our experiments from a two-class problem to practical multi-class scenarios 496 within a baseline metric learning setting where the direct computation of β is not feasible. To test 497 the conjecture that strong train-unseen similarity leads to better alignment and improved Recall@1 498 performance, we analyze the correlation between Recall@1 and Alignment scores. At each step of 499 training, we measure class-wise Alignment and class-wise Recall@1 with unseen classes. After, 500 we compute the correlation between the Recall@1 and Alignment for each unseen class. Table 1 501 demonstrates the consistent tendencies of a positive correlation between Recall@1 and Alignment 502 with a near zero p-value. This matches with our empirical results from a two-class synthetic dataset, where we observed a tendency for higher alignment to be associated with better Recall@1 perfor-504 mance. We use Pearson Correlation to measure the strength and direction of the linear relation-505 ship between Recall@1 and Alignment. For the p-value, we use two-sided test. We use Fisher's 506 Combined Probability Test to combine the p-values. We provide unaggregated seed-wise results in 507 Appendix O.

508 509

510

6 CONCLUSION

511 In this paper, we explored the emergence of Alignment and Local Elasticity in two-layer neural net-512 works, focusing on their behavior when trained in the proportional regime. Our theoretical analysis 513 extends the Conjugate Kernel (CK) framework to classification tasks, providing insights into how 514 neural networks learn feature representations, particularly under Sub-Gaussian data distributions. 515 We demonstrate that both Alignment and Local Elasticity arise simultaneously after just one step 516 of training, especially in cases where data distributions closely resemble the training data. This phenomenon not only helps explain the clustering of representations but also sheds light on why 517 neural networks trained on similar domains serve as effective feature extractors for tasks like metric 518 learning. Furthermore, we validated our theoretical findings through experiments across various se-519 tups. These experiments confirmed the theoretical predictions, showing that neural networks exhibit 520 stronger Alignment and Local Elasticity when evaluated on data distributions closer to the training 521 set. Additionally, we identified a possible relationship between Recall@1, one of the generalization 522 performance metrics for unseen distributions, and Alignment. Our work provides a unified frame-523 work for understanding feature learning in neural networks and opens avenues for further research 524 in metric learning, transfer learning, and other task domains where neural networks are applied as 525 feature extractors for unseen distributions. We believe this work offers valuable insights into the 526 dynamics of neural networks, contributing to the broader understanding of deep learning theory.

527

Reproducibility Statement In section 5, Appendix K, Appendix L, and Appendix O the dataset
 generation methods and hyperparameters for experimental reproduction are documented. The code
 used for data generation and the experimental from this research can be downloaded https://
 anonymous.4open.science/r/emk-2E61. Also, We derived all the proofs line by line.

- 531 532
- 533
- 534
- 535
- 536
- 538
- 539

540 REFERENCES

564

565

566

576

580

581

582

- Emmanuel Abbe, Enric Boix-Adsera, Matthew Brennan, Guy Bresler, and Dheeraj Nagaraj. The
 staircase property: How hierarchical structure can guide deep learning, 2021. URL https:
 //arxiv.org/abs/2108.10573.
- Xiang An, Jiankang Deng, Kaicheng Yang, Jaiwei Li, Ziyong Feng, Jia Guo, Jing Yang, and Tongliang Liu. Unicom: Universal and compact representation learning for image retrieval, 2023. URL https://arxiv.org/abs/2304.05884.
- J. Armstrong. Forecasting by extrapolation: Conclusions from 25 years of research. *Interfaces*, 14:
 52–66, 12 1984. doi: 10.1287/inte.14.6.52.
- Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect, 2021.
- Jimmy Ba, Murat A. Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. Highdimensional asymptotics of feature learning: How one gradient step improves the representation, 2022. URL https://arxiv.org/abs/2205.01445.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the presence of low-dimensional structure: A spiked random matrix perspective. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 17420–17449. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/38a1671ab0747b6ffe4dlc6ef117a3a9-Paper-Conference.pdf.
 - Francis Bach. Breaking the curse of dimensionality with convex neural networks, 2016. URL https://arxiv.org/abs/1412.8690.
- 567 Zhidong Bai and Jack W. Silverstein. Spectral Analysis of Large Dimensional Random Matrices.
 568 Springer New York, 2010. ISBN 9781441906618. doi: 10.1007/978-1-4419-0661-8. URL
 569 http://dx.doi.org/10.1007/978-1-4419-0661-8.
- 570
 571
 572
 Randall Balestriero, Jerome Pesenti, and Yann LeCun. Learning in high dimension always amounts to extrapolation, 2021. URL https://arxiv.org/abs/2110.09485.
- 573 Mathieu Bazinet, Valentina Zantedeschi, and Pascal Germain. Sample compression unleashed :
 574 New generalization bounds for real valued losses, 2024. URL https://arxiv.org/abs/ 2409.17932.
- 577 Daniel Beaglehole, Ioannis Mitliagkas, and Atish Agarwala. Feature learning as alignment: a structural property of gradient descent in non-linear neural networks, 2024. URL https://arxiv.org/abs/2402.05271.
 - Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives, 2014. URL https://arxiv.org/abs/1206.5538.
- Lucas Benigni and Sandrine Péché. Eigenvalue distribution of some nonlinear models of random matrices. *Electronic Journal of Probability*, 26(none), January 2021. ISSN 1083-6489. doi: 10.1214/21-ejp699. URL http://dx.doi.org/10.1214/21-EJP699.
- Lucas Benigni and Sandrine Péché. Largest eigenvalues of the conjugate kernel of single-layered neural networks, 2022. URL https://arxiv.org/abs/2201.04753.
- Peter Bienstman. Mathematics for photonics. Course Syllabus, September 2023. URL https://
 studiekiezer.ugent.be/studiefiche/en/E002640/current. Course size: 4.0
 credits, Study time: 120 hours. Offered in English and Dutch.
- 593 Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with shallow neural networks, 2022. URL https://arxiv.org/abs/2210.15651.

594 595 596	Simone Bombari, Mohammad Hossein Amani, and Marco Mondelli. Memorization and optimiza- tion in deep neural networks with minimum over-parameterization. In Alice H. Oh, Alekh Agar- wal, Danielle Belgrave, and Kyunghyun Cho (eds.). Advances in Neural Information Processing
597	Systems, 2022. URL https://openreview.net/forum?id=x8DNliTBSYY.
598	Simone Bombari Shayan Kiyani and Marco Mondelli Beyond the universal law of robustness:
599	Sharper laws for random features and neural tangent kernels. 2023. URL https://arxiv.
600	org/abs/2302.01629.
601	
602	Stevo Bozinovski. Reminder of the first paper on transfer learning in neural networks, 1976. Infor-
603	<i>matica</i> (<i>Slovenia</i>), 44, 2020. URL https://api.semanticscholar.org/CorpusID:
604	227241910.
605 606 607	Yuan Cao, Quanquan Gu, and Mikhail Belkin. Risk bounds for over-parameterized max- imum margin classification on sub-gaussian mixtures. In M. Ranzato, A. Beygelz-
608 609 610	<pre>imer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 8407-8418. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/ file/46e0eae7d5217c79c3ef6b4c212b8c6f-Paper.pdf.</pre>
611 612	Satrajit Chatterjee. Coherent gradients: An approach to understanding generalization in gradient descent-based optimization 2020 URL https://arxiv.org/abs/2002_10657
013	
614	Shuxiao Chen, Hangfeng He, and Weijie J. Su. Label-aware neural tangent kernel: Toward better
610	generalization and local elasticity, 2020. URL https://arxiv.org/abs/2010.11775.
617	Lenaic Chizat Edouard Ovallon and Francis Bach. On lazy training in differentiable programming
618	2020. URL https://arxiv.org/abs/1812.07956.
619	Frank Cole and Yulong Lu. Score-based generative models break the curse of dimensionality in
620 621	learning a family of sub-gaussian distributions. In <i>The Twelfth International Conference on Learn-</i> <i>ing Representations</i> , 2024. URL https://openreview.net/forum?id=wG12xUSqrI.
622 623 624	Alex Damian, Jason D. Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent, 2022. URL https://arxiv.org/abs/2206.15144.
625 626	Soham Dan, Anirbit Mukherjee, Avirup Das, and Phanideep Gampa. Dynamics of local elasticity during training of neural nets, 2023.
628 629 630	Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. How two-layer neural networks learn, one (giant) step at a time, 2023. URL https://arxiv.org/abs/2305.18270.
631 632 633 634	Jiankang Deng, Jia Guo, Jing Yang, Niannan Xue, Irene Kotsia, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. <i>IEEE Transactions on Pattern Analysis</i> and Machine Intelligence, 44(10):5962–5979, October 2022. ISSN 1939-3539. doi: 10.1109/ tpami.2021.3087709. URL http://dx.doi.org/10.1109/TPAMI.2021.3087709.
635 636 637	Zhou Fan and Zhichao Wang. Spectra of the conjugate kernel and neural tangent kernel for linear- width neural networks, 2020. URL https://arxiv.org/abs/2005.11879.
638 639 640 641	Cong Fang, Hangfeng He, Qi Long, and Weijie J. Su. Exploring deep neural networks via layer- peeled model: Minority collapse in imbalanced training. <i>Proceedings of the National Academy</i> <i>of Sciences</i> , 118(43), October 2021. ISSN 1091-6490. doi: 10.1073/pnas.2103091118. URL http://dx.doi.org/10.1073/pnas.2103091118.
642 643 644	Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learn- ing, 2022. URL https://arxiv.org/abs/2112.15121.
645 646 647	Sebastian Goldt, Bruno Loureiro, Galen Reeves, Florent Krzakala, Marc Mézard, and Lenka Zde- borová. The gaussian equivalence of generative models for learning with shallow neural net- works. <i>Proceedings of the 2nd Mathematical and Scientific Machine Learning Conference, PMLR</i> 145:426-471 (2021), 06 2020. URL https://arxiv.org/pdf/2006.14709.pdf.

674

675

676

677

693

694

- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1.
 MIT Press, 2016.
- Harold Grad. Note on n-dimensional hermite polynomials. Communications on Pure and Applied Mathematics, 2:325-330, 1949. URL https://api.semanticscholar.org/ CorpusID:122776469.
- X. Y. Han, Vardan Papyan, and David L. Donoho. Neural collapse under mse loss: Proximity to and dynamics on the central path, 2022. URL https://arxiv.org/abs/2106.02073.
- Hangfeng He and Weijie J Su. The local elasticity of neural networks. *arXiv preprint arXiv:1910.06943*, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. URL https://arxiv.org/abs/1512.03385.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
 unsupervised visual representation learning, 2020. URL https://arxiv.org/abs/1911.
 05722.
- Hong Hu and Yue M. Lu. Universality laws for high-dimensional learning with random features,
 2022. URL https://arxiv.org/abs/2009.07669.
- Hailang Huang, Zhijie Nie, Ziqiao Wang, and Ziyu Shang. Cross-modal and uni-modal soft-label alignment for image-text retrieval. 03 2024. doi: https://doi.org/10.1609/aaai.v38i16.29789. URL https://arxiv.org/pdf/2403.05261.pdf.
- Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding
 generalization in deep learning, 2022. URL https://arxiv.org/abs/2202.08384.
 - L. Isserlis. On a formula for the product-moment coefficient of any order of a normal frequency distribution in any number of variables. *Biometrika*, 12(1/2):134–139, 1918. ISSN 00063444, 14643510. URL http://www.jstor.org/stable/2331932.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks, 2020. URL https://arxiv.org/abs/1806.07572.
- Arun Jambulapati, Jerry Li, and Kevin Tian. Robust sub-gaussian principal component analysis and width-independent schatten packing. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 15689–15701. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_ files/paper/2020/file/b58144d7e90b5a43edcce1ca9e642882-Paper.pdf.
- Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J. Su. An unconstrained layer-peeled
 perspective on neural collapse, 2022. URL https://arxiv.org/abs/2110.02796.
- Katie Kang, Amrith Setlur, Claire Tomlin, and Sergey Levine. Deep neural networks tend to extrapolate predictably, 2024. URL https://arxiv.org/abs/2310.00873.
- Maurice G Kendall and B Babington Smith. The problem of m rankings. *The annals of mathematical statistics*, 10(3):275–287, 1939.
 - Simon Kornblith, Jonathon Shlens, and Quoc V. Le. Do better imagenet models transfer better?, 2019. URL https://arxiv.org/abs/1805.08974.
- Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)* Workshops, June 2013.
- Ananya Kumar, Aditi Raghunathan, Robbie Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution, 2022. URL https://arxiv.org/abs/2202.10054.

702	Donghwan Lee, Behrad Moniri, Xinmeng Huang, Edgar Dobriban, and Hamed Hassani. Demysti-
703	fying disagreement-on-the-line in high dimensions, 2023. URL https://arxiv.org/abs/
704	2301.13371.
705	

- Christopher Liaw, Abbas Mehrabian, Yaniv Plan, and Roman Vershynin. A simple tool for bounding
 the deviation of random matrices on geometric sets, 2016. URL https://arxiv.org/abs/
 1603.00897.
- Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-margin softmax loss for convolutional neural networks, 2017. URL https://arxiv.org/abs/1612.02295.
- Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song. Sphereface: Deep hypersphere embedding for face recognition, 2018. URL https://arxiv.org/abs/1704.08063.
- Cosme Louart, Zhenyu Liao, and Romain Couillet. A random matrix approach to neural networks, 2017. URL https://arxiv.org/abs/1702.05419.
- Bruno Loureiro, Gabriele Sicuro, Cedric Gerbelot, Alessandro Pacco, Florent Krzakala, 718 and Lenka Zdeborová. Learning gaussian mixtures with generalized linear mod-719 Precise asymptotics in high-dimensions. In M. Ranzato, A. Beygelzimer, els: 720 Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural In-721 formation Processing Systems, volume 34, pp. 10144-10157. Curran Associates, Inc., 722 2021 URL https://proceedings.neurips.cc/paper_files/paper/2021/ 723 file/543e83748234f7cbab21aa0ade66565f-Paper.pdf. 724
- Jiawei Ma, Chong You, Sashank J. Reddi, Sadeep Jayasumana, Himanshu Jain, Felix Yu, Shih Fu Chang, and Sanjiv Kumar. Do we need neural collapse? learning diverse features for fine grained and long-tail classification, 2023. URL https://openreview.net/forum?id=
 5gri-cs4RVq.
- Francesca Mignacco, Florent Krzakala, Yue Lu, Pierfrancesco Urbani, and Lenka Zdeborova. The role of regularization in classification of high-dimensional noisy Gaussian mixture. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 6874–6883. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/mignacco20a.html.
- Behrad Moniri, Donghwan Lee, Hamed Hassani, and Edgar Dobriban. A theory of non-linear feature learning with one gradient step in two-layer neural networks, 2024. URL https://openreview.net/forum?id=MY8SBpUece.
- Ryan O'Donnell. Analysis of boolean functions, 2021. URL https://arxiv.org/abs/ 2105.10386.
- Vardan Papyan, X. Y. Han, and David L. Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40): 24652–24663, September 2020. ISSN 1091-6490. doi: 10.1073/pnas.2015509117. URL http://dx.doi.org/10.1073/pnas.2015509117.
- Jeffrey Pennington and Pratik Worah. Nonlinear random matrix theory for deep learning.
 In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and
 R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran
 Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/
 paper/2017/file/0f3d014eead934bbdbacb62a01dc4831-Paper.pdf.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL https://arxiv.org/abs/2103.00020.
- Jorge Sánchez and Florent Perronnin. High-dimensional signature compression for large-scale image classification. pp. 1665–1672, 06 2011. doi: 10.1109/CVPR.2011.5995504.

756 Mariia Seleznova, Dana Weitzner, Raja Giryes, Gitta Kutyniok, and Hung-Hsu Chou. Neural (tan-757 gent kernel) collapse. arXiv preprint arXiv:2305.16427, 2023. 758 Haozhe Shan and Blake Bordelon. A theory of neural tangent kernel alignment and its influence on 759 training, 2022. URL https://arxiv.org/abs/2105.14301. 760 761 Vidvashankar Sivakumar, Arindam Baneriee, and Pradeep K Ravikumar. Bevond sub-762 gaussian measurements: High-dimensional structured estimation with sub-exponential de-763 In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), Adsigns. 764 vances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 765 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/ 766 file/f3f1b7fc5a8779a9e618e1f23a7b7860-Paper.pdf. 767 G. Szegő. Orthogonal Polynomials. American Math. Soc: Colloquium publ. American Mathemati-768 cal Society, 1975. ISBN 9780821810231. URL https://books.google.co.kr/books? 769 id=ZOhmnsXlcY0C. 770 771 Vincent Szolnoky, Viktor Andersson, Balazs Kulcsar, and Rebecka Jörnsten. On the interpretability 772 of regularisation for neural networks through model gradient similarity. Advances in Neural 773 Information Processing Systems, 35:16319–16330, 2022. 774 Tom Tirer and Joan Bruna. Extended unconstrained features model for exploring deep neural col-775 lapse, 2022. URL https://arxiv.org/abs/2202.08087. 776 777 Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. Overparameterization improves robustness 778 to covariate shift in high dimensions. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman 779 Vaughan (eds.), Advances in Neural Information Processing Systems, 2021. URL https:// 780 openreview.net/forum?id=PxMfDdPnTfV. 781 Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *Chapter 5 of:* 782 Compressed Sensing, Theory and Applications. Edited by Y. Eldar and G. Kutyniok. Cambridge 783 University Press, 2012, 11 2010. URL https://arxiv.org/pdf/1011.3027.pdf. 784 785 Roman Vershynin. High-Dimensional Probability: An Introduction with Applications in Data Sci-786 ence. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 787 2018. 788 C. Vignat. A generalized isserlis theorem for location mixtures of gaussian random vectors. 07 789 2011. URL https://arxiv.org/pdf/1107.2309.pdf. 790 791 Nikhil Vyas, Yamini Bansal, and Preetum Nakkiran. Limitations of the ntk for understanding gen-792 eralization in deep learning, 2022. URL https://arxiv.org/abs/2206.10012. 793 794 Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 796 Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and 797 Wei Liu. Cosface: Large margin cosine loss for deep face recognition, 2018. URL https: 798 //arxiv.org/abs/1801.09414. 799 800 Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through align-801 ment and uniformity on the hypersphere, 2022. URL https://arxiv.org/abs/2005. 10242. 802 803 Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. A survey of transfer learning. Journal of 804 Big data, 3:1–40, 2016. 805 Greg Yang and Edward J. Hu. Feature learning in infinite-width neural networks, 2022. URL 807 https://arxiv.org/abs/2011.14522. 808 Yongyi Yang, Jacob Steinhardt, and Wei Hu. Are neurons actually collapsed? on the fine-grained 809

810	Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep
811	neural networks?, 2014. URL https://arxiv.org/abs/1411.1792.
812	

- Jacob A. Zavatone-Veth, Sheng Yang, Julian A. Rubinfien, and Cengiz Pehlevan. Neural networks
 learn to magnify areas near decision boundaries, 2023.
- Andrew Zhai and Hao-Yu Wu. Classification is a strong baseline for deep metric learning, 2019.
 URL https://arxiv.org/abs/1811.12649.
- Jiayao Zhang, Hua Wang, and Weijie Su. Imitating deep learning dynamics via locally elastic
 stochastic differential equations. *Advances in Neural Information Processing Systems*, 34:6392–6403, 2021.
- Jinxin Zhou, Xiao Li, Tianyu Ding, Chong You, Qing Qu, and Zhihui Zhu. On the optimization landscape of neural collapse under mse loss: Global optimality with unconstrained features, 2022. URL https://arxiv.org/abs/2203.01238.
- Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu.
 A geometric analysis of neural collapse with unconstrained features, 2021. URL https: //arxiv.org/abs/2105.02375.
 - Liu Ziyin, Isaac Chuang, Tomer Galanti, and Tomaso Poggio. Formation of representations in neural networks, 2024. URL https://arxiv.org/abs/2410.03006.

A LIMITATIONS AND FUTURE WORKS

866 In this study, we have focused on non-centered sub-Gaussian training data, but this framework could 867 be extended to more complex distributions, such as Gaussian mixtures. Exploring these broader 868 classes of data distributions would enrich our understanding of the model's generalization capabilities. By the way, we have found that both Alignment and Local Elasticity are more strongly emerged by train-unseen similarity. However, it is necessary to explore how these two phenomena occur si-870 multaneously. Furthermore, replacing the MSE loss with softmax cross-entropy could link this work 871 more directly to Neural Collapse research (Ji et al., 2022), providing new insights into the geometric 872 structures emerging during training. Additionally, studying scenarios where the parameters diverge 873 further from their initialization after the first step of training could offer a long-term perspective on 874 the learning dynamics. Moreover, There seems to be a connection between neural network align-875 ment and the contraction of the Riemannian metric (Zavatone-Veth et al., 2023). Further research 876 into this relationship could unveil deeper insights into the geometry of neural networks. Finally, 877 in this study, the average of the Alignment and Elasticity scores was analyzed, and through mul-878 tiple experiments, the validity of the analysis was supported. Theoretically, this can be extended 879 to concentration as in Loureiro et al. (2021) and Mignacco et al. (2020), and analyzing the condi-880 tions under which the Alignment and Elasticity scores concentrate around the mean is one of the important research directions.

882 883

884

895

896 897

899

903 904

905

B EXTEND CLASSIFICATION SETTINGS TO REGRESSION SETTINGS

We chose a binary classification setup to analyze network learning, ensuring that it aligns with the settings proposed by He & Su (2019). However, our analysis is not limited to classification tasks alone. Inspired by works like Ba et al. (2022), we will incorporate a setting that reflects an regression form to demonstrate that our proof techniques can straightforwardly extend to scenarios involving regression setting. This straightforward adaptability is possible because our analysis applies to any loss or model that satisfies the condition of Proposition 3.1 and Lemma 3.2 in the main text. We argue that this is a key aspect showcasing the extensibility of our study.

⁸⁹² Under all the assumptions stated in our paper, we define a new random variable $a \sim \mathcal{N}(0, \frac{1}{N}I)$, ⁸⁹³ and modify the assumptions from Ba et al. (2022) for regression by replacing the centered Gaussian ⁸⁹⁴ assumption with a non-centered sub-Gaussian assumption, leading to the following problem setup:

$$x_i \sim \text{SG s.t. } \mathbb{E}[x_i] \neq 0, \quad y_i = f^*(x_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma_{\epsilon}^2),$$

$$f^* \text{ is a Lipschitz function, and } \sqrt{\mathbb{E}_x[f^{*2}]} = \Theta(1).$$
(8)

In the above setup, we define the loss as follows:

900
$$\theta = \{W, a\},$$

901
902

$$\{a\}, L(X, y; \theta) = \frac{1}{2n} ||y - \sigma(XW)a||^2$$
$$G = -\frac{\partial L}{\partial W} = -\frac{1}{n} \left[X^\top \left[\left(\frac{1}{\sqrt{N}} (\frac{1}{\sqrt{N}} \sigma(XW)a - y)a^\top \right) \odot \sigma'(XW) \right] \right]$$

In this case, if we define $\alpha = a$ as opposed to the main text, α becomes a Gaussian with zero mean and variance halved, so it follows the same bound structure.

Specifically, based on the sub-Gaussian bound results in Appendix E, \mathbb{A} can be bounded using the second equation of Lemma 14(i) from Ba et al. (2022) and the fact that \mathbb{A} is rank-1, so $||A|| = ||A||_F$. For \mathbb{B} , we can remove $||a_2||_{\infty}$ from equation 20 in our Lemma F.1 proof. For \mathbb{C} , by removing $a_2a_2^T$ from equation 30, we obtain the same bounds as the previous results.

In conclusion, these three bounds satisfy our Proposition 3.1 under the same conditions, and the same conclusion holds even outside of classification tasks. Note that β is defined as in the previous setting.

914 915

916

C ADDITIONAL RELATED WORK

917 ADDITIONAL RELATED WORKS ABOUT NEURAL COLLAPSE Additionally, investigations into the features of neural networks have led to observations suggesting that Neural Collapse does not

918 actually take place internally (Yang et al., 2023), and claims that it does not contribute to under-919 standing generalization (Hui et al., 2022; Ma et al., 2023; Galanti et al., 2022). 920

921 REGARDING FEATURE-GRADIENT ALIGNMENT, Zivin et al. (2024) argues that the alignment be-922 tween features, weights, and gradients naturally facilitates the learning of compact representations. Beaglehole et al. (2024) investigate the Alignment of feature matrices by examining the correlation 923 between feature matrices and the outer product of gradients. Furthermore, He & Su (2019) claim 924 that gradients influence feature structures and Szolnoky et al. (2022); Chatterjee (2020) unveil the 925 relation between the gradients of similar datasets. 926

927 RIEMANNIAN GEOMETRY PERSPECTIVE There is research from a Riemannian geometry per-928 spective related to our result that the closer the data is to training data, i.e., the larger $|\beta^{\perp}\mu|$, the 929 greater the occurrence of alignment and LE. Zavatone-Veth et al. (2023) find out the decrease of 930 determinant of Riemannian metrics in the space i.e. volume decrease around training data. This is 931 related to the strong tendency of the Local Elasticity and Alignment at the point close to the training 932 samples.

933

945

946 947

948

949 950

951

952 953

954

961

963 964 965

966 967

934 NEURAL NETWORK THEORIES BASED ON THE TWO-LAYER ASSUMPTION Several prior stud-935 ies have effectively utilized the feature extractor assumption same to our, to interpret phenomena 936 observed in practical neural networks. For example, Damian et al. (2022) analyzed the efficient 937 generalization and transfer performance of neural networks, while Tripuraneni et al. (2021) used this framework as a tool to study robustness to input distribution shifts. Similarly, Lee et al. (2023) 938 employed it to analyze out-of-distribution inputs, and Bombari et al. (2023) utilized it to investi-939 gate adversarial robustness. These studies focused on understanding phenomena of neural network 940 representations, particularly the hidden representations allowing them to model and explain behav-941 iors observed in practical deep learning scenarios. Based on this body of work, we argue that the 942 assumption of a two-layer network capable of learning hidden representations is a reasonable and 943 effective framework for analyzing neural networks without significant loss of generality. 944

D ADDITIONAL NOTATIONS

 $\|\cdot\|_F$ is the Frobenius norm. $\|\cdot\|_{\infty}$ is the infinity norm. $\|\cdot\|_{\psi_2}$ is orlicz-2 norm $e^{(i)}$ Standard basis vector with 1 at position *i*.

ADDITIONAL INFORMATION OF HERMITE POLYNOMIALS Hermite polynomials can be represented as the following explicit form:

$$H_n(x) = (-1)^n e^{\frac{x^2}{2}} \frac{d^n}{dx^n} e^{-\frac{x^2}{2}}.$$

for $n \in \mathbb{N}_0$. Lastly, there are another expression:

$$H_n(x) = n! \sum_{m=0}^{\lfloor \frac{n}{2} \rfloor} \frac{(-1)^m}{m!(n-2m)!} \frac{x^{n-2m}}{2^m}$$

The probabilist's Hermite polynomials form an orthogonal set with respect to the standard normal 960 weight function $w(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$ on the interval $(-\infty,\infty)$. Their orthogonality condition is given 962 by:

$$\int_{-\infty}^{\infty} H_m(x) H_n(x) \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \delta_{mn} n!,$$

where δ_{mn} is the Kronecker delta, and n! is the factorial of n.

GENERALIZATION OF CENTERED SUB-GAUSSIAN RESULTS TOWARD E NON-CENTERED

⁹⁶⁹ 970

For more detailed explanation and well known results of Sub-Gaussian we used, please refer to 971 Vershynin (2018; 2010).

Lemma E.1. Truncated Gaussian distribution which have support on (a, b) s.t. $a, b \in (-\infty, \infty)$ is Sub-Gaussian.

975Proof. Denote $\mathcal{N}_{(a,b)}(0,\sigma^2)$ is Truncated Gaussian distribution which have support on (a,b) s.t.976 $a, b \in (-\infty, \infty)$. support $(\mathcal{N}_{(a,b)}(0,\sigma^2)) \subset \mathbb{R}^d$. Therefore, $\mathbb{P}(|X| \ge t)$ s.t. $X \sim \mathcal{N}_{(a,b)}(0,\sigma^2)$ have
same tail behavior with Gaussian and Gaussian is Sub-Gaussian.978

Lemma E.2. Sum of non-centered Sub-Gaussian random variable is Sub-Gaussian.

Proof. If the Orlicz 2 norm is bounded $||X||_{\psi_2} < \infty$, then X is Sub-Gaussian. Also, $||\mathbb{E}X||_{\psi_2} \le C||X||_{\psi_2}$, and Sum of centered Sub-Gaussian random variable is Sub-Gaussian. We show $||\sum X_i||_{\psi_2} < \infty$, s.t. X is non-centered Sub-Gaussian.

$$\|\sum X_{i}\|_{\psi_{2}} \leq \|\sum (X_{i} - \mathbb{E}X_{i})\|_{\psi_{2}} + \|\sum \mathbb{E}X_{i}\|_{\psi_{2}}$$

$$\leq \|\sum (X_{i} - \mathbb{E}X_{i})\|_{\psi_{2}} + \sum \|\mathbb{E}X_{i}\|_{\psi_{2}}$$

$$\leq \|\sum (X_{i} - \mathbb{E}X_{i})\|_{\psi_{2}} + C \sum \|X_{i}\|_{\psi_{2}} < \infty$$

$$\Box$$

Lemma E.3. (Operator norm bound for non-centered Sub-Gaussian matrix, generalization of 4.4.5 in Vershynin (2018)) let $A \in \mathbb{R}^{m \times n}$, A[i][j] is independent, non-centered Sub-Gaussian. $\forall t > 0$,

$$||A|| \le CK(\sqrt{m} + \sqrt{n} + t) \text{ w.p. } 1 - \exp(-t^2)$$
Alternatively, $||A|| \le CK(\sqrt{m} + n + t) \text{ w.p. } 1 - \exp(-t^2)$
(10)

 $K = \max_{i,j} ||A[i][j]||_{\psi_2}$ 996 $K = \max_{i,j} ||A[i][j]||_{\psi_2}$

A

Lemma E.4. (Expectation of operator norm for non-centered Sub-Gaussian matrix generalization of 4.4.6 in Vershynin (2018))

$$\mathbb{E}||A|| \le CK(\sqrt{m} + \sqrt{n})$$
Iternatively, $\mathbb{E}||A|| \le CK(\sqrt{m+n}), \text{ and, } \mathbb{E}||A||^2 \le C(m+n)$
(11)

Proof of Lemma E.3 and Lemma E.4. Based on the result of Lemma E.2, one can follow the same proof process of Vershynin (2018)

F ADDITIONAL RESULTS OF SECTION 3.1

For the aforementinoed \mathbb{A} , \mathbb{B} , and \mathbb{C} , we obtain bounds for each operator norm as follows **Lemma F.1.**

$$\mathbb{P}\left(\|\mathbb{A}\| \leq C(\frac{1}{\sqrt{\mathbf{N}}} - C\frac{\sqrt{\mathbf{d}}}{\sqrt{\mathbf{nN}}})\right) \leq 2\left(e^{-c\mathbf{N}} + e^{-c\mathbf{n}}\right) \\
\mathbb{P}\left(\|\mathbb{B}\| \geq \frac{C}{\mathbf{n}\sqrt{\mathbf{Nd}}}(\sqrt{\mathbf{n}} + \sqrt{\mathbf{d}})(\sqrt{\mathbf{n}} + \sqrt{\mathbf{N}})\log\mathbf{N}\right) \leq C\left(e^{-c\mathbf{N}} + e^{-c\mathbf{d}} + \mathbf{N}e^{-c\log^{2}\mathbf{n}} + e^{-(\sqrt{\mathbf{n}} + \sqrt{\mathbf{d}})^{2}}\right) \\
\mathbb{P}\left(\|\mathbb{C}\| \geq \frac{C}{\sqrt{\mathbf{nN}}}(2\sqrt{\mathbf{d}} + \sqrt{\mathbf{n}})\log\mathbf{n}\log\mathbf{N}\right) \leq 2\left(\mathbf{n}e^{-c\mathbf{d}} + \mathbf{n}e^{-c\log^{2}\mathbf{n}} + \mathbf{N}e^{-c\log^{2}\mathbf{n}}\right).$$
(12)

Proof of Lemma F.1 (A). Let us first define $\alpha = a_1 - a_2$. Then, we obtain

A

$$= \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} X^{\top} \left(y_1 a_1^{\top} + y_2 a_2^{\top} \right).$$
(13)

)

1022 Then, we can find an explicit notation of the norm as

1023
1024
$$\|\mathbb{A}\| = \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X^{\top}(y_1 a_1^{\top} + y_2 a_2^{\top})\| = \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X^{\top}y_1(a_1^{\top} - a_2^{\top})\|_{op}$$
(14)

1025
$$= \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X^{\top}y_1\|_2 \|(a_1 - a_2)\|_2 = \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} (y_1^{\top}XX^{\top}y_1)^{1/2} \|\alpha\|_2$$
(14)

1027 $\|\alpha\|_2$ study By definition, $\alpha \sim \mathbf{N}(0, \frac{2}{\mathbf{N}})$, so $\frac{\sqrt{\mathbf{N}}}{2}\alpha[i]$ is a sub-Gaussian. Use Thm 3.3.1 in Vershynin (2018),

$$\mathbb{P}\left(\left|\left\|\frac{\sqrt{\mathbf{N}}}{2}\alpha\right\| - \sqrt{\mathbf{N}}\right| \ge t\right) \le 2e^{-ct^2} \quad \text{let } t = \sqrt{\mathbf{N}}/2$$

$$\mathbb{P}\left(\left\|\alpha\right\|_{2} \le 1\right) \le 2e^{-c\mathbf{N}} \tag{15}$$

 $\mathbb{P}(\|\alpha\|_2 \le 1) \le 2e$

 $(y_1^{\top}XX^{\top}y_1)^{1/2}$ study Note that the U, V matrices resulting from the SVD belong to the *O*-group, so there is no length transformation.

$$y_{1}^{\top}XX^{\top}y_{1} = \|X^{\top}y_{1}\|_{2}^{2} = \|U\Sigma V^{\top}y_{1}\|_{2}^{2} = \|\Sigma V^{\top}y_{1}\|$$

$$= \sum_{i} \sigma_{i}^{2}|V^{\top}y[i]|^{2} \ge \sigma_{\min}^{2} \sum_{i} |V^{\top}y[i]|^{2} = \sigma_{\min}^{2} \|y\|_{2}^{2} = \mathbf{n}\sigma_{\min}^{2}$$
(16)

We get $(y_1^{\top}XX^{\top}y_1)^{1/2} \ge \sqrt{\mathbf{n}}\sigma_{\min}$. σ_{\min} is singular value of X which is a anistropic sub-Gaussian matrix. With the result of Remark 1.2 in Liaw et al. (2016),

$$\mathbb{P}\sigma_{\min} \le (\sqrt{\mathbf{n}} - c\sqrt{\mathbf{d}})) \le e^{-\mathbf{n}}.$$
(17)

1044
1045 Therefore,
$$\mathbb{P}(\|\mathbb{A}\| \le C(\frac{1}{\sqrt{N}} - C\frac{\sqrt{d}}{\sqrt{nN}})) \le 2(e^{-cN} + e^{-cn}).$$

1046 $\mathbb{P}(\|\mathbb{A}\| \le C(\frac{1}{\sqrt{N}} - C\frac{\sqrt{d}}{\sqrt{nN}})) \le 2(e^{-cN} + e^{-cn}).$

Fact F.2 (from Ba et al. (2022)). For $m \in \mathbb{R}^m$, $n \in \mathbb{R}^n$, $M \in \mathbb{R}^{m \times n}$,

$$mn^{\top} \odot M = diag(m) M diag(n) \|mn^{\top} \odot M\| \le \|diag(m)\| \|M\| \|diag(n)\| = \|m\|_{\infty} \|M\| \|n\|_{\infty} n$$
(18)

1051 Lemma F.3. For Sub-Gaussian R.V. a,

$$\mathbb{P}(\|a\|_{\infty} \le t/\sqrt{\mathbf{N}}) \ge 1 - 2\mathbf{N}e^{-ct^2}$$

Proof. We use the Hoeffding inequality such that

$$\mathbb{P}(\|a\|_{\infty} \geq \frac{t}{\sqrt{\mathbf{N}}}) = \mathbb{P}\left(\max_{i} |a_{i}| \geq \frac{t}{\sqrt{\mathbf{N}}}\right) \leq \mathbb{P}\left(\bigcup_{i} \{|a_{i}| \geq \frac{t}{\sqrt{\mathbf{N}}}\}\right) \leq \sum_{i} \mathbb{P}\left(|a_{i}| \geq \frac{t}{\sqrt{\mathbf{N}}}\right)$$

$$\stackrel{\text{i.i.d.}}{=} \mathbf{N}\mathbb{P}\left(|a_{i}| \geq \frac{t}{\sqrt{\mathbf{N}}}\right) = \mathbb{P}(|\sqrt{\mathbf{N}}a_{i}| \geq t) \leq 2\mathbf{N}\exp(-ct^{2})$$
(19)

Fact F.4. Let a sub-Gaussian random variable v s.t. $||v||_{\psi_2} \leq k$, and bounded function σ , then $\sigma(v)$ is Sub-Gaussian, i.e. $||\sigma(v)||_{\psi_2} \leq ||\lambda||_{\psi_2} < \infty$.

Proof of Lemma F.1 (\mathbb{B}).

$$\mathbb{B} = \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} X^{\top} (y_1 a_1^{\top} + y_2 a^{\top}) \odot \sigma'_{\perp} (XW_0)$$
(20)

$$\|\mathbb{B}\| \leq \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X\| \|y_{1}a_{1}^{\top} + y_{2}a_{2}^{\top} \odot \sigma_{\perp}'(XW_{0})\|$$

$$\leq \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X\| \left(\|y_{1}a_{1}^{\top} \odot \sigma_{\perp}'(XW_{0})\| + \|y_{2}a_{2}^{\top} \odot \sigma_{\perp}'(XW_{0})\| \right)$$

$$\leq \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X\| \left(\|y_{1}\|_{\infty} \|\sigma_{\perp}'(XW_{0})\| \|a_{1}\|_{\infty} + \|y_{2}\|_{\infty} \|\sigma_{\perp}'(XW_{0})\| \|a_{2}\|_{\infty} \right)$$
(21)

1078
1079
$$= \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X\| \|\sigma'_{\perp}(XW_0)\|(\|a_1\|_{\infty} + \|a_2\|_{\infty})$$

 $\|\sigma'_{\perp}(XW_0)\|$ study Use the result of D.4 in Fan & Wang (2020), which is hold for orthogonal 1081 columns. X is sampled from continuous support distribution c_1, c_2 . The first vector is linearly 1082 independent with probability 1 due to the continuous support of its distribution. For the second 1083 vector, which is drawn independently, the probability that it lies in the span of the first vector is 0, 1084 as it also has a continuous density. This reasoning extends to n vectors, implying that, with high 1085 probability, they are orthogonal or nearly orthogonal because no vector falls into the span of the 1086 others. Thus, $\forall \mathbb{B} > 0$ following is hold.

$$\mathbb{P}(\{\|\sigma_{\perp}'\| \ge C(\sqrt{\mathbf{n}} + \sqrt{\mathbf{N}})\lambda_{\sigma}\mathbb{B}\}, \mathcal{A}_{\mathbb{B}}) \le 2e^{-c\mathbf{N}}$$
$$\mathcal{A}_{\mathbb{B}} = \{\{\|W_0\| \le \mathbb{B}\}, \{\sum_{i=1}^{\mathbf{N}} (\|W[i]\|^2 - 1)^2 \le \mathbb{B}^2\}\}.$$
(22)

¹⁰⁹² Therefore,

$$\mathbb{P}(\|\sigma_{\perp}^{\prime}\| \ge C(\sqrt{\mathbf{n}} + \sqrt{\mathbf{N}})\lambda_{\sigma}\mathbb{B}) \le 2e^{-c\mathbf{N}} + \mathbb{P}(\mathcal{A}_{\mathbb{B}}^{c})$$
(23)

1096 $\mathbb{P}(\mathcal{A}_{\mathbb{B}})$ study We choose $t = C\sqrt{\frac{d}{N}}, B = C\sqrt{\frac{d}{N}}$.

1098 CASE OF $||W_0|| \le B$ By Lemma E.3,

$$\mathbb{P}(\|\sqrt{\mathbf{N}}W_0\| \ge 2\sqrt{\mathbf{N}} + \sqrt{\mathbf{d}}) \le 2e^{-c\mathbf{N}} \Rightarrow \quad \mathbb{P}(\|W_0\| \ge C\sqrt{\frac{\mathbf{d}}{\mathbf{N}}}) \le 2e^{-c\mathbf{N}}$$
(24)

1103 Therefore, $||W_0|| \leq \mathbb{B}$ at least w.p. $1 - 2e^{-c\mathbf{N}}$

1107 CASE OF $\sum_{i=1}^{N} (\|W[i]\|^2 - 1)^2 \leq \mathbb{B}^2$ By definition, $\|W_0[i]\|^2 = 1$, so $0 \leq \mathbb{B}^2$ w.p. 1. We know $\mathbb{P}(\mathcal{A}_{\mathbb{B}}^c) \leq 2e^{-cN}$.

$$\mathbb{P}(\|\sigma_{\perp}'\| \ge C(\sqrt{\mathbf{n}} + \sqrt{\mathbf{N}})\sqrt{\frac{\mathbf{d}}{\mathbf{N}}}) \le 2e^{-c\mathbf{N}}$$
(25)

1112 Use Lemma F.3, and E.3,

$$\|\sigma_{\perp}'\| \le C\left(\sqrt{\frac{\mathbf{nN}}{\mathbf{d}}} + \sqrt{\frac{\mathbf{N}^2}{\mathbf{d}}}\right) \qquad \text{w.p. } 1 - C(e^{-c\mathbf{N}} + e^{-c\mathbf{d}}) \qquad (26)$$

$$\|a\|_{\infty} \le \frac{\iota}{\sqrt{\mathbf{N}}} \qquad \qquad \text{w.p. } 1 - 2\mathbf{N}e^{-ct^2} \qquad (27)$$

$$||X|| \le \sqrt{\mathbf{n}} + \sqrt{\mathbf{d}} + t'$$
 w.p. $1 - 2e^{-ct'^2}$. (28)

In summary, we get

Proof of Lemma F.1 (C). We know that σ' is bounded, so $\|\sigma'\|_F \le \lambda_{\sigma} \sqrt{nN}$

$$\mathbb{C} = -\frac{1}{\mathbf{nN}} X^{\top} \sigma(XW_0) \left(a_1 a_1^{\top} + a_2 a_2^{\top} \right) \odot \sigma'(XW_0), \tag{30}$$

ans we can bound the norm as follows

$$\begin{aligned} \|\mathbb{C}\| &\leq \frac{1}{\mathbf{nN}} \|X\| (\|\sigma a_1 a_1^\top \odot \sigma'\| + \|\sigma a_2 a_2^\top \odot \sigma'\|) \\ &\leq \frac{1}{\mathbf{nN}} \|X\| (\|\sigma a_1\|_{\infty} \|a_1\|_{\infty} \|\sigma'\|_F + \|\sigma a_2\|_{\infty} \|a_2\|_{\infty} \|\sigma'\|_F) \\ &\leq \frac{\lambda_{\sigma}}{\sqrt{-\mathbf{N}}} \|X\| (\|\sigma a_1\|_{\infty} \|a_1\|_{\infty} + \|\sigma a_2\|_{\infty} \|a_2\|_{\infty}) \end{aligned}$$
(31)

 $\leq \frac{1}{\sqrt{\mathbf{nN}}} \|X\| (\|\sigma a_1\|_{\infty} \|a_1\|_{\infty} + \|\sigma a_2\|_{\infty} \|a_2\|_{\infty} \|a_2\|_$

Control of $\|\sigma a\|_{\infty}$ Let $t = \sqrt{\mathbf{d}}$. Given X s.t. $\mathbb{P}(|X[i] - \sqrt{\mathbf{d}}| \ge \sqrt{\mathbf{d}}) \le 2e^{-ct^2}$, consider one 1144 element $\sigma(X[j]^\top W_0)a = \sum_i^{\mathbf{N}} a_i \sigma(x[j]^\top W_0[i])$.

1146 We know $a[i], \sqrt{\mathbf{n}}W_0[i]$ is an independent centered sub-Gaussian, and use Fact F.4, then 1147 $\sigma\left(\frac{X[j]^{\top}}{\sqrt{\mathbf{N}}}\sqrt{\mathbf{N}}W_0\right)a$ is sub-exponential and mean is zero, since $||a_i\sigma(x[j]^{\top}W_0[i])||_{\psi_1} \leq ||a_i||_{\psi_2}||\sigma(x[j]^{\top}W_0[i])||_{\psi_2} < \infty$. Apply the Bernstein inequality for the sub-exponential,

$$\mathbb{P}(|\sigma(X[j]^{\top}a)| \ge \log \mathbf{n} \text{ given } \{ |X[i] - \sqrt{\mathbf{d}}| \ge \sqrt{\mathbf{d}} \}) \le 2e^{-c \log^2 \mathbf{n}}.$$
(32)

.)

1152 For every element $\|\sigma(XW_0)a\|_{\infty} \le \log \mathbf{n}$ w.p. $1 - [2\mathbf{n}e^{-c\log^2 \mathbf{n} + 2\mathbf{n}e^{-c\mathbf{d}}}]$

1153 By Lemma F.3 $\mathbb{P}(||a||_{\infty} \le t/\sqrt{N}) \ge 1 - 2Ne^{-ct^2}$, and Lemma E.3 with $t = \sqrt{d}$

$$\mathbb{P}\left(\|\mathbb{C}\| \ge \frac{C}{\sqrt{\mathbf{n}}\mathbf{N}} (2\sqrt{\mathbf{d}} + \sqrt{\mathbf{n}}) \log \mathbf{n} \log \mathbf{N}\right) \le 2\left(\mathbf{n}e^{-c\mathbf{d}} + ne^{-c\log^{2}\mathbf{n}} + \mathbf{N}e^{-c\log^{2}\mathbf{n}}\right).$$
(33)

Remark F.5. In the proportional regime, as $n, d, N \rightarrow \infty$, these quantities can be interchanged to a constant. Thus, Lemma F.1 is reformulated as follows

$$\mathbb{P}(\|\mathbb{A}\| \le \kappa/\sqrt{\mathbf{n}}) \le Ce^{-c\mathbf{n}})$$

$$\mathbb{P}\left(\|\mathbb{B}\| \ge \frac{C\log\mathbf{N}}{\mathbf{n}}\right) \le C(e^{-c\mathbf{n}} + \mathbf{n}e^{-c\log^{2}\mathbf{n}})$$

$$\mathbb{P}\left(\|\mathbb{C}\| \ge \frac{C\log^{2}\mathbf{N}}{\mathbf{n}}\right) \le C(\mathbf{n}e^{-c\mathbf{n}} + \mathbf{n}e^{-c\log^{2}\mathbf{n}})$$
(34)

Proof of Proposition 3.1. Using $||G_0 - \mathbb{A}|| = ||\mathbb{B} + \mathbb{C}|| \le ||\mathbb{B}|| + ||\mathbb{C}||$ and Lemma F.5

$$\mathbb{P}\left(\|G_0 - \mathbb{A}\| \ge C \frac{\log^2 \mathbf{n}}{\mathbf{n}}\right) \le \mathbb{P}\left(\|G_0 - \mathbb{A}\| \ge C(\frac{\log n}{n} + \frac{\log^2 \mathbf{n}}{\mathbf{n}})\right) \le Cne^{-c\log^2 \mathbf{n}}.$$
 (35)

1172 Therefore, almost surely, in the proportional limit,

$$\|G_0 - \mathbb{A}\| \le C \frac{\log^2 \mathbf{n}}{\mathbf{n}} = \frac{\kappa}{\sqrt{\mathbf{n}}} \frac{C}{\kappa} \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} \le \|\mathbb{A}\| \frac{C}{\kappa} \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} \le \kappa' \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} (\|G_0\| + \|G_0 - \mathbb{A}\|).$$
(36)

1178 G ADDITIONAL RESULTS OF SECTION 3.2

A. $M_a \triangleq \max_{1 \le i \le \mathbf{N}} |a_i| \le \frac{C \log^{1/2} \mathbf{n}}{\sqrt{\mathbf{n}}}$ w.p $1 - 2ne^{-c \log \mathbf{n}}$

B. $M_b \triangleq \max_{1 \le i \le \mathbf{n}} |\langle \tilde{X}[i], \beta \rangle| \le C \log^{1/2} \mathbf{n}, w.p. \ 1 - 2\mathbf{n}e^{-c \log \mathbf{n}}$

1180 Lemma G.1. Given dataset \hat{D} , $\tilde{\hat{D}}$

C. $M_{W_0} \triangleq \sup_{k \ge 1} ||(W_0 W_0^\top)^{\circ k}|| \le C$

 $D. ||\tilde{X}|| \le C\sqrt{\mathbf{n}}$

E. $\sqrt{\mathbf{N}}||G|| = O_{\mathbb{P}}(1)$ F. $||A^{\circ k}|| \le ||A||^k$ Proof. For A, B, C, and D, we employ proof techniques adapted from Moniri et al. (2024). For A, B, by hoeffding inequality $\mathscr{P}(|X_i| \ge t) \le 2e^{-ct^2}$ for $t = \log^{1/2} \mathbf{n}$, and use $a_i, \langle \tilde{X}[i], \beta \rangle$ is Sub-Gaussian. For C, refer Moniri et al. (2024). For D, by Lemma E.3 and the proportional regime. For E, by Lemma F.1 $||G|| \le ||\mathbb{A}|| + ||\mathbb{B}|| + ||\mathbb{C}|| = O_{\mathbb{P}}(\frac{1}{\sqrt{n}} + \frac{\log n}{n} + \frac{\log^2 n}{n}) = O_{\mathbb{P}}(\frac{1}{\sqrt{n}})$ For F, refer Bai & Silverstein (2010) Corollary A.21. **Corollary G.2.** By Proposition 3.1 and D, E in Lemma G.1, we have w.p. 1 - o(1). $||\tilde{X}G^{\top} - \mu_1 \tilde{X}\beta \alpha^{\top}|| = O(\frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}})$ (37)*Remark* G.3. Remark $W_1 = W_0 + \eta \sqrt{\mathbf{n}}G$, so $\tilde{X}W_1 = \tilde{X}W_0 + \eta \sqrt{\mathbf{n}}\tilde{X}G$. *Proof of Lemma 3.2.* k = 1 is trivial with above statements. We follow Moniri et al. (2024) for $k \ge 2$. We need to show $\exists C > 0$, w.p. 1-o(1) $||(\tilde{X}G^{\top})^{\circ k} - c_1^k (\tilde{X}\beta)^{\circ k} (\alpha^{\circ k})^{\top}|| \le C^k \mathbf{n}^{-\frac{k}{2}\log^{2k} \mathbf{n}}$ (38) $(\tilde{X}G^{\top})^{\circ k} = (\tilde{X}G^{\top} - c_1\tilde{X}\beta\alpha^{\top} + c_1\tilde{X}\beta\alpha^{\top})^{\circ k}$ $=\sum_{i=1}^{\kappa} \binom{k}{j} c_1^{k-j} \operatorname{diag}(\tilde{X}\beta)^{\circ(k-j)} (\tilde{X}G^{\top} - c_1 \tilde{X}\beta\alpha)^{\circ j} \operatorname{diag}(\alpha)^{\circ(k-j)} \quad \text{by binomial theorem}$ $+ c_1^k \operatorname{diag}(\tilde{X}\beta)^{\circ k} \operatorname{diag}(\alpha)^{\circ k}$ (39)

Thus,
$$(\tilde{X}G^{\top})^{\circ k} - c_1^k (\tilde{X}\beta)^{\circ k} (\alpha^{\circ k})^{\top} = \sum_{j=1}^k \binom{k}{j} c_1^{k-j} \operatorname{diag}(\tilde{X}\beta)^{\circ (k-j)} (\tilde{X}G^{\top} - c_1 \tilde{X}\beta\alpha)^{\circ j} \operatorname{diag}(\alpha)^{\circ (k-j)}$$

$$(40)$$

We have to norm bound RHS of equation 40

1234 In summary, w.p. 1 - o(1)

$$||(\tilde{X}G^{\top})^{\circ k} - c_1^k (\tilde{X}\beta)^{\circ k} (\alpha^{\circ k})^{\top}|| \le C \sum_{j=1}^k \left(\frac{\log^{1/2} \mathbf{n}}{\sqrt{\mathbf{n}}}\right)^{k-j} \left(\frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}}\right)^j$$
(42)

Remark G.4. The definition of gradient G and the size of the norm are different between Moniri et al. (2024) and our paper, but both produce the same results as above, up to scaling factor $\frac{1}{\sqrt{N}}$.

1242 H STUDY OF EXPECTATION OF HERMITE POLYNOMIAL

The inner product between a random Gaussian vector x and a vector u, v, where u, v corresponds to a column of the weight matrix W or β , is substituted into the variable of a Hermite polynomial and its expectation is derived.

We have analyzed various macroscopic results regarding the feature space of a neural network using
Hermite polynomials and different activation functions. We have cited previously known facts, while
our derived results are presented without explicitly marking them as new. We believe these findings
will strengthen our paper and aid in the analysis of dynamics across different feature spaces.

1252 H.1 EXPECTATION OF A PRODUCT OF TWO HERMITE POLYNOMIALS

Here is the result of the expectation of the product of two Hermite polynomials, utilizing the orthog-onality of Hermite polynomials.

Lemma H.1 (Orthogonality of Hermite polynomials from Lemma C.1 Moniri et al. (2024)). See also derivation in Chapter 11.2 O'Donnell (2021).

1258 1259 1260 Let (Z_1, Z_2) be jointly Gaussian with $\mathbb{E}[Z_1] = \mathbb{E}[Z_2] = 0$, $\mathbb{E}[Z_1^2] = \mathbb{E}[Z_2^2] = 1$, and $\mathbb{E}[Z_1 Z_2] = \rho$. Then for any $k_1, k_2 \in \{0, 1, \dots, \}$

$$\mathbb{E}[H_{k_1}(Z_1)H_{k_2}(Z_2)] = k_1!\rho^{k_1}\mathbf{1}_{k_1=k_2}$$

1263 In the other form, for $d \in \mathbb{N}$, $Z \sim \mathcal{N}(0, I_d)$, $a, b \in \mathbb{S}^{d-1}$,

$$\mathbb{E}[H_{k_1}(Z^{\top}a)H_{k_2}(Z^{\top}b)] = k_1!(a^{\top}b)^{k_1}\mathbf{1}_{k_1=k_2}$$

Fact H.2. Let $W \in \mathbb{R}^{d \times N}$ s.t. $\forall i \ W[i] \in \mathbb{S}^{d-1}$. For $Z \sim \mathcal{N}(0, I)$,

$$\mathbb{E}_{Z \sim \mathcal{H}(0,1)}[H_j(W^\top Z)H_k(W^\top Z)^\top] = k!(W^\top W)^{\circ j}\mathbf{1}_{j=k}$$
(43)

1268 1269 1270

1282

1267

1261 1262

1264 1265

$$\mathbb{E}_{Z \sim \mathcal{N}(0,1)}[H_j(W^{\top}Z)^{\top}H_k(W^{\top}Z)] = k! \sum ||W[i]||^{2j} \mathbf{1}_{j=k} = k! N \mathbf{1}_{j=k}$$
(44)

1271 Proof. We apply H_j element-wise. By Lemma H.1, we can acquire the above result.

The following remark presents a modified condition of Lemma H.1 for the case where $a, b \notin \mathbb{S}^{d-1}$ in Lemma H.1. In this case, the variances of $Z^{\top}a$ and $Z^{\top}b$ are not equal to 1, and the covariance may exceed the bounds [-1, 1]. Under this condition, we will compute the expectation of the product of two Hermite polynomials as in Lemma H.1.

1277 *Remark* H.3 (the modified condition of Lemma H.1). For $d \in \mathbb{N}$, $u, v \in \mathbb{R}^d$, $Z \sim \mathcal{N}(0, I_d)$, 1278

1279
$$Z_1 = \langle u, Z \rangle \sim \mathcal{N}(0, ||u||_2^2), Z_2 = \langle v, Z \rangle \sim \mathcal{N}(0, ||v||_2^2)$$

1280 Then, Z_1, Z_2 is $\rho \triangleq \langle \frac{u}{||u||}, \frac{v}{||v||} \rangle$ - correlated 1281

$$corr(Z_1, Z_2) = \frac{\mathbb{E}[Z_1 Z_2]}{\sqrt{V(Z_1)}\sqrt{V(Z_2)}} = \frac{\mathbb{E}_Z \langle u, Z \rangle \langle v, Z \rangle}{||u|| ||v||}$$
$$= \frac{\mathbb{E}_g \sum_i \sum_j u_i v_j Z_i Z_j}{||u|| ||v||} = \frac{\sum_i \sum_j u_i v_j \mathbb{E}_Z[Z_i Z_j]}{||u|| ||v||}$$
$$= \frac{\langle u, v \rangle}{||u|| ||v||}$$
(45)

$$=\frac{(u,v)}{||u|| ||v||}$$

Additionally,

$$\begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} ||u||^2 & \langle u, v \rangle \\ \langle v, u \rangle & ||v||^2 \end{pmatrix} \right)$$
 (46)

We first introduce Isserlis' theorem, which is essential for the proof. This theorem allows the expectation of the product of centered Gaussian random variables to be expressed as a product of covariances, making the computation feasible. **Theorem H.4** (Isserlis' Theorem (Isserlis, 1918; Vignat, 2011)). Let $X = (X_1, \dots, X_d)$ Gaussian random vector s.t. $\mathbb{E}[X] = 0$, and let $A = \{\alpha_1, \dots, \alpha_N\}$ be set of integers s.t. $1 \le \alpha_i \le d, \forall i$. Denote $X_A = \prod_{\alpha_i \in A} X_{\alpha_i}$, and $X_{\emptyset} = 1$. Let $\prod(A)$ denote partitions of A into disjoint pairs and $\sigma \in \prod(A)$ is pair. (47)

$$\mathbb{E}[X_A] = \sum_{\sigma \in \prod(A)} \prod_{(i,j) \in \sigma} \mathbb{E}[X_{\alpha_i} X_{\alpha_j}] \mathbf{1}_{\mathrm{d} \text{ is even}}.$$
(47)

Now, we generalize the assumptions from the previous works so that Lemma H.1 holds for arbitrary vectors as Remark H.3. This could allow the weights of the networks to become analyzable when they go beyond the assumption of lying on the unit spheres.

Theorem H.5 (Generalization of Lemma H.1 for centered Gaussian distribution). For $d \in \mathbb{N}$, $u, v \in \mathbb{R}^d$, $g \sim \mathcal{N}(0, I_d)$, $\langle u, g \rangle \sim \mathcal{N}(0, ||u||_2^2)$, $\langle v, g \rangle \sim \mathcal{N}(0, ||v||_2^2)$.

1309 1310

1311

1312 1313 1314

1328 1329

1341 1342 1343

1349

1301 1302

$$\mathbb{E}_{g}[H_{j}(u^{\top}g)H_{k}(v^{\top}g)]$$

$$= \frac{j!\langle u, v \rangle^{j}}{||u||^{2}||v||^{2}} \mathbf{1}_{j=k} - \frac{(||u||^{2} - 1)(||v||^{2} - 1)}{||u||^{2}||v||^{2}} \mathbb{E}_{g}[(v^{\top}g)^{k}(u^{\top}g)^{j}] + \frac{(||v||^{2} - 1)}{||v||^{2}} \mathbb{E}_{g}[H_{j}(u^{\top}g)(v^{\top}g)^{k}] + \frac{(||u||^{2} - 1)}{||u||^{2}} \mathbb{E}_{g}[H_{k}(v^{\top}g)(u^{\top}g)^{j}]$$

$$(48)$$

Remark H.6. The same results can be derived as in Lemma H.1 when the variance is 1 in Thm. H.5.

Proof of Theorem H.5. (Generalize Chapter 11.2 O'Donnell (2021)'s derivation to non unit variance)

1320 $\mathbb{E}_{z \sim \mathcal{N}(0,\sigma^2)}[e^{tz}]$ study

First, we study about $\mathbb{E}_{g \sim \mathcal{N}(0,\sigma^2)}[e^{tg}]$ in order to analysis non unit variance case.

$$\mathbb{E}_{g \sim \mathcal{H}(0,\sigma^2)}[e^{tg}] = \frac{1}{\sqrt{2\pi\sigma}} \int e^{tg} e^{-\frac{g^2}{2\sigma^2}} dg$$

$$= \frac{1}{\sqrt{2\pi\sigma}} e^{\frac{1}{2}t^2} \int \exp(-\frac{(g-\sigma^2 t)^2}{2\sigma^2}) \quad \text{complete square}$$

$$= e^{\frac{1}{2}t^2}$$
(49)

1330 $\mathbb{E}_{Z,Z'}[\exp(sZ+tZ')]$ study

1331 Studying $\mathbb{E}_{Z,Z'}[\exp(sZ + tZ')]$, we can derive what we need to show.

$$\begin{aligned} & \text{1333} \\ & \text{1334} \\ & \text{IS35} \\ & \text{IS36} \\ & \text{IS36} \\ & \text{IS37} \\ & \text{IS37} \\ & \text{IS38} \\ & \text{IS39} \\ \end{aligned} \\ & \text{IS39} \\ & \text{IS39} \\ & \text{IS40} \\ & \text{IS40} \end{aligned} \\ & \text{IS37} \\ & \text{IS38} \\ & \text{IS39} \\ & \text{IS3$$

Therefore,

$$\exp(\langle u, v \rangle st) = \mathbb{E}_g[\exp(su^\top g - \frac{1}{2}s^2 ||u||^2) \exp(tv^\top g - \frac{1}{2}t^2 ||v||^2)]$$

1344
1345Facts for proof : one can verify below propositions with simple calculations.1346
1346Let $P_j(z) + z^j = H_j(z), C_u = ||u||^2 - 1, a > 0$ 1347
1348Let $f(s) = \exp(sz - \frac{1}{2}s^2), \bar{f}(s) = \exp(sz - \frac{1}{2}as^2)$

7.3.A. By Taylor expansion, $\exp(\langle u, v \rangle st) = \sum_{j=0}^{\infty} \frac{1}{j!} \langle u, v \rangle^j s^j t^j$.

1350 7.3.B. By Taylor expansion,
$$\bar{f}(s) = \sum_{j=0}^{\infty} \frac{1}{j!} \bar{f}^{(n)}(0) s^j$$

1352 7.3.C.
$$\bar{f}^{(n)}(0) = H_n(z) + C_u P_n(z)$$

1353

By using the fact that $\exp(\langle u, v \rangle st) = \mathbb{E}_g[\exp(su^\top g - \frac{1}{2}s^2||u||^2)\exp(tv^\top g - \frac{1}{2}t^2||v||^2)]$, we can eliminate the different orders of *s t* by a Taylor expansion and equating all monomials of the resulting polynomials.

$$j!\langle u, v \rangle^{j} \mathbf{1}_{j=k} = \mathbb{E}_{g} \Big[(H_{j}(u^{\top}g) + P_{j}(u^{\top}g)C_{u})(H_{j}(v^{\top}g) + P_{j}(v^{\top}g)C_{v}) \Big]$$

$$= \mathbb{E}_{g} \Big[(H_{j}(u^{\top}g) + (H_{j}(u^{\top}g) - (u^{\top}g)^{j})C_{u})(H_{j}(v^{\top}g) + (H_{j}(v^{\top}g) - (v^{\top}g)^{j})C_{v}) \Big]$$

$$= ||u||^{2} ||v||^{2} \mathbb{E}_{g} [H_{j}(u^{\top}g)H_{j}(v^{\top}g)] + (||u||^{2} - 1)(||v||^{2} - 1)\mathbb{E}_{g} [(v^{\top}g)^{j}(u^{\top}g)^{j}]$$

$$- ||u||^{2} (||v||^{2} - 1)\mathbb{E}_{g} [H_{j}(u^{\top}g)(v^{\top}g)^{j}] - ||v||^{2} (||u||^{2} - 1)\mathbb{E}_{g} [H_{j}(v^{\top}g)(u^{\top}g)^{j}]$$

$$(51)$$

Therefore,

$$\mathbb{E}_{g}[H_{j}(u^{\top}g)H_{j}(v^{\top}g)] = \frac{j!\langle u, v \rangle^{j}}{||u||^{2}||v||^{2}} \mathbf{1}_{j=k} - \frac{(||u||^{2} - 1)(||v||^{2} - 1)}{||u||^{2}||v||^{2}} \mathbb{E}_{g}[(v^{\top}g)^{j}(u^{\top}g)^{j}] + \frac{(||v||^{2} - 1)}{||v||^{2}} \mathbb{E}_{g}[H_{j}(v^{\top}g)(v^{\top}g)^{j}] + \frac{(||u||^{2} - 1)}{||u||^{2}} \mathbb{E}_{g}[H_{j}(v^{\top}g)(u^{\top}g)^{j}]$$
(52)

Note that the result of Lemma H.7 can be applied for concrete calculation, and conclude the proof.
1374
1375

Lemma H.7. For
$$d \in \mathbb{N}$$
, $u, v \in \mathbb{R}^d$, $g \sim \mathcal{N}(0, I_d)$, $\bar{Z}_1 = \langle u, g \rangle$, $\bar{Z}_2 = \langle v, g \rangle$.

$$\begin{pmatrix} \bar{Z}_1 \\ \bar{Z}_2 \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} ||u||^2 & \langle u, v \rangle \\ \langle v, u \rangle & ||v||^2 \end{pmatrix}\right)$$
(53)

 X_{α_i} is defined at Thm. H.4

$$\mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[H_{j}(\bar{Z}_{1})\bar{Z}_{2}^{k}] = j! \sum_{m=0}^{\lfloor \frac{j}{2} \rfloor} \frac{(-1)^{m}}{m!(j-2m)!2^{m}} \sum_{\sigma \in \prod(\{\{\bar{Z}_{1}\}\times j-2m\}\cup\{\{\bar{Z}_{2}\}\times k\}\})} \prod_{(p,q)\in\sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j+k-2m \text{ is even}} \\ \mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \sum_{\sigma \in \prod(\{\{\bar{Z}_{1}\}\times j\}\cup\{\{\bar{Z}_{2}\}\times k\}\})} \prod_{(p,q)\in\sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j+k \text{ is even}} \\ \mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \sum_{\sigma \in \prod(\{\{\bar{Z}_{1}\}\times j\}\cup\{\{\bar{Z}_{2}\}\times k\}\})} \prod_{(p,q)\in\sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j+k \text{ is even}} \\ \mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \sum_{\sigma \in \prod(\{\{\bar{Z}_{1}\}\times j\}\cup\{\{\bar{Z}_{2}\}\times k\}\})} \prod_{(p,q)\in\sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j+k \text{ is even}} \\ \mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \sum_{\sigma \in \prod(\{\{\bar{Z}_{1}\}\times j\}\cup\{\{\bar{Z}_{2}\}\times k\}\})} \prod_{(p,q)\in\sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j+k \text{ is even}} \\ \mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \sum_{\sigma \in \prod(\{\{\bar{Z}_{1}\}\times j\}\cup\{\{\bar{Z}_{2}\}\times k\}\})} \prod_{(p,q)\in\sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j+k \text{ is even}} \\ \mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \sum_{\sigma \in \prod(\{\{\bar{Z}_{1}\}\times j\}\cup\{\{\bar{Z}_{2}\}\times k\}\})} \prod_{(p,q)\in\sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j+k \text{ is even}} \\ \mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \sum_{\sigma \in \prod(\{\bar{Z}_{1},\bar{Z}_{2}\}\times k\}} \prod_{(p,q)\in\sigma} \mathbb{E}_{\bar{Z}_{2},\bar{Z}_{2}}[\bar{Z}_{1},\bar{Z}_{2},\bar{Z}$$

Proof. By explicit formula of Hermite polynomials

$$\mathbb{E}_{\bar{Z}_1,\bar{Z}_2}[H_j(\bar{Z}_1)(\bar{Z}_2)^k] = j! \sum_{m=0}^{\lfloor \frac{j}{2} \rfloor} \frac{(-1)^m}{m!(j-2m)!2^m} \mathbb{E}_{\bar{Z}_1,\bar{Z}_2}[\bar{Z}_1^{j-2m}\bar{Z}_2^k]$$
(55)

1397 Therefore, we need to figure out $\mathbb{E}_{\bar{Z}_1,\bar{Z}_2}[\bar{Z}_1^p \bar{Z}_2^q]$. We know \bar{Z}_1,\bar{Z}_2 is mean zero Gaus-1398 sian, so we can apply Thm. H.4 with $A = \{\{\bar{Z}_1\} \times p\} \cup \{\{\bar{Z}_2\} \times q\}\}, \mathbb{E}[\bar{Z}_1^p \bar{Z}_2^q] = \sum_{\sigma \in \prod(A)} \prod_{(\tau,\upsilon) \in \sigma} \mathbb{E}[X_{\alpha_\tau} X_{\alpha_\upsilon}] \cdot \mathbf{1}_{p+q \text{ is even}}$

1402 Corollary H.8 (Corollary of Lemma H.7). Remark $Z_1 \sim \mathcal{N}(0, ||u||^2)$ For the case k = 0,

$$\mathbb{E}_{\bar{Z}_1}[\bar{Z}_1^j] = \|u\|^j (j-1)!! \mathbf{1}_{j \text{ is even}}$$
(56)

Proof.

=

H.2 EXPECTATION OF A PRODUCT OF TWO HERMITE POLYNOMIALS—GENERALIZATION TOWARD NON-CENTERED GAUSSIAN

 $\sum_{\sigma \in \prod(\{\bar{Z}_1\} \times j\})} \prod_{(p,q) \in \sigma} \|u\|^2 \mathbf{1}_{\mathbf{j} \text{ is even}} = \sum_{\sigma \in \prod(\{\bar{Z}_1\} \times j\})} \|u\|^j \mathbf{1}_{\mathbf{j} \text{ is even}} = (j-1)!! \|u\|^j \mathbf{1}_{\mathbf{j} \text{ is even}}$

 $\mathbb{E}_{\bar{Z}_1,\bar{Z}_2}[\bar{Z}_1^j\bar{Z}_2^k] = \mathbb{E}_{\bar{Z}_1}[\bar{Z}_1^j] = \sum_{\sigma \in \prod(\{\bar{Z}_1\} \times j\})} \prod_{(p,q) \in \sigma} \mathbb{E}[X_{\alpha_p} X_{\alpha_q}] \mathbf{1}_{j \text{ is even}}$

We will change Theorem H.5 and Lemma H.7 to adopt a generalized Gaussian assumption with a mean of zero.

Lemma H.9 (Taylor expansion of Hermite polynomials from Lemma C.2 Moniri et al. (2024)). For any $k_1, k_2 \in \{0, 1, \dots, \}$ and $x, y \in \mathbb{R}$,

$$H_k(x+y) = \sum_{j=0}^k \binom{k}{j} x^j H_{k-j}(y).$$
 (58)

(57)

Theorem H.10 (Generalization of Thm. H.5 for any Gaussian distribution). For $d \in \mathbb{N}$, $u, v \in \mathbb{R}^d$, $\xi \sim \mathcal{N}(0,1), g \sim \mathcal{N}(\mu, \Sigma), Z_1 = \langle u, g \rangle \sim \mathcal{N}(\mu^\top u, u^\top \Sigma u), Z_2 = \langle v, g \rangle \sim \mathcal{N}(\mu^\top v, v^\top \Sigma v).$

Proof of Theorem H.10. By reparametrization i.e. $Z_1 = \sqrt{u^{\top}\Sigma u}\xi + u^{\top}\mu, Z_2 = \sqrt{v^{\top}\Sigma v}\xi + v^{\top}\mu,$ and Lemma H.9,

$$H_j(\sqrt{u^{\top}\Sigma u}\xi + u^{\top}\mu) = \sum_{\alpha=0}^j \binom{j}{\alpha} (u^{\top}\mu)^{\alpha} H_{j-\alpha}(\sqrt{\mu^{\top}\Sigma u}\xi).$$
(60)

$$\mathbb{E}_{g}[H_{j}(u^{\top}g)H_{k}(v^{\top}g)] = \mathbb{E}_{\xi}[H_{j}(\sqrt{u^{\top}\Sigma u}\xi + u^{\top}\mu)H_{k}(\sqrt{v^{\top}\Sigma v}\xi + v^{\top}\mu)]$$
$$= \mathbb{E}_{\xi}\Big[\sum_{\alpha=0}^{j} \binom{j}{\alpha}(u^{\top}\mu)^{\alpha}H_{j-\alpha}(\sqrt{\mu^{\top}\Sigma u}\xi)\Big]\Big[\sum_{\beta=0}^{k} \binom{k}{\beta}(v^{\top}\mu)^{\beta}H_{k-\beta}(\sqrt{\mu^{\top}\Sigma v}\xi)\Big]$$
(61)

$$=\sum_{\alpha=0}^{j}\sum_{\beta=0}^{k} \binom{j}{\alpha} \binom{k}{\beta} (u^{\top}\mu)^{\alpha} (v^{\top}\mu)^{\beta} \mathbb{E}_{\xi} [H_{j-\alpha}(\sqrt{\mu^{\top}\Sigma u}\xi)H_{k-\beta}(\sqrt{\mu^{\top}\Sigma v}\xi)]$$

Use same proof technique Thm. H.5, with $\begin{pmatrix} \sqrt{u^{\top}\Sigma u}\xi\\ \sqrt{v^{\top}\Sigma v}\xi \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0\\0 \end{pmatrix}, \begin{pmatrix} u^{\top}\Sigma u & u^{\top}\Sigma v\\ v^{\top}\Sigma u & v^{\top}\Sigma v \end{pmatrix}\right)$

$$\mathbb{E}_{\xi}[H_{j-\alpha}(\sqrt{u^{\top}\Sigma u}\xi)H_{k-\beta}(\sqrt{v^{\top}\Sigma v}\xi)] \\
= \frac{(j-\alpha)!(u^{\top}\Sigma v)^{j-\alpha}}{u^{\top}\Sigma uv^{\top}\Sigma v}\mathbf{1}_{j-\alpha=k-\beta} - \frac{(u^{\top}\Sigma u-1)(v^{\top}\Sigma v-1)}{u^{\top}\Sigma uv^{\top}\Sigma v}\mathbb{E}_{g}[(\sqrt{u^{\top}\Sigma u}\xi)^{j-\alpha}(\sqrt{v^{\top}\Sigma v}\xi)^{k-\beta}] \\
+ \frac{(v^{\top}\Sigma v-1)}{v^{\top}\Sigma v}\mathbb{E}_{g}[H_{j-\alpha}(\sqrt{u^{\top}\Sigma u}\xi)(\sqrt{v^{\top}\Sigma v}\xi)^{k-\beta}] + \frac{(u^{\top}\Sigma u-1)}{u^{\top}\Sigma u}\mathbb{E}_{g}[(\sqrt{u^{\top}\Sigma u}\xi)^{j-\alpha}H_{k-\beta}(\sqrt{v^{\top}\Sigma v}\xi)] \\$$
(62)

In summery,

$$\mathbb{E}_{g}[H_{j}(u^{\top}g)H_{k}(v^{\top}g)] = \sum_{\alpha=0}^{j}\sum_{\beta=0}^{k} {j \choose \alpha} {k \choose \beta} (u^{\top}\mu)^{\alpha} (v^{\top}\mu)^{\beta} \\
\times \left[\frac{(j-\alpha)!(u^{\top}\Sigma v)^{j-\alpha}}{u^{\top}\Sigma uv^{\top}\Sigma v} \mathbf{1}_{j-\alpha=k-\beta} - \frac{(u^{\top}\Sigma u-1)(v^{\top}\Sigma v-1)}{u^{\top}\Sigma uv^{\top}\Sigma v} \mathbb{E}_{\xi}[(\sqrt{u^{\top}\Sigma u}\xi)^{j-\alpha}(\sqrt{v^{\top}\Sigma v}\xi)^{k-\beta}] \\
+ \frac{(v^{\top}\Sigma v-1)}{v^{\top}\Sigma v} \mathbb{E}_{\xi}[H_{j-\alpha}(\sqrt{u^{\top}\Sigma u}\xi)(\sqrt{v^{\top}\Sigma v}\xi)^{k-\beta}] + \frac{(u^{\top}\Sigma u-1)}{u^{\top}\Sigma u} \mathbb{E}_{\xi}[(\sqrt{u^{\top}\Sigma u}\xi)^{j-\alpha}H_{k-\beta}(\sqrt{v^{\top}\Sigma v}\xi)]$$
(63)
$$\Box$$

The following Corollary which calculates the Expectation of the Power of a Gaussian Random Variable can be derived using the binomial expansion with the reparametrization technique and Corollary H.8. It corresponds to the case k = 0 in Lemma H.7.

Corollary H.11 (Corollary of Lemma H.7). Given β , let Gaussian Random Variable $Z \sim \mathcal{N}(\mu^{\top}\beta, \beta^{\top}\Sigma\beta)$, then expectation of power of Z is

(64)

 $\mathbb{E}_{Z}(Z)^{k} = \sum_{t=0}^{k} \binom{k}{t} (\mu^{\top}\beta)^{k-t} \mathbb{E}_{\bar{Z} \sim \mathcal{H}(0,\beta^{\top}\Sigma\beta)}[\bar{Z}^{t}]$

The following corollary, which computes the Gaussian expectation of Hermite polynomials, is derived from the explicit form of Hermite polynomials and Corollary H.11. It corresponds to the case k = 0 in Theorem H.10.

 $=\sum_{t=0}^{k} \binom{k}{t} (\mu^{\top}\beta)^{k-t} (t-1)!! \cdot (\beta^{\top}\Sigma\beta)^{\frac{t}{2}} \mathbf{1}_{t \text{ is even }}.$

Corollary H.12 (Corollary of Theorem H.10). For $d \in \mathbb{N}$, given $w \in \mathbb{R}^d$, $x \sim \mathcal{N}(\mu, \Sigma)$,

$$\mathbb{E}_{x}[H_{n}(w^{\top}x)] = \sum_{m=0}^{\lfloor \frac{n}{2} \rfloor} \sum_{i=0}^{n-2m} \frac{(-1)^{m}(i-1)!! n!}{2^{m} m!(n-2m)!} \binom{n-2m}{i} (w^{\top}\mu)^{n-2m-i} (w^{\top}\Sigma w)^{\frac{i}{2}} \mathbf{1}_{i \text{ is even}}$$
(65)

I DETAIL OF ALIGNMENT ANALYSIS

Proof of Theorem 4.2. Let $\mathfrak{a}_k \triangleq c_1^k c_k \eta^k$.

When $x \sim \mathcal{N}(\mu, \Sigma)$, we approximate F to dominant term F_l , then $\mathbb{E}_{x,x',\theta}[F_l(x)^\top F_l(x')]$ $= \mathbb{E}_{x,x',\theta}[(F_0(x) + \sum_{l=1}^{l} \mathbf{n}_k(\beta^\top x)^k \alpha^{\circ k})^\top (F_0(x') + \sum_{l=1}^{l} \mathbf{n}_k(\beta^\top x')^k \alpha^{\circ k})^\top]$ $= \mathbb{E}_{x,x',\theta} \left| \langle F_0(x), F_0(x') \rangle \right.$ $+ \langle F_0(x), \sum_{k=1}^{l} \mathbf{a}_k(\beta^\top x')^k \alpha^{\circ k} \rangle + \langle F_0(x'), \sum_{k=1}^{l} \mathbf{a}_k(\beta^\top x)^k \alpha^{\circ k} \rangle + \sum_{k=1}^{l} \sum_{j=1}^{l} \mathbf{a}_k \mathbf{a}_j(\beta^\top x)^k (\beta^\top x')^j \sum_{r} \alpha[r]^{j+k} \right]$ $\triangleq \mathscr{A} + \mathscr{B} + \mathscr{C}$ (66)I.1 $\mathscr{A}: \mathbb{E}_{\theta} [\mathbb{E}_x[\sigma(W_0^{\top}x)]^{\top} \mathbb{E}_{x'}[\sigma(W_0^{\top}x')]]$ STUDY Let $\Box_{n,m,i} \triangleq \frac{(-1)^m (i-1)!! n!}{2^m m! (n-2m)!} {n-2m \choose i}.$ By Corollary H.12, $\mathscr{A} = \mathbb{E}_{\theta} \Big[\sum_{a=1}^{\mathbf{N}} \sum_{n=1}^{\infty} \sum_{m=0}^{\lfloor \frac{n}{2} \rfloor} \sum_{i=0}^{n-2m} \sum_{c=1}^{\infty} \sum_{n=0}^{\lfloor \frac{n}{2} \rfloor} \sum_{a=0}^{o-2p} c_n c_o \beth_{n,m,i} \beth_{o,p,q} \Big]$ (67) $\times (W_0[a]^{\top} \mu)^{n-2m-i} (W_0[a]^{\top} \Sigma W_0[a])^{i/2} \mathbf{1}_{i \text{ is even}}$ $\times (W_0[a]^\top \mu)^{o-2p-q} (W_0[a]^\top \Sigma W_0[a])^{q/2} \mathbf{1}_{q \text{ is even}}$ We know, $W_0 \sim \mathbb{R}^{\mathbf{d} \times \mathbf{N}}$, $W_0[i] \sim \text{Unif}(\mathbb{S}^{\mathbf{d}-1})$. Denote $w \triangleq W_0[a], t = n + o - 2m - 2p - i - q, k = \frac{i+q}{2}$ $\mathbb{E}[(W_0[a]^{\top}\mu)^{n-2m-i}(W_0[a]^{\top}\Sigma W_0[a])^{i/2}(W_0[a]^{\top}\mu)^{o-2p-q}(W_0[a]^{\top}\Sigma W_0[a])^{q/2}]$ (68) $= \mathbb{E}[(w^{\top}\mu)^t (w^{\top}\Sigma w)^k]$ Use covariance matrix property, which is diagonalizable i.e. $\Sigma = Q\Lambda Q$. $w^{\top}\Sigma w = w^{\top}Q\Lambda Q^{\top}w$. Let $z = Q^{\top}w$ and $\tilde{\mu} = Q^{\top}\mu$ then $w^{\top}\Sigma w = z^{\top}\Lambda z = \sum_{i}\lambda_{i}z_{i}^{2}$. By symmetry, $z \sim \text{Unif}(\mathbb{S}^{d-1})$ Therefore, using multi-index notation, where $|\alpha| = \sum_{i=1}^{n} \alpha_i, \alpha_i \ge 0$, and $\binom{t}{\alpha} = \frac{t!}{\alpha_1!\alpha_2!\cdots\alpha_n!}$ $\mathbb{E}_w[(w^\top \mu)^t (w^\top \Sigma w)^k] = \mathbb{E}_z[(z^\top \tilde{\mu})^t (\sum_i \lambda_i z_i^2)^k]$ $=\mathbb{E}_{z}[(\sum_{i=1}^{k} \binom{t}{\alpha} \prod_{i=1}^{d} (z_{i}\tilde{\mu}_{i})^{\alpha_{i}})(\sum_{j>1=1}^{k} \binom{k}{\beta} \prod_{i=1}^{d} (z_{j}^{2}\lambda_{j})^{\beta_{j}})$ (69) $=\sum_{|\alpha|=4}\sum_{|\alpha|=1}\binom{t}{\alpha}\binom{k}{\beta}\prod_{i=1}^{d}\prod_{j=1}^{d}\tilde{\mu}_{i}^{\alpha_{i}}\lambda_{j}^{\beta_{j}}\mathbb{E}_{z}[z_{i}^{\alpha_{i}}z_{j}^{2\beta_{j}}]$ The term related to μ of \mathscr{A} is associated with the random value W[a]. Therefore, taking expecta-tion on network parameters, \mathscr{A} only depends on unseen distribution parameter μ, Σ without train distribution. In summery, let $S(r, s, i, j) = \mathbb{E}_{z}[z_{i}^{r} z_{i}^{s}],$

 $R(n,m,i,o,p,q,\alpha,\beta,c_n,c_o,\mathbf{N}) = \mathbf{N}c_nc_o \beth_{n,m,i} \beth_{o,p,q} \mathbf{1}_{i,q \text{ are even}} \binom{n+o-2m-2p-i-q}{\alpha} \binom{\frac{i+q}{2}}{\beta}, \text{ which are deterministic function, then}$

 $\mathscr{A} = \sum_{n=1}^{\infty} \sum_{m=0}^{\lfloor \frac{n}{2} \rfloor} \sum_{i=0}^{n-2m} \sum_{o=1}^{\infty} \sum_{p=0}^{\lfloor \frac{o}{2} \rfloor} \sum_{q=0}^{o-2p} \sum_{|\alpha|=n+o-2m-2p-i-q} \sum_{|\beta|=\frac{i+q}{2}}$ $R(n, m, i, o, p, q, \alpha, \beta, c_n, c_o, \mathbf{N}) \prod_{l=1}^{d} \prod_{i=1}^{d} \tilde{\mu}_l^{\alpha_l} \lambda_j^{\beta_j} S(\alpha_l, \beta_j, l, j)$ I.2 $\mathscr{B}: 2\sum_{k=1}^{l} \mathfrak{B}_k \langle \mathscr{B}_k, \mathbb{E}_{\alpha}[a^{\circ k}] \rangle$ STUDY $\mathcal{B}_k: \mathbb{E}_{x,x'}[\sigma(W_0^\top x)(\beta^\top x')^k] \text{ study} \\ \text{Let } Z_1 = \langle W_0[a], x \rangle, Z_2 = \langle \beta, x' \rangle, \end{cases}$ then $Z_1|W_0[a] \sim \mathcal{N}(W_0[a]^\top \mu, W_0[a]^\top \Sigma W_0[a]), Z_2 \sim \mathcal{N}(\beta^\top \mu, \beta^\top \Sigma \beta).$ Therefore, by Corollary H.11 and H.12 $\mathscr{B}_k[a] = \mathbb{E}_{\theta} \sum_{i=1}^{\infty} c_j \mathbb{E}_{Z_1} H_j(Z_1) \mathbb{E}_{Z_2}(Z_2)^k$ $=\sum_{i=1}^{\infty}\sum_{j=1}^{\lfloor\frac{j}{2}\rfloor}\sum_{i=1}^{j-2m}\sum_{k=0}^{k}\binom{j-2m}{i}\binom{k}{t}\frac{(-1)^{m}c_{j}(i-1)!!}{2^{m}m!(j-2m)!}\mathbf{1}_{k-t\text{ is even }}\mathbf{1}_{i\text{ is even }}\mathbf{1}_{i\text{$ $\mathbb{E}_{\theta}\left[(W_0[a]^{\top} \mu)^{j-2m-i} (W_0[a]^{\top} \Sigma W_0[a])^{\frac{i}{2}} \right] (\mu^{\top} \beta)^t \cdot (\beta^{\top} \Sigma \beta)^{\frac{k-t}{2}}$ $=\sum_{i=1}^{\infty}\sum_{m=0}^{\lfloor\frac{j}{2}\rfloor}\sum_{i=0}^{j-2m}\sum_{t=0}^{k}\sum_{|\alpha|=j-2m-i}\sum_{|\beta|=\frac{i}{\pi}}\binom{j-2m}{i}\binom{k}{t}\binom{j-2m-i}{\alpha}\binom{\frac{j}{2}}{\beta}\frac{(-1)^{m}c_{j}(i-1)!!}{2^{m}m!(j-2m)!}$ $\prod_{u=1}^{d} \prod_{u=1}^{d} \tilde{\mu}_{u}^{\alpha_{u}} \lambda_{v}^{\beta_{v}} S(\alpha_{u}, \alpha_{v}, u, v) \ (\mu^{\top} \beta)^{t} \cdot (\beta^{\top} \Sigma \beta)^{\frac{k-t}{2}} \mathbf{1}_{k-t \text{ is even }} \mathbf{1}_{i \text{ is even }}$ The term related to μ of $\mathscr{B}_k[a]$ is associated with the random value W[a]. Therefore, $\forall a, \mathscr{B}_k[a]$ depends on unseen distribution parameter μ, Σ and $\beta^{+}\mu$ with same value. $\mathbb{E}_{\alpha}[a^{\circ k}]$ study We know $a_1[i], a_2[i] \sim \mathcal{N}(0, \frac{1}{N})$, so $\alpha[i] \triangleq (a_1 - a_2)[i] \sim \mathcal{N}(0, \frac{2}{N})$. Therefore, by centered gaussian moments,

$$\mathbb{E}_{\alpha}[\alpha[r]^{k}] = \frac{(k)!}{w^{\frac{k}{2}}(\frac{k}{2})!} (\frac{2}{N})^{\frac{k}{2}} \mathbf{1}_{k \text{ is even}}$$
(72)

(70)

(71)

Since $\mathbb{E}_{\alpha}[\alpha[r]^{k}]$ is nonzero only when k is even, and even condition of k-t is exist in $\mathscr{B}_{k}[a]$, taking the absolute value of $\beta^T \mu$ within \mathscr{B}_k produce equivalent results.

Therefore, in $\mathscr{B} = 2 \sum_{k=1}^{l} \mathfrak{a}_k \langle \mathscr{B}_k, \mathbb{E}_{\alpha}[a^{\circ k}] \rangle$ is depends on $\mu, \Sigma, |\beta^{\top} \mu|$ and $\beta^{\top} \Sigma \beta$

1610 I.3
$$\mathscr{C}: \sum_{k=1,j=1}^{l} \mathfrak{a}_k \mathfrak{a}_j \mathbb{E}_{\alpha}[\sum_r \alpha[r]^{j+k}] \mathscr{C}_{j,k}$$
 STUDY

 $\mathscr{C}_{j,k}$: $\mathbb{E}_{x,x'}[(\beta^{\top}x)^k(\beta^{\top}x')^j]$ study Let $Z_1 \triangleq \beta^\top x \sim \mathcal{N}(\beta^\top \mu, \beta^\top \Sigma \beta)$, same as Z'_2 . Using Corollary H.11, $\mathbb{E}Z_{1}^{j}Z_{2}^{k} = \sum_{j}^{j}\sum_{k}^{k} \binom{j}{s} \binom{k}{t} (j-s-1)!!(k-t-1)!!(\mu^{\top}\beta)^{s+t} (\beta^{\top}\Sigma\beta)^{\frac{j-s+k-t}{2}} \mathbf{1}_{j-s, k-t \text{ are even}}$ (73)

Therefore, the term related to μ in $\mathscr{C}_{i,k}$ is only dominated by the discriminative data β , independent of the randomly initialized parameters.

 $\sum_{r} \mathbb{E}_{\alpha}[\alpha[r]^{j+k}]$ study

$$\sum_{r} \mathbb{E}_{\alpha}[\alpha[r]^{j+k}] = \frac{\mathbf{N}(j+k)!}{2^{\frac{j+k}{2}}(\frac{j+k}{2})!} \left(\frac{2}{\mathbf{N}}\right)^{\frac{j+k}{2}} \mathbf{1}_{j+k \text{ and } i \text{ is even}}$$
(74)

Since $\mathbb{E}_{\alpha}[\alpha[r]^{j+k}]$ is nonzero only when j+k is even, and even condition of j-s and k-t are exist in $\mathscr{C}_k[a]$, so s + t is even in this conditions, taking the absolute value of $\beta^T \mu$ within \mathscr{C}_k produce equivalent results.

1628 Therefore, in
$$\mathscr{C} = \sum_{k=1,j=1}^{l} \alpha_k \alpha_j \mathbb{E}_{\alpha} [\sum_r \alpha[r]^{j+k}] \mathscr{C}_{j,k}$$
 is depends on $|\beta^{\top}\mu|$ and $\beta^{\top} \Sigma \beta$
1630

DETAIL OF LOCAL ELASTICITY ANALYSIS J

Proof of Theorem 4.3.

$$\mathbb{E}_{x,\theta}||F_l(x) - F_0(x)||^2 = \mathbb{E}_{x,\theta}[\sum_{k=1}^l c_1^k c_k \eta^k (\beta^\top x)^k (\alpha^{\circ k})]^\top [\sum_{m=1}^l c_1^m c_m \eta^m (\beta^\top x)^m (\alpha^{\circ m})].$$
(75)

For $\mathbb{E}[(x^{\top}\beta)^{k+m}]$, by Corollary H.11

$$\mathbb{E}_{x}[(x^{\top}\beta)^{k+m}] = \mathbb{E}_{z\sim n(0,1)}[(\beta^{\top}\mu + \sqrt{\beta^{\top}\Sigma\beta}z)^{k+m}]$$

$$= \sum_{i=0}^{k+m} {\binom{k+m}{i}} (\beta^{\top}\mu)^{k+m-i} (\beta^{\top}\Sigma\beta)^{\frac{i}{2}} \mathbb{E}[z^{i}]$$

$$= \sum_{i=0}^{k+m} {\binom{k+m}{i}} (\beta^{\top}\mu)^{k+m-i} (\beta^{\top}\Sigma\beta)^{\frac{i}{2}} (i-1)!! \mathbf{1}_{i \text{ is even}}$$
(76)

Remark $\mathfrak{a}_k \triangleq c_1^k c_k \eta^k$. Finally, $\mathbb{E}_{x,\theta}||F_l(x) - F_0(x)||^2$ $=\sum_{k=1}^{l}\sum_{m=1}^{l}c_{1}^{k+m}c_{k}c_{m}\eta^{k+m}\sum_{i=0}^{k+m}\binom{k+m}{i}(\beta^{\top}\mu)^{k+m-i}(\beta^{\top}\Sigma\beta)^{\frac{i}{2}}(i-1)!!\mathbf{1}_{\mathrm{i}\,\mathrm{is\,even}}\mathbb{E}_{\alpha}[\alpha^{\circ k^{\top}}\alpha^{\circ m}]$ $=\sum_{k=1}^{l}\sum_{m=1}^{l}\sum_{i=0,mm}^{k+m} a_{k}a_{m}\binom{k+m}{i}(\beta^{\top}\mu)^{k+m-i}(\beta^{\top}\Sigma\beta)^{\frac{i}{2}}(i-1)!!\sum_{r}\mathbb{E}_{\alpha}[\alpha[r]^{k+m}].$ (77)

Taking Expectation over Network parameters, one can acquire

$$=\sum_{k=1}^{l}\sum_{m=1}^{l}\sum_{i=0}^{k+m}\mathbf{n}_{k}\mathbf{n}_{m}\binom{k+m}{i}(\beta^{\top}\mu)^{k+m-i}(\beta^{\top}\Sigma\beta)^{\frac{i}{2}}(i-1)!!\frac{\mathbf{N}(k+m)!}{2^{\frac{k+m}{2}}(\frac{k+m}{2})!}(\frac{2}{\mathbf{N}})^{\frac{k+m}{2}}\mathbf{1}_{k+m \text{ and } i \text{ is even}}$$
(78)

Therefore, k + m - i is even. For clarity, we use absolute values,

$$=\sum_{k=1}^{l}\sum_{m=1}^{l}\sum_{i=0}^{k+m}\mathfrak{a}_{k}\mathfrak{a}_{m}\binom{k+m}{i}\frac{\mathbf{N}(k+m)!}{2^{\frac{k+m}{2}}(\frac{k+m}{2})!}\left(\frac{2}{\mathbf{N}}\right)^{\frac{k+m}{2}}(i-1)!!|\beta^{\top}\mu|^{k+m-i}(\beta^{\top}\Sigma\beta)^{\frac{i}{2}}\mathbf{1}_{k+m \text{ and } i \text{ is even}}$$
(79)

For clearity, we define constant

$$\kappa_{LE}(k,m,i,\mathbf{N},c_1,c_k,c_m,\eta) \triangleq \mathbf{x}_k \mathbf{x}_m \binom{k+m}{i} \frac{\mathbf{N}(k+m)!}{2^{\frac{k+m}{2}}(\frac{k+m}{2})!} \left(\frac{2}{\mathbf{N}}\right)^{\frac{k+m}{2}} (i-1)!!$$

which depends on constant $k, m, i, \mathbf{N}, c_1, c_k, c_m, \eta$.

$$=\sum_{k=1}^{l}\sum_{m=1}^{l}\sum_{i=0}^{k+m}\kappa_{LE} |\beta^{\top}\mu|^{k+m-i}(\beta^{\top}\Sigma\beta)^{\frac{i}{2}} \mathbf{1}_{k+m \text{ and } i \text{ is even}}$$
(80)

 κ_{LE} depends only on the constants $k, m, i, N, c_1, c_k, c_m, \eta$, and is independent of the parameters of the data distribution.

ADDITIONAL INFORMATION OF EXPERIMENT 1, 2 Κ

ADDITIONAL RESULTS FOR ALIGNMENT AND ELASTICITY K.1





Figure K.2: Experiment 2 The observation of Alignment and Elasticity (y-axis), derived from the LHS of Thms. 1.1 and 1.2, across different values of $\beta^{\top}\mu$ (x-axis) with varying R.



Figure K.3: **Experiment 1** The x-axis is displayed on a logarithmic scale. Observation of Alignment and LE (y-axis) derived from the LHS of Thm. 4.2, 4.3 across different *e* (blue, lower x-axis) and $\beta^{\top}\mu$ (red, upper x-axis) values.

1782 K.2 ADDITIONAL RESULTS FOR RECALL@1 1783



Figure K.4: Recall@1 (y-axis) measurement of **Exper 1, 2** of features across different $\beta^{\perp}\mu$ values 1805 (x-axis). The blue line represents the clustering performance measured using the features in their 1806 initialized state, the orange line reflects the performance after one step of training, and the green line 1807 indicates the improvement, i.e., the difference between the two. For Setup 1 (top), the x-axis is on a 1808 logarithmic scale, whereas for Setup 2 (bottom), the x-axis is on a linear scale. 1809

K.2.1 INNER PRODUCT RECALL@1 OF EXPERIMENT 1 1812

In this experiment we use Recall@1 with Inner Product similarity. Figure K.5. Similar trends are observed in the Recall@1 of the Inner Product similarity as in the Cosine similarity. The Recall@1 1815 of the Inner Product similarity is also maximized when the alignment is high.



1810 1811

1813

1814

1836 1837 1838 K.3 EMPIRICAL VALIDATION OF THE LINEAR RELATIONSHIP BETWEEN GENERATED DATA PARAMETER e AND $\beta^{\top}\mu$ IN EXPERIMENT 1

1839 As shown in Figure K.6, we observe a positive, linear relationship between e and $\beta^{\top}\mu$ as e is varied. 1840 This confirms the validity of our test data generation method based on e.



Figure K.6: We calculated β from Training Datasets 1, 2, and 3, and then computed $\beta^{\top}\mu$ by adjusting *e* to determine μ in the test data. The x-axis represents *e*, and the y-axis shows the values of $\beta^{\top}\mu$.

4 K.4 ROTATION MATRIX GENERATION PROCESS OF Setup 2

To generate a set of rotation matrices with diverse magnitudes of rotation, we constructed an algorithm that samples k = 300 random matrices, each formed by adding i.i.d. Gaussian noise matrix of varying variance to the identity matrix I. The process ensures the generation of rotation matrices with varying extents of rotation, from slight to more substantial deviations from the identity matrix.

1860 The rotation matrices are generated as follows:

- 1. A matrix is initialized as $I + \epsilon \cdot M$, where M is a i.i.d. standard random Gaussian matrix.
- 2. Using the QR decomposition, we orthogonalize this matrix to ensure it forms a valid rotation matrix.
 - 3. Finally, if the determinant of the resulting matrix is negative, we flip the sign of the first column to maintain a determinant of +1, ensuring it is a valid rotation.

In summary, this method provides a collection of matrices that progressively deviate from I, allowing us to observe and sample rotations of increasing magnitude.

1871 K.5 Additional Discussion of Recall@1 evaluation for Expr 1

1873 The Recall@1 results of Expr 1 setting indicate three phases in Recall@1 outcomes.

1874 The first phase: The learning process fails to improve performance either because the training and 1875 evaluation data are too distant, as predicted by our theory, or because e is too small for the fea-1876 ture extractor to achieve separation. We interpret that either of these factors contributes to the lack 1877 of performance improvement. The second phase: Performance improves as the similarity between 1878 training and evaluation data becomes appropriate, allowing better Recall@1 after training. It is note-1879 worthy that the improvement also increases along with larger e. This indicates not only the increased 1880 e leading to greater distances between evaluation features but also the Recall@1 improvement with e as our theory. The third phase: Effective feature separation has already occurred; thus, even with 1881 sufficiently close training data, the learning process does not enhance Recall@1 performance. 1882

We conclude that the first and third phases represent unsuitable configurations for retrieval tasks, while the second phase provides a dataset which effectively supports training for retrieval tasks and is explainable by our theory.

1886

1849

1850

1851 1852

1855

1861

1862 1863

1864

1865

1866

1867

1870

- L ADDITIONAL SETTINGS FOR EXPERIMENT 3, 4
- 1888 1889

In Table 2, we provide the detailed parameters for the experiments.

Experiment 3 **Experiment 4** Model ResNet18 ResNet50 Multi-class Classification Learning Task **Binary Classification** Norm Softmax (Zhai & Wu, 2019) Loss Mean Squared Error Epoch 20 30 96 (Full GD) **Batch Size** 75 (use 3 classes with 25 samples) Optimizer Adam SGD Learning Rate 0.001 CARS: 0.01 / CUB: 0.001

Table 2: Comparison between Experiment 3 and Experiment 4

¹⁹⁰¹ M Additional Information for Setup 3

For gradient stability and fair evaluation, all classes are truncated to include only 48 images. The batch size is set to 96 for Full gradient descent. Remark that, to align the experimental setup with our theoretical setting, two classifier heads and sign flipped label 1, -1 is used.

N ADDITIONAL RESULTS OF EXPERIMENT 3

The performance of the two classifier heads during training is shown in Figure N.1. The results of the empirical validation without Kendall's W aggregation are presented in subsection N.1 for the model trained with the CARS196 dataset and in subsection N.2 for the model trained with the CUB200 dataset. Consistent with Kendall's W calculations and theoretical analyses, in most cases, we observe that LE, alignment, and $|\beta^{\top}\mu|$ individually rise and fall in similar trends during training. The gray line represents metrics calculated on the entire dataset, while the colored lines denote individual test classes. Since classes were randomly sampled per seed, the same color represents the same class only within a single seed.



(e) 1st seed for the CUB (f) 2nd seed for the CUB (g) 3rd seed for the CUB (h) 4th seed for CUB dataset train dataset train dataset train

Figure N.1: Classification accuracy measured with training data. As two classifier heads were used in the theoretical setup, two accuracy values are plotted for each setting.

36

1907 1908 1909

1890

1892

1894

1895

1897 1898

1899 1900

1903

1904

1905

1906

1930 1931

1932

1933 1934 1935

1938 1939











Figure N.11: The 1st seed of the CUB dataset training. Results were computed using the features of five randomly selected classes from the CUB dataset's test set.



Figure N.12: The 2nd seed of the CUB dataset training. Results were computed using the features of five randomly selected classes from the CAR dataset's test set.



Figure N.13: The 2nd seed of the CUB dataset training. Results were computed using the features of five randomly selected classes from the CUB dataset's test set.



ADDITIONAL RESULTS OF EXPERIMENT 4 Ο

We provide non aggregated data for the experiment 4 in this section. The data is presented in the form of tables. The tables are as follows:

2165	R@1 v. Align	p-value	Recall@1	Recall@2	Recall@4	Recall@8	final Avg. Align
0400	0.3195	0.0013	93.4694	96.6056	97.9707	98.7947	829.4791
2166	0.2386	0.0180	93.7523	96.4949	97.9461	98.7332	858.0315
2167	0.1052	0.3025	94.0844	96.7409	98.1798	98.9792	857.6394
0400	0.2864	0.0043	93.3096	96.4457	98.0445	98.8439	827.9364
2168	$0.2374 {\pm} 0.0942$	0.0000	$93.6539 {\pm} 0.3404$	96.5718±0.1311	98.0353±0.1050	$98.8378 {\pm} 0.1046$	$843.2716 {\pm} 16.8294$
2169							

Table 3: Measurement from CARS196 trained model. R@1 v. Align is Pearson correlation between Recall@1 and final Avg. Align. Recall@k and final average Alignment is measured after training.

0170							
2173	R@1 v. Align	p-value	Recall@1	Recall@2	Recall@4	Recall@8	final Avg. Align
2174	0.2925	0.0032	68.0621	78.7643	86.5294	91.6948	1060.0284
0.185	0.2308	0.0209	68.6867	79.4564	87.6097	92.2687	1088.8525
2175	0.3498	0.0004	67.9946	79.2539	86.7995	92.4038	1050.2296
2176	0.2769	0.0053	67.7583	78.8150	87.0695	92.2181	1090.8598
0177	$0.2875 {\pm} 0.0491$	0.0000	$68.1254{\pm}0.3962$	$79.0724 {\pm} 0.3374$	87.0020 ± 0.4613	$92.1464 {\pm} 0.3111$	$1072.4926 {\pm} 20.4613$
21//							

Table 4: Measurement from CUB200 trained model. R@1 v. Align is Pearson correlation between Recall@1 and final Avg. Align. Recall@k and final average Alignment is measured after training.