052

053

054

000

How Classifiers Extract General Features for Downstream Tasks: An Asymptotic Analysis in Two-Layer Models

Anonymous Authors¹

Abstract

Neural networks learn effective feature representations through intermediate layers, enabling feature transfer without additional training for new tasks. However, the conditions for successful feature transfer remain underexplored. In this paper, we investigate feature transfer in classifier-trained networks, focusing on clustering in unseen distributions. In binary classification, we find that higher similarity between training and unseen distributions improves Cohesion and Separability, while Separability further requires unseen data to be assigned to different training classes. In multiclass classification, our analysis shows that the feature extractor maps input point based on their similarity to training classes, i.e. that unrelated training classes to input have negligible impact on feature extraction. We validate our theoretical findings in synthetic dataset and demonstrate practical applicability utilizing ResNet and variations of CAR, CUB, SOP, ISC, and ImageNet datasets.

1. Introduction

Neural networks have the remarkable ability to adapt to specific tasks, learning representations through penultimate layers. Training these intermediate layers is crucial for neural network generalization (Damian et al., 2022). Also, these layers can extract semantically meaningful and transferable features from new data, enabling feature transfer for new tasks (Yosinski et al., 2014; Kornblith et al., 2019). A wide range of techniques, from open set clustering (Roth et al., 2020; Huang et al., 2024) to vision-language models (Li et al., 2023) and language models (Brown et al., 2020; Kojima et al., 2023), leverage feature transfer for downstream tasks. However, the specific conditions where features can be effectively transferred remain underexplored. Among various applications, classification based visual open-set clustering (Musgrave et al., 2020) serves as a fundamental benchmark for evaluating whether a feature extractor can generalize to unseen data. Typically, this task involves classifier training on one set of classes and then testing it on disjoint classes to assess whether the extracted features form cohesive and separable class-wise clusters on unseen data (Wang et al., 2018; Seidenschwarz et al., 2021; Deng et al., 2022). Given this context, we aim to investigate feature clustering with the following research questions:

Can we capture the presences of feature learning in classification and identify the conditions where features cluster effectively on new distributions?

To address this question, we analyze a two-layer nonlinear network network trained with a single large gradient descent step on a mean-squared classification loss in the *proportional regime* (in section 2). The proportional regime intuitively represents a scenario where the network width and the size of the dataset are of similar scales, aligning with common practices in model scaling (Ba et al., 2022), and they are known to effectively capture the phenomena occurring during the actual training process, as demonstrated in studies such as Mei & Montanari (2020); Moniri et al. (2024). We capture that the dominant part of the trained feature is composed of random initialization and *spikes* (Def. 3.4) associated with the training classes (section 3). Leveraging dominant features, we identify conditions for effective clustering on new distributions (section 4).

In a binary classification setting, we assess the intra-class *cohesion* and inter-class *separability* of trained features in a numerical-analytical manner representing the clustering population risks (Def. 4.3) (Clémençcon, 2011; Papa et al., 2015; Li & Liu, 2021) and goals for clustering performance (Liu et al., 2017). As a result, *Cohesion* increases as the *train-unseen similarity* (in Def. 4.1) grows larger. Meanwhile, for *Separability*, if classes classes are *assigned* (Notes 4.2, E.1) to different training classes, *Separability* increases as the *train-unseen similarity* grows larger; otherwise, it decreases, as illustrated in Figure 1.

 ¹Anonymous Institution, Anonymous City, Anonymous Region,
 Anonymous Country. Correspondence to: Anonymous Author
 <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.



100

104



Figure 1: Mapping data from the input space (left) to the learned feature space (right). Training classes are shown as balls, and unseen classes as dashed lines (a, b, p, n). *Cohesion*: Strong *cohesion* occurs for a, p, n, which have high similarity to the training classes compared to *b*. *Separability of* a, n: a and n, *assigned* to different training class, demonstrate high Separability. *Separability of* a, p: a and p, *assigned* to the same training class, exhibit low Separability.

In the multi-class classification setting, we analyze the *spikes* of features and find that *spikes* map new inputs based on a linear combination of randomly initialized classifier heads' weight with coefficients that represent the similarity of the training classes. Therefore, the more *spikes* aligned with the input data the greater their contribution to feature extraction, enhancing the expressiveness of the features.

In the experiments, we empirically observe train-unseen similarity, cohesion, Separability, and recall@1 under our theoretical assumptions in synthetic datasets. As a result, we confirm that the theoretical interpretation aligns with the actual findings (subsection 5.2). Additionaly, we explore practical metric learning settings and find evidence support-086 ing the validity of our analysis results in a practical setup 087 (subsection 5.4). In most cases, we observe that clustering 088 performance is higher when the unseen classes share the 089 same sementic domain as the training classes. Moreover, 090 adding semantically relevant training classes improves per-091 formance, whereas adding unrelated training classes does 092 not lead to performance improvement. 093

094 Our contributions are summarized into following:

- We analyze the classifier feature, providing insights into how feature extractors operate:
 - Higher train-unseen similarity increases cohesion.
 - Higher train-unseen similarity increases separability between data assigned to different classes but reduces it otherwise.
 - Expressiveness of feature improves with an increased number of *spikes* non-orthogonal to input.
- We generalize the distribution assumption of prior works and present novel proof techniques for classifier analysis.
- The theoretical results are validated through diverse experiments, including synthetic and real-world datasets.

1.1. Related Works

Metric Learning and Open Set Clustering Metric learning is proposed to cluster visually similar unseen classes using classification or triplet loss (Movshovitz-Attias et al., 2017; Zhai & Wu, 2019; Boudiaf et al., 2021). Several recent approaches have focused on increasing the number of classes in the training data to improve clustering. One approach adds virtual classes (Chen et al., 2018; Qian et al., 2020; Gu et al., 2021). Another approach suggested leveraging a larger number of classes induced from Schuhmann et al. (2021) to achieve state-of-the-art performance (An et al., 2023). This aligns with our analysis, which suggests that performance improves as the number of relevant classes in clustering increases.

Neural Collapse (NC) and Unconstrained Layer-Peeled Model (ULPM) Recent studies have introduced the concept of Neural Collapse (Papyan et al., 2020) to explain the emergence of intra-class features and feature-weight alignment in trained neural networks. Several studies propose the ULPM to understand training dynamics of NC treating features and weights as unconstrained free variables (Fang et al., 2021; Zhu et al., 2021; Ji et al., 2022; Tirer & Bruna, 2022). However, ULPM, unlike the two layer network model we use, assumes the free variable features, which limits analyzability about input distribution and, consequently, prevents studying feature transferability.

Feature Learning in Two-Layer Networks Many works (Louart et al., 2017; Goldt et al., 2020; Hu & Lu, 2022) study the Conjugate Kernel (CK), which enables the analysis of the structure of the first layer in two-layer networks. Ba et al. (2022); Moniri et al. (2024); Ba et al. (2023) argue that feature learning aids in reducing the population risk when evaluated on distributions same to the training data. Unlike these studies, we claim that the CK feature learning model not only explains this generalization but also enables the analysis of features from non-identical distributions, facilitating a deeper understanding of feature transfer.

Additional related works are provided in Appendix A.

2. Problem Statement

Notations Let $\|\cdot\|$ be L^2 or the operator norm. Let \odot be the Hadamard product. Let $A^{\circ k}$ be the Hadamard power. Let C, c > 0 and $\kappa \in \mathbb{R}$ be constants that may change from line to line. Define $[d] \triangleq \{1, 2, \cdots, d\}$. For o, O, Θ notations we follow Moniri et al. (2024)

Training Data We define data for one vs. one classification with $\#_{cls}$ classes. The number of problem $\#_P \triangleq \frac{\#_{cls}(\#_{cls}-1)}{2}$. Let $\#_{cls}$ be the number of training classes, and let $\mathscr{C}_1, \dots, \mathscr{C}_{\#_{cls}}$ represent the class-conditional distri-

110 butions of the training data. Define the training dataset 111 as $\mathcal{D} = (X, Y)$, where $X \in \mathbb{R}^{\mathbf{n} \times \mathbf{d}}$, $Y \in [\#_{cls}]^{\mathbf{n}}$, 112 $X = (\{x \sim \mathscr{C}_1\} \times m \cup \cdots \cup \{x \sim \mathscr{C}_{\#_{cls}}\} \times m)$, where 113 $\#_{cls}m = \mathbf{n}$ and m is the number of instances per class. Let 114 $\tilde{\mathcal{D}} = (\tilde{X}, \tilde{Y})$ an i.i.d. copy of \mathcal{D} . 115

Network Structure We consider two-layer networks. The 116 initial weight of the first layer, $W_0 \in \mathbb{R}^{\mathbf{d} \times \mathbf{N}}$, is initialized as 117 $W_0[i] \sim Unif(\mathbb{S}^{\mathbf{d}-1})$ for $i \in [\mathbf{d}]$. We denote W obtained 118 119 via a single step of gradient descent. The initial weights of 120 the second layer, $a_{ij} \in \mathbb{R}^{N}$ for $i, j \in [\#_{cls}]$ s.t. i < j, are initialized as $a_{ij} \sim \mathbf{N}(0, \frac{1}{\mathbf{N}}I)$. We define the initialized fea-121 122 ture as $F_0(x) \triangleq \sigma(W_0^{\top} x)$ and the one-step trained feature 123 as $F(x) \triangleq \sigma(W^{\top}x)$. The network output is defined as the 124 following $\#_P$ -dimensional vector: $(F(x)^{\top}a_{ij})|_{ij}$. 125

Proportional Regime We consider the two-layer neural networks in the proportional regime. n, d, and N are sample size, data and feature dimension, respectively. We perform our analysis under d/n, $N/n \rightarrow c$ as n, d, $N \rightarrow \infty$.

126

127

128

129

130 131

132

133

134

135

136

137

138 139

140 141 142

143

144

145

146

147

148

149

Optimization Problem Denote the set of all network parameters as $\theta = \{W, a_{12}, \cdots, a_{\#_P-1, \#_P}\}$. Let X_{ij} be a matrix in $\mathbb{R}^{2m \times d}$, where the first m rows contain samples $x \sim c_i$ and the last m rows contain samples $x \sim c_j$. Let $y \triangleq [1, 1, \ldots, 1, -1, \ldots, -1]^\top \in \mathbb{R}^{2m}$ be a vector consisting of m ones followed by m negative ones. To classify the given data, we use the Mean Squared Error,

$$L(x, y; \theta) = \frac{1}{2n} \sum_{i < j}^{c} \|y - \sigma(X_{ij}W)a_{ij}\|^2.$$
(1)

The weight update formula for the first layer is given by $W = W_0 + G$, where $G \triangleq -\frac{\partial L}{\partial w} = \sum_{i < j} G_{ij}$, s.t.

$$G_{ij} = -\frac{1}{n} \left[X_{ij}^T [(\sigma(X_{ij}W)a_{ij} - y)a_{ij}^T \odot \sigma'(X_{ij}W)] \right].$$
(2)

Now, we introduce the assumptions for theoretical analysis.

150 Assumption 2.1 (Activation Function). Let $\sigma(x)$ be an 151 element-wise activation s.t. $\sigma, \sigma', \sigma''$ is bounded by λ_{σ} 152 almost surely. It admits a Hermite decomposition i.e. 153 $\sigma(z) = \sum_{k=0}^{\infty} c_k H_k(z)$, where $c_k = \frac{1}{k!} \mathbb{E}[\sigma(z) H_k(z)]$ for 154 standard gaussian z. We assume $c_0 = 0, c_1 > 0$ and 155 $c_k^2 k! \leq C k^{-3/2-w}$, for constants C, w > 0. For example, 156 Shifted ReLU $\max(x, 0) - \frac{1}{\sqrt{2\pi}}$ satisfies this condition.

Assumption 2.2 (Training Data). Let the class-conditional training data distributions \mathscr{C}_i be non-centered Sub-Gaussians (Vershynin, 2018; Cao et al., 2021; Cole & Lu, 2024). This distribution family is suitable for classification, including distributions with limited support that are separable. It is an extension of the Gaussian assumption of Ba et al. (2022).

3. Feature Decomposition

This section analyzes the learning dynamics during a single gradient descent step. First, we demonstrate that the gradient with respect to the W_0 exhibits an almost Rank- $\#_P$ property within the proportional regime. Subsequently, we prove that the learned features can be predominantly expressed as Rank- $\#_P$ components, establishing the dominant components for subsequent analyses.

Gradient Decomposition We decompose the gradient (equation 2) using Hermite decomposition, which allows us to extract the essential rank-one matrix structure for each *ij*-th classification problem. Note that $\sigma' = c_1 + \sigma'_{\perp}$.

$$G_{ij} = \frac{c_1}{n} X_{ij}^T y a_{ij}^T + \frac{1}{n} X_{ij}^T y a_{ij}^T \odot \sigma'_{\perp}(X_{ij} W_0) - \frac{1}{n} X_{ij}^T \sigma(X_{ij} W_0)(a_{ij} a_{ij}^T) \odot \sigma'(X_{ij} W_0)$$
(3)
$$\triangleq \mathbb{A}_{ij} + \mathbb{B}_{ij} + \mathbb{C}_{ij}.$$

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

Theorem 3.1 (Approximation of Gradient). Under the assumptions in section 2, and when **n** satisfies $\frac{1}{2} > \kappa \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}}$, the following holds w.p. $1 - C(\mathbf{n}e^{-c\log^2 \mathbf{n}} + e^{-c\mathbf{n}})$:

$$\|G - \sum_{i < j} \mathbb{A}_{ij}\| \le \kappa \frac{\log^2 \mathbf{n}}{\mathbf{n}}.$$
 (4)

Feature Decomposition Now we utilize $\sum_{i < j} \mathbb{A}_{ij}$ to decompose the feature extractor. We decompose the one-step trained feature function $F(x) = \sigma((W_0 + G)^{\top} x)$, which serves as a key step in deriving our main analysis. For the proof, please refer to Appendix J.

Definition 3.2 (Data-Label Covariance). Data-Label Covariance for X_{ij} is defined as $\beta_{ij} = \frac{1}{n} X_{ij}^{\top} y \in \mathbb{R}^{\mathbf{d}}$.

Theorem 3.3 (Decomposition of Trained Features). Under the assumptions in section 2, let $F_0 = \sigma(\tilde{X}W_0)$, $L \triangleq \log n$, $F_0^L = \sum_{k=1}^L c_k H_k(\tilde{X}W_0)$, and $spike_L = \sum_{k=1}^L c_k^k c_k (\tilde{X} \sum_{i < j} \beta_{ij} a_{ij}^T)^{ok}$. With probability 1 - o(1),

$$F = F_0^L + spike_L + \Delta. \tag{5}$$

Moreover, $\|\text{spike}_L\|$ is greater than \sqrt{n} , $\|F_0^L\| = \Theta(\sqrt{n})$, and $\|\Delta\| = o(\sqrt{n})$.

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



Figure 2: Numerical Observation of Cohesion and Separability. Plot of *Cohesion* and Heatmap of *Separability* calculated by adjusting $\beta^{\top}\mu_1$ and $\beta^{\top}\mu_2$.

Definition 3.4 (Dominant Feature
$$F_L = F_0^L + \text{spike}_L$$
).

$$F_L(x) \triangleq \sum_{k=1}^{L} c_k [H_k(\tilde{X}W_0) + c_1^k (\sum_{i < j} (\beta_{ij}^\top x) a_{ij}^T)^{\circ k}].$$
(6)

Using the feature decomposition conducted so far, the next section analyzes clustering risk and explores the conditions for effective clustering of unseen data.

4. Feature Analysis

175

176

177 178 179

180 181 182

183

184

185

186 187 188

189

190

213

214

215

216

217

4.1. Clustering Risk Analysis in binary classification

¹⁹¹ In this section, we analyze clustering risks. We show **train** ¹⁹² (β)-**unseen** (μ) similarity governs the the clustering pop-¹⁹³ ulation risk i.e. *Cohesion* and *Separability* of F_L from ¹⁹⁴ Definition 3.4 under condition 4.4. We derive *cohesion* ¹⁹⁵ and *separability* of F_L for two "unseen" class-conditional ¹⁹⁶ distributions.

Definition 4.1 (*Train-Unseen Similarity*). Given Train Data-Label Covariance β in Definition 3.2 and mean of Unseen distribution μ , *Train-Unseen Similarity* is defined as $\beta^{\top}\mu$.

Note 4.2 (Explanation of *assignment* and $\beta^{\top}\mu$). β_{ij} represents the normal vector of the linear decision boundary, i.e. the direction determining class *i* vs. *j* based on the sign of its inner product with data. Therefore, the sign of $\beta^{\top}\mu$ indicates the class *assignment* of unseen data with μ .

206 Definition 4.3 (*Cohesion* and *Separability*). We define the
207 clustering risks based on similarity between feature vectors
208 using inner products.

209 Cohesion measures the expected similarity between i.i.d. 210 features of the same class over network parameters θ and 211 data $x, x' \sim c_1$, i.e. 212

$$\mathbb{E}_{\theta}[\mathbb{E}_{x \sim c_1} F(x)^T \mathbb{E}_{x' \sim c_1} F(x')]$$

Separability measures the expected dissimilarity between independent features of different classes over θ , $x \sim c_1$ and $x' \sim c_2$ i.e.

$$-\mathbb{E}_{\theta}[\mathbb{E}_{x \sim c_1} F(x)^T \mathbb{E}_{x' \sim c_2} F(x')].$$

Condition 4.4. We fix n, d, N large enough. Under assumptions 2.1, 2.2, let $c_i = \mathcal{N}(\mu_i, I_d)$ for $i \in [2]$ be the class conditional distributions. Define $\rho_{k,k'}^{(1)} > 0, \rho_{k,k'}^{(2)}(\cos(\mu_1, \mu_2)), \rho_{k,k',r}^{(3)} > 0, \rho_{k,k',r,r'}^{(4)} > 0$ as functions of N, d. Note that $\rho_{k,k'}^{(2)}$ increases as $\cos(\mu_1, \mu_2)$ grows. Exact definitions are in Def. K.1. The Shifted ReLU, as stated in Assumption 2.1, is used as the activation.

Proposition 4.5 (Cohesion). Following condition 4.4, the Cohesion of F_L for c_i , $i \in [2]$ is given by:

$$\sum_{k=1,k'=1}^{L} c_{k} c_{k'} \left\{ \begin{array}{l} \rho_{k,k'}^{(1)} \|\mu\|^{k+k'} \\ + 2 \sum_{r'=0}^{k'} \rho_{k,k',r'}^{(3)} |\mu^{T}\beta|^{k'-r'} \|\beta\|^{r'} \|\mu\|^{k} \\ + \sum_{r,r'=(0,0)}^{(k,k')} \rho_{k,k',r,r'}^{(4)} |\mu^{T}\beta|^{k+k'-r-r'} \|\beta\|^{r+r'}. \end{array} \right.$$

Proposition 4.6 (Separability). Following condition 4.4, the Separability of F_L for c_1, c_2 is given by:

$$-\sum_{k=1,k'=1}^{L} c_{k}c_{k'} \left[\begin{array}{c} \rho_{k,k'}^{(2)}(\cos(\mu_{1},\mu_{2})) \|\mu_{1}\|^{k} \|\mu_{2}\|^{k'} \\ +\sum_{r=0}^{k} \rho_{k,k',r}^{(3)} |\mu_{1}^{T}\beta|^{k-r} \|\beta\|^{r'} \|\mu_{2}\|^{k'} \\ +\sum_{r'=0}^{k'} \rho_{k,k',r'}^{(3)} |\mu_{2}^{T}\beta|^{k'-r'} \|\beta\|^{r'} \|\mu_{1}\|^{k} \\ +\sum_{r,r'=(0,0)}^{(k,k',r,r')} \rho_{k,k',r,r'}^{(4)} (\mu_{1}^{T}\beta)^{k-r} (\mu_{2}^{T}\beta)^{k'-r'} \|\beta\|^{r+r'} \\ \end{array} \right]$$

$$\tag{8}$$

The proofs of Propositions 4.5 and 4.6 are provided in Appendix K. We numerically analyze the results of propositions 4.5 and 4.6 to investigate Cohesion and Separability further. For this numerical observations, we set $\|\mu_1\| = \|\mu_2\| = \|\beta\| = 1, \ \mu_1 = -\mu_2 \in \mathbb{R}^{320000}$ and $L = \log_{10} \mathbf{n}$. We calculate equation 7 and equation 8 by adjusting $\mu_1^T \beta$ and $\mu_2^T \beta$, as shown in Figure 2, which demonstrates the Cohesion and Separability of F_L . Cohesion increases when the $|\mu^T\beta|$ increases. Separability increases when $\mu_1^T \beta$ and $\mu_2^T \beta$ grow with opposite signs and decreases when they grow with the same sign. Moreover, we observe that this phenomenon is governed by the last term of equation 7, 8 (related to $\rho^{(4)}$), as shown by separately computing this term and the others numerically in Appendix B. Additionally, under the theoretical setup, we observe that our hypothesis tends to hold over a wider range as n increases (please refer to Appendix B).

The analytical results in equation 7 and equation 7 can be explained as follows. With $\rho^{(4)} > 0$, the last term inside the bracket of *Cohesion* in equation 7 increases in value as *Train-Unseen Similarity* grows. The last term of *Separability* is influenced by $(\mu_1^T \beta)^{k-r} (\mu_2^T \beta)^{k'-r'}$. Provided that k - r and k' - r' are odd, this term implies that if the *Train-Unseen Similarities* have opposite signs and increase, then this term improves; otherwise, if the signs are the same and increase, *Separability* decreases. According to the analysis



Figure 3: As shown in equation 6, after one step of training with spike $\beta_1, \beta_2, \beta_3, \beta_4$, the inner product between input x_i and β_i acts as the coefficient in the linear combination of a_i , forming the *spikes* structure of the feature.

in Appendix H, the first coefficient c_1 of Shifted ReLU is a large positive value, and subsequent Hermite coefficients approach zero while oscillating around it. Thus, we hypothesize that the positive part is likely to dominate $\sum c_k c_{k'}$, but further work is needed to confirm this.

4.2. Spike Component Analysis

232

233

234

235

236 237

238

239

240

241

242 243

244

245

247

248

249

250

251

252

253

254

255

256

257

258

259

261

262

263

264

265

266

267

268

269

270

271

272

273

274



Figure 4: When trained along the directions β_2 and β_3 , we observe significant changes in the feature space distance as x_1 and x_2 vary, compared to β_1, β_4 .

In this section, based on the previous feature decomposition and extend it to examine the impact of a multi-class classifier's spike structure on unseen data clustering. We examine the spike structure in $F_L = F_0^L + \text{spike}_L$ and its influence on feature mapping. This examination allows us to explore the impact of the training data's structure β on the feature generation of unseen data. The spike structure inside the Hadamard power involves the linear combination coefficient $\beta_{ij}^{\top}x$ and the random initialized classifier head a_{ij} (equation 3.4). Thus, the feature extraction is closely linked to the inner product between β_{ij} and the input point x. If the direction of x is not orthogonal to β_{ij} , then spike of β_{ij} involve feature extraction.



Figure 5: Comparison of log average slope between Theory and Two-layer Networks. \blacksquare Midpoint (β_1) \blacksquare Interpolation (β_2) = Extrapolation (β_3) = Orthogonal (β_4) . The intersection implies learning intersecting β .

Conversely, when x is orthogonal to β_{ij} , the impact of spike β_{ij} is eliminated. To validate this, we define following four spikes, given test input $x_1, x_2 \in \mathbb{S}^{d-1}(\sqrt{\mathbf{d}})$, $\beta_1 = \frac{x_1+x_2}{2}, \beta_2 = \frac{x_1+3x_2}{4}, \beta_3 = \frac{-x_1+5x_2}{4}$ and β_4 , a random vector orthogonal to x_1, x_2 . Then, the magnitudes are adjusted to $\sqrt{\mathbf{d}}$. By definition, β_1, β_4 cannot contribute to feature extraction because they are Midpoint or Orthogonal, while β_2 and β_3 can distinguish the two inputs. For illustration see Figure 3.

Now, we demonstrate this explanation using the approximated features F_L and the two layer neural network F with the four disjoint sub-classification problem 1 defined as follows: We generated four classification problems by creating Gaussian training data with means β_i and $-\beta_i$, and a covariance of 0.1I for $\mathbf{n}, \mathbf{d}, \mathbf{N} = 2^{11}$, enabling the networks to learn β_i as their *spike*. F is trained by this data and F_L is calculated by its definition. We observed the feature distance between $F(x_1), F(x_2)$ and between $F_L(x_1), F_L(x_2)$ for $\binom{4}{k}$ combinations of β_i in this problem by varying the angle between x_1, x_2 . Please refer to Figure 4 and 21 for results. It can be observed that the feature from β_1 and β_4 hardly captures variations in the angle of test input x_1, x_2 within the data space. In contrast, the feature from β_2 and β_3 is highly sensitive to such variations, suggesting that it effectively preserves the structural changes in the input data. Both $F_L(x_1)$ and $F(x_1)$ exhibit the same trends, which supports the validity of our feature approximation. To aggregate these combinatorial results, we measure the log of the average slope, which indicates that features with sensitive changes tend to have larger values, as shown in Figure 5.

As a result in Figure 5, we observe that when multiple β s are used in training, features are more sensitive to changes

¹Instead of studying all combinations for 8 classes classification, we simplify the task by grouping four pairs, performing only four combinations of classifications.

Submission and Formatting Instructions for ICML 2025



Figure 6: Examples of training datasets (Data 1, 2, 3) and evaluation data Eval 1, 2.

in distance within the data space. Meanwhile, the Midpoint β_1 and Orthogonal spike β_4 seem ineffective for feature extraction, even when learned alongside other spikes. Experiments show that learning representations with unrelated classes limits expressiveness, while related classes enhance the model's ability to capture fine-grained features of unseen data. This trend is consistently observed in real-world datasets in Expr V, VI at subsection 5.4. Additionally, to clarify the effect of the spikes, we compute F_0^L and spike_L separately as shown in Figure 22. The results show that the spike_L created by β_1 and β_4 embeds x_1 and x_2 as the same feature. Therefore, it confirms that the distinction between x_1 and x_2 created by the model trained with β_1 and β_4 is due to the random feature F_0^L .

5. Experiments

284 285 286

287

288

289

290

291

292

293

294

295

296

297

298

299

300 301

302 303

304

305

306

317

318

Remark 5.1. *recall*@ $l \triangleq \mathbb{E}_{x_i, y_i} \mathbf{1}_{y_i = \hat{y}_{i, 1-NN}}$. $\hat{y}_{i, 1-NN}$ is class of the closest feature to x_i . This is a feasible measure for evaluating whether new classes form clusters.

307 In this section, we conduct seven experimental setups to vali-308 date our theoretical results. First, in Experiments I, II and III, 309 we utilize a synthetic dataset to confirm that, as discussed 310 in subsection 4.1, Cohesion, Separability are determined by 311 the Train-unseen similarity. Second, to demonstrate how our 312 theoretical explanations can provide intuition in practical 313 settings, we conduct Experiments IV, V, VI, and VII. For 314 this purpose, we analyze the open-set clustering problem 315 using fine-grained real image datasets. 316

5.1. Setup for Theory Vaildation: Expr. I, II, III

319 We use three types of different non-centered Sub-Gaussian 320 distributions as training datasets that are symmetric about 321 the origin. For the evaluation, we introduce two distribution 322 i.e. Eval 1, Eval 2 with translation parameter e and rotation 323 parameter $R \in \mathcal{R} \subseteq SO(n)$ to control the *train-unseen* 324 similarity $\beta^{\top}\mu$. e.g. as e increases from 0 towards 1, $\beta^{\top}\mu$ 325 increases, and as R approaches the identity matrix $I, \beta^{\top} \mu$ increases. For illustration of the data, see Figure 6. For 327 detail, refer to subsection D.1. We follow the condition 328 described in section 2 and subsection 4.1. 329

Now we explain Expr. I, II, III. For each experiment, we utilize all datasets 1, 2, 3, with distinct Eval data usage. Expr. I uses two Eval 1 data with translation parameter $e_1 \in [-0.9, 0.9]$ and $e_2 = -e_1$, so they are *assigned* to opposite training classes (say pos-neg). Experiments II and III are based on two Eval 2 data distributions, each parameterized by a small-angle random rotation matrix $R \in \mathcal{R}$. In Experiment II, considering the case where the datasets are *assigned* to opposite classes, the first distribution uses R and the second distribution is origin symmetry of the first distribution. In Experiment III, considering the situation where the datasets are *assigned* to the same class (say pospos), the first distribution uses R and the second uses R^{\top} to slightly rotate given means.

5.2. Results of Theory Vaildation: Expr. I, II, III

In this experiment, we examine the relationships between the train-unseen similarity (i.e. $\beta^{\top}\mu$), Cohesion, Separability that we discussed in subsection 4.1 and Recall@1 to evaluate performance using practical measures. All test data are generated symmetrically, so for simplicity in visualization, we report the measurement for a single class. For Expr I, we present a summary of the results in Figure 8. We observe that for large values of $|\beta^{\top}\mu|$, strong *Cohesion* and Separability occur across all datasets. For Expr II and III, in accordance with the Separability structure observed in subsection 4.1, when the signs of $\beta^{\top}\mu_1, \beta^{\top}\mu_2$ are opposite (Expr II), we observed an increase in Separability, whereas in the other case (Expr III), we observed a decrease Figure 7. For recall@1, we observed a similar trend as Separability. These results correspond to our theoretical findings. For individual graphs, refer to Appendix D.

5.3. Setup for Practical Vaildation: Expr. IV, V, VI, VII

We designed experiments to examine whether these insights are also applicable to clustering performance in image datasets and practical neural networks. In these scenarios, we utilize *train-unseen similarity* to conceptualize semantic similarity between training and unseen classes (Expr. IV). The number of non-orthogonal *spikes* is interpretable as the number of semantically similar or dissimilar training classes

Submission and Formatting Instructions for ICML 2025



Figure 7: Data 1 evaluated in the Eval 2 setup. Upper row: In Expr II, all metrics increase as $|\beta^{\top}\mu|$ increases. Lower row: In Expr III, where two test classes are *assigned* to a single train class, recall@1 and Separability tend to decrease as $|\beta^{\top}\mu|$ increases. This aligns with our predictions. The red line — represents the values after one step training. Tje blue line — represents the values from initialization.



350

351

352

353 354 355

356

358

359

360 361

362

363

365

366

367

368

369

370

371 372

Figure 8: Summary of Expr. I. D_i denotes Data *i* and C, S denote *Cohesion* and *Separability*. Dark and large points indicate low $|\beta^{\top}\mu|$ values, while the opposite indicates high values. All measurements increase with respect to $|\beta^{\top}\mu|$. We scaled using the absolute value at the 85th percentile.

(Expr. V, VI). Additionally, we validate whether removing
the duplicatively *assigned* unseen classes improve clustering risk compared to random removal, as suggested by the
results of *Separability* (Expr. VII).

For this investigation, we used the benchmark datasets
CAR(Vehicle) (Krause et al., 2013), CUB(Bird) (Wah et al.,
2011), SOP(Product) (Song et al., 2015), and ISC (Clothing) (Liu et al., 2016), referred to as *Domain*. Additionally,
we utilized ImageNet subsets corresponding to the domains
Vehicle, Bird, Product, and Clothing, denoted as I(V), I(B),
I(P), and I(C), referred to as *sub In1k* for extra classes. Also,

we performed experiments on the whole classes ImageNet by sampling 100 instances per class (say *subsampled whole In1k*). Details are in Appendix N. The objective function and most experimental configurations followed the approach outlined in Zhai & Wu (2019), which is a seminal baseline. We use ResNet18 and ResNet50 (He et al., 2015). In addition to the randomly initialized networks in the main text, we conducted experiments with pre-trained networks common in feature learning, and results are included in Appendix E. The two setups exhibited similar trends.

5.4. Results of Practical Vaildation: Expr. IV, V, VI, VII



Figure 9: Expr. IV, *recall@1* measurements. Most cases show the highest performance when the domain of the Train and Test corresponds.

For **Expr. IV**, we trained with each *Domain* dataset (CAR, CUB, SOP, and ISC train datasets) and *Domain+sub In1k* dataset (CAR+I(V), CUB+I(B), SOP+I(P), and ISC+I(C)), and then measured how each model well cluster on all of the test datasets (CAR, CUB, SOP, ISC test datasets). As shown in Figure 9, we verify whether clustering the test dataset related to the train classes is more effective than clustering unrelated data, analogous to result in subsection 4.1.



395

396

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

425

426

427

428

Figure 10: Expr V in ResNet50(init). The pink , red , and blue bars represent *Domain*, *Domain+sub In1k*, *Domain+subsampled whole In1k*, respectively.

In **Expr. V**, we measured the clustering performance for corresponding test datasets after learning the *Domain*, *Domain+sub In1k*, and *Domain+subsampled whole In1k*. We find that adding classes from the entire ImageNet dataset during training, rather than including only related classes, does not significantly improve clustering (Figure 10, 32).



Figure 11: Expr VI, Recall@1 values for the CAR, CUB, SOP, and ISC datasets are shown with dashed lines - - for ResNet18 and solid lines — for ResNet50.

429 In Expr. VI, experiments are conducted by dividing the 430 Domain datasets into four steps to observe the impact of 431 increasing the number of related classes on recall@1 perfor-432 mance (Figure 11). From Step 0 to Step 3, 25%, 50%, 75%, 433 and 100% of the Domain dataset classes are sequentially 434 added for training. The added classes are randomly selected, 435 and each experiment is repeated three times. For the number 436 of classes, refer to Table 6. Furthermore, we observed that 437 some results of Expr. V align with those of Expr. VI, as 438 discussed in detail in subsection E.1. 439

For **Expr. VII**, in evaluation, removing duplicatively *assigned* of unseen classes resulted in a $1.73 \pm 2.87\%$ improvement in recall@1 compared to random removal of same amount of unseen classes, with max improve: 13.65%, min decrease: -3.28%, a success rate: 79% and $p = 9.40 \times 10^{-7}$. This suggest that duplicate *assignments* hinder clustering, which aligns with our theory. Details are in subsection E.2.

6. Conclusion

In this study, we explored the feature learning dynamics of a two-layer classifier in the proportional regime to uncover the mechanisms underlying feature transferability. Specifically, we analyzed the conditions where the learned features of unseen classes form cohesive and separable cluster. Our theoretical analysis extends the Conjugate Kernel framework to classification tasks. As a result, our numerical-analytical theory demonstrated that feature *cohesion* increases with greater similarity between training and unseen data, while feature *separability* is influenced not only by similarity but also by avoiding duplicate class assignments in binary classification. Additionally, we showed that only when the spikes are non-orthogonal to the input, do they get involved in feature extraction. In addition to validation on synthetic datasets, we observed that our theory offers valuable insights even when applied to real-world datasets.

Our empirical findings suggest that clustering performance improves when the test data share the same semantic domain as the training data. Furthermore, adding semantically relevant classes to the training set leads to performance gains, whereas introducing unrelated classes has little effect. Contrary to existing research that focuses on performance improvement through large-scale learning on broad domains (Brown et al., 2020; An et al., 2023), our study provides evidence that only certain relevant knowledge, closely related to the domain, influences feature transfer. This evidence mirrors classical problems in the field of artificial intelligence, such as the frame problem and the installation problem. Specifically, AI agents do not require all available knowledge to solve a given problem; only specific, detailed knowledge is necessary. Dennett (1984) states about this as follows: "People in AI ... take the shortcut of installing all that an agent has to know to solve a problem. This may, of course, be a dangerous shortcut." We hope that our study may remind the AI community of the longstanding principle that it may not be the scale of the data that matters. We have also discussed the limitations and future research directions related to the Hermite expansion approximation and general results for cohesion and separability in Appendix F.

440 Impact Statement

441 This paper presents work aimed at advancing the field of 442 Machine Learning. In this research, we analyze the poten-443 tial for clustering performance improvement through the 444 classification training of a large number of highly granular 445 classes. Such an approach may lead to a reduction in the 446 level of personal data masking required for fine-grained data 447 differentiation, which could trigger new ethical discussions 448 regarding privacy protection. Additionally, to effectively 449 implement this approach, there may be a tendency to collect 450 more data, which can have significant implications for the 451 scale and scope of data collection, as well as data manage-452 ment practices. 453

References

454

455

456

457

458

459

460

- An, X., Deng, J., Yang, K., Li, J., Feng, Z., Guo, J., Yang, J., and Liu, T. Unicom: Universal and compact representation learning for image retrieval, 2023. URL https://arxiv.org/abs/2304.05884.
- 461 Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., Wu, D.,
 and Yang, G. High-dimensional asymptotics of feature
 learning: How one gradient step improves the representation, 2022. URL https://arxiv.org/abs/2205.
 01445.
- Ba, J., Erdogdu, M. A., Suzuki, T., Wang, Z., and Wu, D.
 Learning in the presence of low-dimensional structure: A
 spiked random matrix perspective. In Oh, A., Naumann,
 T., Globerson, A., Saenko, K., Hardt, M., and Levine, S.
 (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 17420–17449. Curran Associates,
 Inc., 2023.
- 474
 475
 476
 476
 476
 477
 478
 479
 478
 479
 478
 479
 478
 479
 479
 478
 479
 479
 479
 479
 470
 470
 470
 471
 472
 473
 474
 474
 474
 475
 475
 476
 477
 478
 479
 479
 479
 479
 470
 470
 470
 470
 471
 471
 472
 473
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
 474
- 480 Bellet, A. and Habrard, A. Robustness and generalization for 481 metric learning. *Neurocomputing*, 151:259–267, March 482 2015. ISSN 0925-2312. doi: 10.1016/j.neucom.2014.
 483 09.044. URL http://dx.doi.org/10.1016/j. 484 neucom.2014.09.044.
- Bienstman, P. 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.
- Boudiaf, M., Rony, J., Ziko, I. M., Granger, E., Pedersoli,
 M., Piantanida, P., and Ayed, I. B. A unifying mutual information view of metric learning: cross-entropy vs.

pairwise losses, 2021. URL https://arxiv.org/ abs/2003.08983.

- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners, 2020. URL https:// arxiv.org/abs/2005.14165.
- Cao, Y., Gu, Q., and Belkin, M. Risk bounds for over-parameterized maximum margin classification on sub-gaussian mixtures. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (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.
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. Emerging properties in self-supervised vision transformers, 2021. URL https: //arxiv.org/abs/2104.14294.
- Chang, Y., Hu, C., and Turk, M. Manifold of facial expression. In 2003 IEEE International SOI Conference. Proceedings (Cat. No.03CH37443), pp. 28–35, 2003. doi: 10.1109/AMFG.2003.1240820.
- Chen, B., Deng, W., and Shen, H. Virtual class enhanced discriminative embedding learning, 2018. URL https://arxiv.org/abs/1811.12611.
- Chopra, S., Hadsell, R., and LeCun, Y. Learning a similarity metric discriminatively, with application to face verification. In 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05), volume 1, pp. 539–546 vol. 1, 2005. doi: 10.1109/CVPR.2005.202.
- Clémençcon, S. On u-processes and clustering performance. In Shawe-Taylor, J., Zemel, R., Bartlett, P., Pereira, F., and Weinberger, K. (eds.), *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips. cc/paper_files/paper/2011/file/ ald0c6e83f027327d8461063f4ac58a6-Paper. pdf.
- Cole, F. and Lu, Y. Score-based generative models break the curse of dimensionality in learning a fam-

- 495 ily of sub-gaussian distributions. In *The Twelfth In-*496 *ternational Conference on Learning Representations*,
 497 2024. URL https://openreview.net/forum?
 498 id=wG12xUSqrI.
- Damian, A., Lee, J. D., and Soltanolkotabi, M. Neural networks can learn representations with gradient descent, 2022. URL https://arxiv.org/abs/2206.15144.
- Deng, J., Guo, J., Yang, J., Xue, N., Kotsia, I., and Zafeiriou,
 S. Arcface: Additive angular margin loss for deep face recognition. *IEEE Transactions on Pattern Anal ysis 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.
- 512 Dennett, D. Cognitive wheels: The frame problem of ai. 01
 513 1984.
- 515 El-Nouby, A., Neverova, N., Laptev, I., and Jégou, H. Train516 ing vision transformers for image retrieval, 2021. URL
 517 https://arxiv.org/abs/2102.05644.
- Ermolov, A., Mirvakhabova, L., Khrulkov, V., Sebe, N., and
 Oseledets, I. Hyperbolic vision transformers: Combining
 improvements in metric learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7409–7419, 2022.
- Fan, Z. and Wang, Z. Spectra of the conjugate kernel and neural tangent kernel for linear-width neural networks, 2020. URL https://arxiv.org/abs/2005.11879.
- Fang, C., He, H., Long, Q., and Su, W. J. Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training. *Proceedings of the National Academy of Sciences*, 118(43), October 2021. ISSN 1091-6490. doi: 10.1073/pnas.2103091118. URL http: 534 //dx.doi.org/10.1073/pnas.2103091118.
- Galanti, T., György, A., and Hutter, M. On the role of
 neural collapse in transfer learning, 2022. URL https:
 //arxiv.org/abs/2112.15121.

547

548

- Goldt, S., Loureiro, B., Reeves, G., Krzakala, F., Mézard,
 M., and Zdeborová, L. The gaussian equivalence of generative models for learning with shallow neural networks. *Proceedings of the 2nd Mathematical and Scientific Machine Learning Conference, PMLR 145:426-471 (2021)*,
 06 2020. URL https://arxiv.org/pdf/2006.
 14709.pdf.
 - Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. *Deep learning*, volume 1. MIT Press, 2016.

- Gu, G., Ko, B., and Kim, H.-G. Proxy synthesis: Learning with synthetic classes for deep metric learning, 2021. URL https://arxiv.org/abs/2103.15454.
- Han, X. Y., Papyan, V., and Donoho, D. L. Neural collapse under mse loss: Proximity to and dynamics on the central path, 2022. URL https://arxiv.org/abs/ 2106.02073.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition, 2015. URL https://arxiv.org/abs/1512.03385.
- Hu, H. and Lu, Y. M. Universality laws for high-dimensional learning with random features, 2022. URL https://arxiv.org/abs/2009.07669.
- Huai, M., Xue, H., Miao, C., Yao, L., Su, L., Chen, C., and Zhang, A. Deep metric learning: The generalization analysis and an adaptive algorithm. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 2535–2541. International Joint Conferences on Artificial Intelligence Organization, 7 2019. doi: 10.24963/ijcai.2019/352. URL https: //doi.org/10.24963/ijcai.2019/352.
- Huang, H., Nie, Z., Wang, Z., and Shang, Z. Cross-modal and uni-modal soft-label alignment for image-text retrieval, 03 2024. URL https://arxiv.org/pdf/ 2403.05261.pdf.
- Hui, L., Belkin, M., and Nakkiran, P. Limitations of neural collapse for understanding generalization in deep learning, 2022. URL https://arxiv.org/abs/2202. 08384.
- Isserlis, L. 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.
- Ji, W., Lu, Y., Zhang, Y., Deng, Z., and Su, W. J. An unconstrained layer-peeled perspective on neural collapse, 2022. URL https://arxiv.org/abs/2110.02796.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners, 2023. URL https://arxiv.org/abs/2205.11916.
- Kornblith, S., Shlens, J., and Le, Q. V. Do better imagenet models transfer better?, 2019. URL https://arxiv. org/abs/1805.08974.
- Krause, J., Stark, M., Deng, J., and Fei-Fei, L. 3d object representations for fine-grained categorization. In Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops, June 2013.

- Kuchibhotla, A. K. and Chakrabortty, A. Moving beyond sub-gaussianity in high-dimensional statistics: applications in covariance estimation and linear regression. *Information and Inference: A Journal of the IMA*, 11(4):1389–1456, June 2022. ISSN 2049-8772. doi: 10.1093/imaiai/iaac012. URL http://dx.doi.org/ 10.1093/imaiai/iaac012.
- Lee, K.-C., Ho, J., Yang, M.-H., and Kriegman, D. Videobased face recognition using probabilistic appearance manifolds. In 2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings., volume 1, pp. I–I, 2003. doi: 10.1109/CVPR. 2003.1211369.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Boot-strapping language-image pre-training with frozen image encoders and large language models, 2023. URL https://arxiv.org/abs/2301.12597.
- Li, S. and Liu, Y. Sharper generalization bounds for clustering. In Meila, M. and Zhang, T. (eds.), Proceedings of the 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pp. 6392–6402. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/ v139/li21k.html.
- Liaw, C., Mehrabian, A., Plan, Y., and Vershynin, R. A simple tool for bounding the deviation of random matrices on geometric sets, 2016. URL https://arxiv.org/abs/1603.00897.
- Liu, W., Wen, Y., Yu, Z., and Yang, M. Large-margin softmax loss for convolutional neural networks, 2017. URL https://arxiv.org/abs/1612.02295.
- 586 Liu, W., Wen, Y., Yu, Z., Li, M., Raj, B., and Song, L.
 587 Sphereface: Deep hypersphere embedding for face recog588 nition, 2018. URL https://arxiv.org/abs/
 589 1704.08063.

- Liu, Z., Luo, P., Qiu, S., Wang, X., and Tang, X. Deepfashion: Powering robust clothes recognition and retrieval
 with rich annotations. In *Proceedings of the IEEE Con- ference on Computer Vision and Pattern Recognition*(CVPR), June 2016.
- Louart, C., Liao, Z., and Couillet, R. A random matrix
 approach to neural networks, 2017. URL https://
 arxiv.org/abs/1702.05419.
- Mei, S. and Montanari, A. The generalization error of random features regression: Precise asymptotics and double descent curve, 2020. URL https://arxiv.org/ abs/1908.05355.

- Moniri, B., Lee, D., Hassani, H., and Dobriban, E. A theory of non-linear feature learning with one gradient step in two-layer neural networks, 2024. URL https://openreview.net/forum?id=MY8SBpUece.
- Movshovitz-Attias, Y., Toshev, A., Leung, T. K., Ioffe, S., and Singh, S. No fuss distance metric learning using proxies, 2017. URL https://arxiv.org/abs/1703. 07464.
- Musgrave, K., Belongie, S., and Lim, S.-N. A metric learning reality check, 2020. URL https://arxiv.org/ abs/2003.08505.
- O'Donnell, R. Analysis of boolean functions, 2021. URL https://arxiv.org/abs/2105.10386.
- Papa, G., Clémençon, S., and Bellet, A. Sgd algorithms based on incomplete u-statistics: Large-scale minimization of empirical risk. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and Garnett, R. (eds.), Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/67e103b0761e60683e83c559be18d40c-Paper.pdf.
- Papyan, V., Han, X. Y., and Donoho, D. L. 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.
- Qian, Q., Shang, L., Sun, B., Hu, J., Li, H., and Jin, R. Softtriple loss: Deep metric learning without triplet sampling, 2020. URL https://arxiv.org/abs/ 1909.05235.
- Roth, K., Milbich, T., Sinha, S., Gupta, P., Ommer, B., and Cohen, J. P. Revisiting training strategies and generalization performance in deep metric learning. In *International Conference on Machine Learning*, pp. 8242–8252. PMLR, 2020.
- Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., and Komatsuzaki, A. Laion-400m: Open dataset of clipfiltered 400 million image-text pairs, 2021. URL https: //arxiv.org/abs/2111.02114.
- Seidenschwarz, J., Elezi, I., and Leal-Taixé, L. Learning intra-batch connections for deep metric learning, 2021. URL https://arxiv.org/abs/2102.07753.

605 606 607 608 609 610 611 612 613 614 615	Sohoni, N., Dunnmon, J., Angus, G., Gu, A., and Ré, C. No subclass left behind: Fine-grained robustness in coarse-grained classification problems. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 19339–19352. Curran Associates, Inc., 2020. URL https://proceedings.neurips. cc/paper_files/paper/2020/file/ e0688d13958a19e087e123148555e4b4-Paper. pdf.	 Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. How transferable are features in deep neural networks?, 2014. URL https://arxiv.org/abs/1411.1792. Zavatone-Veth, J. A., Yang, S., Rubinfien, J. A., and Pehlevan, C. Neural networks learn to magnify areas near decision boundaries, 2023. Zhai, A. and Wu, HY. Classification is a strong baseline for deep metric learning, 2019. URL https://arxiv.org/abs/1811.12649.
616 617 618 619	Song, H. O., Xiang, Y., Jegelka, S., and Savarese, S. Deep metric learning via lifted structured feature em- bedding, 2015. URL https://arxiv.org/abs/ 1511.06452.	Zhou, J., Li, X., Ding, T., You, C., Qu, Q., and Zhu, Z. On the optimization landscape of neural collapse under mse loss: Global optimality with unconstrained features, 2022. URL https://arxiv.org/abs/2203.01238.
620 621 622 623 624	Szegő, G. Orthogonal Polynomials. American Math. Soc: Colloquium publ. American Mathematical Society, 1975. ISBN 9780821810231. URL https://books. google.co.kr/books?id=ZOhmnsXlcYOC.	Zhou, J., Wang, P., and Zhou, DX. Generalization analysis with deep relu networks for metric and similarity learn- ing, 2024. URL https://arxiv.org/abs/2405. 06415.
625 626 627 628 629	Talwalkar, A., Kumar, S., and Rowley, H. Large-scale manifold learning. In 2008 IEEE Conference on Computer Vision and Pattern Recognition, pp. 1–8, 2008. doi: 10.1109/CVPR.2008.4587670.	Zhu, Z., Ding, T., Zhou, J., Li, X., You, C., Sulam, J., and Qu, Q. A geometric analysis of neural collapse with unconstrained features, 2021. URL https://arxiv. org/abs/2105.02375.
630 631 632 633	Tirer, T. and Bruna, J. Extended unconstrained features model for exploring deep neural collapse, 2022. URL https://arxiv.org/abs/2202.08087.	
634 635 636 637 638	Vershynin, R. Introduction to the non-asymptotic analysis of random matrices. <i>Chapter 5 of: Compressed Sensing,</i> <i>Theory and Applications. Edited by Y. Eldar and G. Ku-</i> <i>tyniok. Cambridge University Press, 2012,</i> 11 2010. URL https://arxiv.org/pdf/1011.3027.pdf.	
639 640 641 642 643	Vershynin, R. <i>High-Dimensional Probability: An Introduc-</i> <i>tion with Applications in Data Science.</i> Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.	
644 645 646	Vignat, C. A generalized isserlis theorem for location mixtures of gaussian random vectors, 07 2011. URL https://arxiv.org/pdf/1107.2309.pdf.	
648 649	Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The caltech-ucsd birds-200-2011 dataset. 2011.	
650 651 652 653 654	Wang, H., Wang, Y., Zhou, Z., Ji, X., Gong, D., Zhou, J., Li, Z., and Liu, W. Cosface: Large margin cosine loss for deep face recognition, 2018. URL https://arxiv. org/abs/1801.09414.	
655 656 657 658 659	Yang, Y., Steinhardt, J., and Hu, W. Are neurons actually collapsed? on the fine-grained structure in neural representations, 2023. URL https://arxiv.org/abs/2306.17105.	

A. Additional Related Works 660

677

680

681

682 683

684

685

686 687

688 689

690

691

692 693

694 695

696 697

698 699

700 701

704 705

706

709

713

714

661 Feature Transferability in Deep Metric Learning The explanation for how Deep Metric Learning learns transferable 662 features towards unseen data remains insufficient. Chopra et al. (2005) suggested that CNNs' robustness to geometric 663 distortions enables the creation of generalizable features. This explanation has been replaced in transformer-based research 664 by the idea that, without the inductive biases of CNNs, transformers are less constrained and thus capable of extracting 665 generalizable features (El-Nouby et al., 2021; Caron et al., 2021). Additionally, following the manifold hypothesis (Chang 666 et al., 2003; Lee et al., 2003; Talwalkar et al., 2008; Goodfellow et al., 2016), Liu et al. (2018); Ermolov et al. (2022) 667 explained that normalized softmax for metric learning works well because hyperspherical/hyperbolic feature space and the 668 data lies on a manifold. However, these studies do not provide a detailed analysis of how features are learned and transferred 669 through classification. 670

671 Neural Collapse (NC) and Features learned by Classifiers There exist studies exploring Neural Collapse (NC) and 672 features learned by classifiers that cannot be explained under the free variable assumption. Hui et al. (2022) argue that 673 NC does not manifest on test data. Sohoni et al. (2020); Yang et al. (2023) claim that even on training data, NC is not 674 fully realized, with critical fine-grained structures concealed. Notably, Yang et al. (2023) utilized a two-layer network to 675 analyze training data features. Regarding NC on novel data, Galanti et al. (2022) statistically analyze NC in transfer learning, 676 suggesting that NC generalizes not only to new samples within training classes but also to unseen classes with empirical observations. However, their analysis is constrained by focusing on general function spaces rather than specific neural 678 network architectures. 679

MSE for Classification 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).

Generalization Bound for Metric Learning Research on the generalization bounds of metric learning related to the U-process we use is also ongoing (Bellet & Habrard, 2015; Huai et al., 2019; Zhou et al., 2024). However, these studies do not analyze the exact feature learning structure.

B. Empirical Insights into High-Dimensional Asymptotics

In asymptotic analysis, $n, d, N \rightarrow \infty$ is crucial for observe result. Please see Figure 12, Figure 13 for the cohesion and Separability in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} . As the dimension increases, the range where cohesion and Separability align with our expectations expands.

For component analysis, please see Figure 14, Figure 15, Figure 16, Figure 17, Figure 18, Figure 19

C. Additional Observation of Multi Classes Feature Analysis

See Figure 21 for multi-directional training result. For F_0^L , and spike term depiced in Figure 22, Figure 23.

D. Additional Results of two-classes Experiments

D.1. Additional setup for Experiment I, II, III

We set $d = n = N = 2^{11}$ and use Shifted ReLU. We repeat each experiment with 3 different initializations of the neural network parameters.

Training Datasets (Data 1) two uniform distributions over a radius- \sqrt{d} ball, (Data 2) two multi-dimensional element-wise truncated Gaussian distributions, and (Data 3) two uniform distributions over a radius- \sqrt{d} sphere, symmetric about the origin ². The two means of training class are denoted as v and -v, respectively. For Data 1, $3 v \triangleq 2r \cdot \mathbf{u}$, with $\mathbf{u} \sim \text{Unif}(\mathbb{S}^{d-1})$. For Data 2, one class has support on $[1, \infty)$ across all dimensions, while the other class has support on $(-\infty, -1]$.

710 **Evaluation Datasets** Eval 1, 2 use the projected Gaussian distribution, which is projected onto the mean direction of 711 one training data v, as defined in equation 9. For Eval 1, we translate mean of projected Gaussian distribution with e, and 712

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



Figure 12: Cohesion in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right), with the computed range expanding from top to bottom.

Submission and Formatting Instructions for ICML 2025



Figure 13: Separability in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right), with the computed range expanding from top to bottom.



Figure 14: Component analysis of Cohesion in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right) in range [-100, 100], top: the dominant last component, bottom: sum of the other terms.



Figure 15: Component analysis of Cohesion in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right) in range [-500, 500], top: the dominant last component, bottom: sum of the other terms.

for Eval 2, we Rotate mean of projected Gaussian distribution with $R \in \mathcal{R}$ and fixed *e*. We generate 300 distinct rotation matrices \mathcal{R} using the process in Appendix O. The projected gaussian distribution is sampled as follows,

$$z - \frac{z^{\top} \nu \nu}{\|\nu\|^4} + \nu, \quad \text{where} \quad z \sim \mathcal{N}(0, cI).$$
(9)

For Eval 1, $\nu \triangleq ev$, c = 1 and for Eval 2, $\nu \triangleq Rev$, $c = 10^{-1}$ with e = 0.01 for Data 2 experiment and e = 0.008 for Data 1 and 3 experiments, $R \in SO(d)$.

D.2. Comprehensive Results of All Experiments

The overall experimental results for *Cohesion* and *Separability* are shown in Figure 24. The results for Eval 1 experimental settings are presented in linear scale in Figure 25 and in logarithmic scale in Figure 26. Additionally, as presented in Figure 7, experiments for Eval 2 settings on Data 2 and 3 are shown in linear scale in Figure 27, with results for *Cohesion*, *Separability*, and Recall@1 (IP). Furthermore, results for Recall@1 (cos) are presented in linear scale in Figure 28. All observed results align with the theoretical predictions.

E. Additional Results of Real-world dataset Experiments

Figure 29 summarizes the experimental results and the purpose of the experiment. Expr. IV is in Figure 30, 31, 1. Expr. V is e in Figure 32, Table 2. Expr. VI is in Figure 33. Expr. VII is in Figure 34, 35, 36, 37, Table 3, and 4.

E.1. Relation between Expr. V and VI

On the other hand, certain results from Expr. V align with those from Expr. VI. As shown in Table 5, for datasets such as CAR and CUB, the number of additional classes introduced by the *sub In1k* dataset is significantly larger compared to SOP. For these data, inclusion of the additional *sub In1k* dataset contributes to improved *recall@1* performance when trained using a Random Initialized Network. Meanwhile, the performance of the pre-trained network is not significantly affected by the additional dataset. We attribute this to the fact that the pre-trained model is additionally re-trained using the same ImageNet dataset *sub In1k*. These findings suggest that further research on the behavior of pre-trained networks is necessary.

E.2. Expr. VII: Removing Duplicately Assigned Eval Classes

In **Expr. VII**, as suggested by the theoretical results on *Separability*, we validated whether eliminating duplicate in the *assignments* improves performance. To clarify, we will provide an example of duplicate *assignment* at Note E.1.

Note E.1 (Example of duplicate assignment). For two train classes $\mathscr{C}_1^{(train)}$, $\mathscr{C}_2^{(train)}$ and two test classes $c_1^{(test)}$, $c_2^{(test)}$, if most instances of $c_1^{(test)}$ and $c_2^{(test)}$ are classified as $\mathscr{C}_1^{(train)}$, both test classes are assigned to $\mathscr{C}_1^{(train)}$, resulting in duplication. Conversely, if $c_1^{(test)}$ is classified as $\mathscr{C}_2^{(train)}$ and $c_2^{(test)}$ as $\mathscr{C}_1^{(train)}$, they are assigned without duplication.

To validate, we introduce treatment and control groups. For treatment group, we eliminate duplicate in the textitassignments for the train classes, i.e., for each unseen class, the most frequently classified training class is aggregated, and the classes are randomly removed to ensure that the selected training classes become unique (2). For the control group, we performed random selection of the same number of classes of treatment group (1). These two groups are evaluated using *recall@1*. This process was repeated five times, and the average was reported. The experimental results are presented in 34, 35, 36, 37, Table 3, and 4. A total of 64 experiments are conducted, of which 51 demonstrated performance improvements: the estimated success rate is 79%. There is a $1.73\% \pm 2.87\%$ average improvement in recall@1, with a maximum improvement of 13.65%, a minimum decrease of -3.28%. These findings suggest that the duplicate reduction treatment group outperforms the randomly removed group with a binomial test p-value of 9.40×10^{-7} .

F. Limitations and Future Work

While our study provides valuable insights into feature learning and transferability, several important directions remain for future research. First, while the Hermite approximation aided our feature analysis, it posed numerical challenges due to the discrepancy between polynomials and nonlinear neural networks. Specifically, the need for extremely high-dimensional approximations Figure 2 and the lack of precise scaling alignment between the approximation and the neural networks in

Algorithm 1 Random Sampling	
Input: Number if unseen classes u , number of classes $ L $	
Output: Sampled class set S _{random}	
Set $S_{random} \leftarrow random.sample(\{0, 1, \dots, u-1\}, L)$	
return S _{random}	
Algorithm 2 Duplicated againment reduction compliant	
Algorithm 2 Duplicated assignment reduction sampling	
Input: Model f, unseen data loader \mathcal{D} , number of train classes C_{train} , number of u	unseen classes C_{unseen}
Output: Sampled class set S_{nondup}	
Initialize counter matrix counter $\leftarrow 0^{\text{Cunseen} \times \text{Ctrain}}$	
for (img, label) in ${\mathcal D}$ do	
$\texttt{pred} \gets f(\texttt{img})$	Predicted class indices
Update counter: counter[label,pred] += 1	
end for	
$top1_index \leftarrow argsort(counter, dim = 1, descending = True)[, 0]$	
unique_label unique(top1_index)	
Initialize $S_{\text{nondun}} \leftarrow \emptyset$	
for each label ℓ in unique_label do	
$I_{\ell} \leftarrow \{i \mid \texttt{top1_index}[i] = \ell\}$	Indices corresponding to label ℓ
$i_{\text{sample}} \leftarrow \text{random.sample}(I_{\ell}, 1)$	Select one random index
$S_{\text{nondup}} \leftarrow S_{\text{nondup}} \cup \{i_{\text{sample}}\}$	
end for	
return Snordun	

961 finite dimensions Figure 4.

These limitations highlight the need for alternative approximation techniques or analytical approaches. Second, the relationship between semantic similarity and train-unseen similarity requires further theoretical exploration. Third, an important direction for future research is expanding the concepts of cohesion and Separability to multi-class softmax classification problems, incorporating normalization and temperature scaling to better align with practical settings or Neural Collapse research. Finally, recently Zavatone-Veth et al. (2023) suggest neural networks tend to compress the feature space around training data while expanding the regions between decision boundaries. We consider this phenomenon appears closely related to the train-unseen similarity-driven cohesion and Separability observed in our study. Investigating this connection through the lens of Riemannian geometry could yield novel insights into the fundamental structure of learned representations.



Figure 16: Component analysis of Cohesion in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right) in range [-1000, 1000], top: the dominant last component, bottom: sum of the other terms.



Figure 17: Component analysis of Separability in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right) in range [-500, 500], top: the dominant last component, bottom: sum of the other terms.



Figure 18: Component analysis of Separability in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right) in range [-500, 500], top: the dominant last component, bottom: sum of the other terms.

Submission and Formatting Instructions for ICML 2025

Figure 19: Component analysis of Separability in \mathbb{R}^{2000} , \mathbb{R}^{20000} , \mathbb{R}^{320000} (left to right) in range [-1000, 1000], top: the dominant last component, bottom: sum of the other terms.

Submission and Formatting Instructions for ICML 2025

by spikes, so when using only the β_1 or β_4 spikes, the two features are always mapped to the same position.

Submission and Formatting Instructions for ICML 2025

1422 Figure 25: Expr. I: translation(*e*) variation case (linear scale). — is after one step training. — is from initialization. As the *train-unseen similarity* increases, both cohesion and Separability become larger due to pos-neg setup.

¹⁴⁷⁷ Figure 26: Expr. I: translation(*e*) variation (log scale). — is after one step training. — is from initialization. As the ¹⁴⁷⁸ *train-unseen similarity* increases, both cohesion and Separability become larger due to pos-neg setup.

Figure 27: Expr. II, Expr. III: rotation(R) variation (linear scale). — is after one step training. — is from initialization. Expr. II is pos-neg. Expr. III is pos-pos.

Submission and Formatting Instructions for ICML 2025

Figure 31: Expr. IV on ResNet18, ResNet50 with Domain + In(S) e.g. CAR+I(V), CUB+I(B), SOP+I(P), ISC+I(C)

Figure 32: Expr. V, additional results, it is represented as follows Domain Domain + Related Subset of In1k Domain
+ Whole In1k subsampled Adding unrelated classes for training does not significantly affect the performance.

Figure 33: Expr VI, it is represented as follows: ResNet18 - - , ResNet50 —, Dataset car, cub, sop, isc. As the steps increased and related classes were added, performance generally improved consistently.

			Tat	ole 1: Table	results for Expr. IV	1			
					1				
R	esNet18 (R	andomly I	(nitialized)		R	esNet50 (F	Randomly	Initialized))
	CAR	CUB	SOP	ISC		CAR	CUB	SOP	ISC
CAR+I(V)	0.3922	0.0847	0.3126	0.2079	CAR+I(V)	0.3280	0.0879	0.3226	0.2000
CAR	0.2383	0.0685	0.2766	0.1994	CAR	0.2067	0.0495	0.2611	0.1583
I(V)	0.1117	0.0618	0.2610	0.1793	I(V)	0.1048	0.0459	0.2670	0.1410
CUB+I(B)	0.1456	0.1205	0.3117	0.2067	CUB+I(B)	0.0755	0.0527	0.2303	0.1414
CUB	0.1432	0.1089	0.3179	0.1998	CUB	0.0626	0.0393	0.1950	0.1081
I(B)	0.0973	0.0640	0.2658	0.1703	I(B)	0.0456	0.0358	0.1954	0.1074
SOP+I(P)	0.1753	0.0748	0.3720	0.3304	SOP+I(P)	0.1662	0.0829	0.3812	0.2934
SOP	0.1754	0.0876	0.3790	0.3306	SOP	0.1725	0.0743	0.3750	0.2754
I(P)	0.1405	0.0586	0.3129	0.2327	I(P)	0.0940	0.0422	0.2716	0.1697
ISC+I(C)	0 1409	0.0613	0.3295	0.4870	I(C) ISC+I(C)	0.1090	0.0550	0.3001	0.5318
ISC	0.1328	0.0685	0.3338	0.4887	ISC	0.1022	0.0503	0.2699	0.458
	0.0908	0.0003	0.3330	0.1823	ISC I(C)	0.0625	0.0303	0.2099	0 1446
(C)	0.0700	0.0471	0.2403	0.1025	1(0)	0.0025	0.0412	0.2274	0.1440
Res	sNet18 (Im	ageNet 1K	Pretrained	ł)	Resl	Net50 (Ima	ageNet 1K	Pretrained	l)
	CAR	CUB	SOP	ISC		CAR	CUB	SOP	ISC
CAR + I(V)	0.8610	0.1131	0.4104	0.2133	$CAR \perp I(V)$	0 9081	0.1268	0.4192	0.1805
	0.8690	0.1101	0.4104	0.2155	CARTI(V)	0.0078	0.1200	0.4192	0.1603
	0.4210	0.1608	0.3900	0.1951		$\frac{0.9078}{0.4013}$	0.1648	0.3943	0.1075
I(V)	0.4210	0.1098	0.4018	0.2307	I(V)	0.4013	0.1040	0.4613	0.2330
	0.3474	0.5269	0.4743	0.2171	CUP	0.2031	0.5057	0.4380	0.1095
	0.3470	0.3300	0.4072	0.2327		0.3073	0.2227	0.4794	0.2205
I(D)	0.5771	0.5400	0.5002	0.2278	I(D)	0.3212	0.3337	0.4/81	0.1840
SOF + I(F)	0.4075	0.1303	0.4773	0.2027	SOP + I(r)	0.4002	0.2204	0.0307	0.3702
SOP	0.5802	0.1499	0.4827	0.5201	SOP L(D)	0.4000	0.2200	0.0270	0.3700
I(P)	0.4003	0.2076	0.4838	0.2569	I(P)	0.3547	0.2208	0.4602	0.2337
ISC+I(C)	0.2420	0.09/6	0.4616	0.7098	ISC+I(C)	0.2301	0.1207	0.5376	0.8/18
ISC	0.2130	0.0847	0.4550	0.7115	ISC	0.2230	0.1274	0.5390	0.8/10
I(C)	0.3738	0.2227	0.4994	0.2457	I(C)	0.3655	0.2311	0.5167	0.2413
		,	Table 2: Ta	able results	of performance for	Expr. V.			
_					_				
R	esNet18 (R	andomly I	nitialized)		R	esNet50 (F	Randomly	Initialized)	
	CAR	CUB	SOP	ISC		CAR	CUB	SOP	ISC
D	0.2383	0.1089	0.3790	0.4887	D	0.2067	0.0393	0.3750	0.4581
D+I(Sub)	0.3922	0.1205	0.3720	0.4870	D+I(Sub)	0.3280	0.0527	0.3812	0.5318
D+I	0.3074	0.1404	0.3591	0.4532	D+I	0.3276	0.0968	0.3726	0.4992
				·					
Res	sNet18 (Im	ageNet 1K	Pretrained	1)	Resl	Net50 (Ima	ageNet 1K	Pretrained	l)
	CAR	CUB	SOP	ISC		CAR	CUB	SOP	ISC
D	0.8680	0.5366	0.4827	0.7115	D	0.9078	0.5778	0.6276	0.8710
D+I(Sub)	0.8610	0 5289	0 4775	0 7098	D+I(Sub)	0.9081	0.5657	0.6367	0.8718
	0.7604	0.5257	0.4766	0.6907		0.7602	0.4690	0.6260	0.0710
D+I	0.7604	0.5357	0.4/66	0.089/	D+I	0.7603	0.4089	0.0360	0.8481

	m '	m	D 1	•	m 1
Test	Train	Treatment	Random	Δ	Total
	CAR	35.26	34.49	0.77	23.83
CAP Test	I(V)	29.02	26.96	2.06	11.17
CAR lest	CAR+I(V)	49.27	46.73	2.54	39.22
	In	39.51	36.48	3.03	25.70
	CUB	20.01	18.49	1.52	10.89
CUP Test	I(B)	16.36	14.75	1.61	6.40
COD lest	CUB+I(B)	19.58	18.39	1.19	12.05
	In	32.16	28.74	3.42	21.49
	ISC	60.64	59.45	1.19	48.87
ISC Test	I(C)	60.93	57.78	3.15	18.23
ISC Test	ISC+I(C)	59.59	59.11	0.48	48.70
	In	45.01	46.92	-1.91	24.75
	SOP	43.58	42.96	0.62	37.90
SOD Test	I(P)	49.76	48.45	1.31	31.29
SOP lest	SOP+I(P)	42.57	42.12	0.45	37.20
	In	51.84	54.03	-2.19	38.82
Average In	provement			1.20	

Table 3: Expr. VII from (Randomly Initialized)

ResNet18 (Randomly Initialized)

ResNet50 (Randomly Initialized)

0.875

Success Rate

	T :		D 1	•	T (1
lest Irain		Treatment	Random	Δ	Total
	CAR	30.45	29.95	0.50	20.67
CAD Test	I(V)	24.49	22.16	2.33	10.48
CAR lest	CAR+I(V)	42.25	42.67	-0.42	32.80
	In(CAR)	51.69	42.39	9.30	30.06
	CUB	13.24	15.84	-2.60	3.93
CAD Test	I(B)	21.20	16.65	4.55	3.58
CAR lest	CUB+I(B)	16.30	13.66	2.64	5.27
	In	48.10	39.59	8.51	28.06
	ISC	60.63	59.41	1.22	45.81
	I(C)	53.67	51.22	2.45	14.46
CAR lest	ISC+I(C)	67.46	66.88	0.58	53.18
	In	44.60	44.85	-0.25	22.85
	SOP	44.02	43.34	0.68	37.50
	I(P)	44.93	45.22	-0.29	27.16
CAR lest	SOP+I(P)	43.51	43.70	-0.19	38.12
	In	59.49	59.47	0.02	42.93
Average Im	provement			1.81	
Success Ra	te			0.6875	

)
)

		1			
Test Train		Treatment	Random	Δ	Total
	CAR	90.90	90.33	0.57	86.80
CAD Test	I(V)	64.51	65.03	-0.52	42.10
CAR lest	CAR+I(V)	90.06	88.79	1.27	86.10
	In(CAR)	71.77	73.08	-1.31	26.00
	CUB	66.19	63.12	3.07	53.66
CAD Test	I(B)	48.67	46.90	1.77	34.00
CAR lest	CUB+I(B)	64.48	63.89	0.59	52.89
	In	44.95	39.30	5.65	30.32
	ISC	78.81	77.15	1.66	71.15
CAD Test	I(C)	70.48	66.47	4.01	24.57
CAR lest	ISC+I(C)	78.58	77.35	1.23	70.98
	In	32.65	35.78	-3.13	13.85
	SOP	52.45	51.81	0.64	48.27
CAD Test	I(P)	66.72	66.81	-0.09	48.38
CAR lest	SOP+I(P)	51.34	51.01	0.33	47.75
	In	46.31	46.95	-0.64	30.66
Average Im	provement			0.94	
Success Ra	te			0.6875	

ResNet18 (ImageNet 1K Pretrained)

ResNet50 (ImageNet 1K Pretrained)

Test	Test Train		Random	Δ	Total
	CAR	93.78	93.57	0.21	90.77
CAD Test	I(V)	70.12	63.34	6.78	40.13
CAR lest	CAR+I(V)	94.45	93.34	1.11	90.81
	In(CAR)	84.20	77.43	6.77	32.51
	CUB	71.44	68.51	2.93	57.78
CAD Test	I(B)	47.78	46.19	1.59	33.37
CAR lest	CUB+I(B)	70.63	67.15	3.48	56.56
	In	75.96	62.32	13.64	35.53
	ISC	91.35	90.49	0.86	87.10
CAD Test	I(C)	68.62	71.90	-3.28	24.13
CAR lest	ISC+I(C)	91.60	90.59	1.01	87.18
	In	39.54	35.39	4.15	8.68
	SOP	68.40	68.07	0.33	62.75
CAD Test	I(P)	66.24	64.09	2.15	64.02
CAR lest	SOP+I(P)	68.83	68.40	0.43	63.66
	In	59.94	54.78	5.16	28.87
Average Im	provement	•		2.96	
Success Ra	te			0.9375	

2145 G. Additional Notations

2158 2159

2169 2170

2171 2172

2176

2178

2180 2181

2182

2195

2196 2197

The operator diag(·) creates a matrix with the elements of the input vector placed along the diagonal. Let $\mathbf{1}_{\text{condition}}$ be 1 if the condition is true and 0 otherwise. Let m! be factorials of m. Let n!! be double factorial. We define (-1)!! = 0!! = 1. For $o_{\mathbb{P}}, O_{\mathbb{P}}, \Theta_{\mathbb{P}}$ notations we follow Moniri et al. (2024) $\|\cdot\|_F$ 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. $\lfloor n/2 \rfloor$ denotes the floor of n/2. $\Gamma(z)$ is the Gamma function.

Additional information of Hermite Polynomials We employ the probabilist's Hermite polynomials (Szegő, 1975; Bienstman, 2023; Moniri et al., 2024). We denote $H_k(x)$ as k-th Hermite polynomial.

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 .

Hermite polynomials can be represented as the following explicit form:

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

2160 2161 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 weight function $\phi(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$ on the interval $(-\infty, \infty)$. Their orthogonality condition is given by:

$$\int_{-\infty}^{\infty} H_m(x) H_n(x) \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = n! \mathbf{1}_{m=n}.$$

2173 H. hermite coef of shifted ReLU

 $^{2174}_{2175}$ One of the activation function that satisfy our condition 2.1 is shifted ReLU,

$$\sigma(x) = \max(0, x) - \frac{1}{\sqrt{2\pi}}.$$

2179 This allow hermite decomposition with coefficient is calculated as

$$c_n = \frac{1}{n!} \mathbb{E}_z[\sigma(z)H_n(z)]$$

21832184Then for the zero-th coefficient is calculated as

$$c_0 = \mathbb{E}_z[\sigma(z) \times 1] = \mathbb{E}_z[\max(0, x)] - \frac{1}{\sqrt{2\pi}}$$
$$= \int_0^\infty x\phi(x)dx - \frac{1}{\sqrt{2\pi}} = 0$$
(10)

By the way, if $n \neq 0$, $\mathbb{E}[\frac{1}{\sqrt{2\pi}} \times H_n] = \frac{1}{\sqrt{2\pi}} \mathbb{E}[1 \times H_n] = \frac{1}{\sqrt{2\pi}} \mathbb{E}[H_0 \times H_n] = 0$ by orthogonality. Thus, shift is only effects on n = 0.

The coefficient c_n of the expansion of Shifted-ReLU is defined as:

$$c_{n} = \begin{cases} 0, & \text{if } n = 0, \\ \sum_{m=0}^{\lfloor n/2 \rfloor} \frac{(-1)^{m} \cdot 2^{\frac{n-2m}{2} - m} \cdot \Gamma\left(\frac{n-2m+2}{2}\right)}{m! \cdot (n-2m)! \cdot \sqrt{2\pi}}, & \text{otherwise.} \end{cases}$$
(11)

We directly calculated equation 11 and obtained the following result in Figure 38.

Figure 38: Hermite Coefficient of Shifted ReLU

I. Proof of Theorem 3.1

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.

We will omit the subscript ij since it does not cause any confusion untill equation 35. The following statements hold for $\forall ij$. For the aforementinoed \mathbb{A} , \mathbb{B} , and \mathbb{C} , we obtain bounds for each operator norm as follows

Lemma I.1.

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

Proof of Lemma I.1 (\mathbb{A}). We obtain

$$\mathbb{A} = \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} X^{\top} y a^{\top}.$$
(13)

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

$$\|\mathbb{A}\| = \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X^{\top} y a^{\top}\| = \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X^{\top} y\|_2 \|a\|_2 = \frac{c_1}{\mathbf{n}\sqrt{\mathbf{N}}} (y^{\top} X X^{\top} y)^{1/2} \|a\|_2$$
(14)

 $||a||_2$ study By definition, $a \sim N(0, \frac{1}{N})$, so $\sqrt{N}\alpha[i]$ is a sub-Gaussian. Use Thm 3.3.1 in Vershynin (2018),

$$\mathbb{P}\left(\left|\|\sqrt{\mathbf{N}}\alpha\| - \sqrt{\mathbf{N}}\right| \ge t\right) \le 2e^{-ct^2} \quad \text{let } t = \sqrt{\mathbf{N}} \\
\mathbb{P}(\|\alpha\|_2 \le 1) \le 2e^{-c\mathbf{N}}$$
(15)

 $(y^{\top}XX^{\top}y)^{1/2}$ study Note that the U, V matrices resulting from the SVD belong to the O-group, so there is no length transformation. $y^{\top}XX^{\top}y = \|X^{\top}y\|_{2}^{2} = \|U\Sigma V^{\top}y\|_{2}^{2} = \|\Sigma V^{\top}y_{1}\|$ $=\sum_i \sigma_i^2 |V^\top y_i|^2 \ge \sigma_{\min}^2 \sum_i |V^\top y_i|^2 = \sigma_{\min}^2 ||y||_2^2 = \mathbf{n}\sigma_{\min}^2$ (16)

We get $(y^{\top}XX^{\top}y)^{1/2} \ge \sqrt{\mathbf{n}}\sigma_{\min}$. 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)

Therefore, $\mathbb{P}(\|\mathbb{A}\| \le C(\frac{1}{\sqrt{\mathbf{N}}} - C\frac{\sqrt{d}}{\sqrt{\mathbf{nN}}})) \le 2(e^{-c\mathbf{N}} + e^{-c\mathbf{n}}).$

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

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

Lemma I.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

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

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

$$\square$$

$$(19)$$

Fact I.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. $\|\sigma(v)\|_{\psi_2} \le \|\lambda\|_{\psi_2} < \infty.$

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

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

$$\|\mathbb{B}\| \leq \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X\| \|y a^{\top} \odot \sigma'_{\perp} (XW_0)\|$$

$$\leq \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X\| \|y a^{\top} \odot \sigma'_{\perp} (XW_0)\|$$

$$\leq \frac{1}{\mathbf{n}\sqrt{\mathbf{N}}} \|X\| \|y\|_{\infty} \|\sigma'_{\perp} (XW_0)\| \|a\|_{\infty}$$

$$(21)$$

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

 $= \frac{1}{n\sqrt{N}} \|X\| \|\sigma'_{\perp}(XW_0)\|\|a\|_{\infty}$

$$\mathbb{P}(\{\|\sigma_{\perp}'\| \ge C(\sqrt{\mathbf{n}} + \sqrt{\mathbf{N}})\lambda_{\sigma}B\}, \mathcal{A}_B) \le 2e^{-c\mathbf{N}}$$
$$\mathcal{A}_B = \{\{\|W_0\| \le B\} \cup \{\sum_{i=1}^{\mathbf{N}} (\|W_{0,i}\|^2 - 1)^2 \le B^2\}\}.$$
(22)

• •

Therefore,

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

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

2365 By Lemma I.3 $\mathbb{P}(||a||_{\infty} \le t/\sqrt{N}) \ge 1 - 2Ne^{-ct^2}$, and Lemma L.3 with $t = \sqrt{d}$

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

Remark I.5. In the proportional regime, as $\mathbf{n}, \mathbf{d}, \mathbf{N} \to \infty$, these quantities can be interchanged to a constant. Thus, Lemma I.1 is reformulated as follows 2373 $\mathbb{P}(||\mathbb{A}|| \le \kappa/\sqrt{\mathbf{n}}) \le Ce^{-c\mathbf{n}})$

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

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

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

Also, for gradient, we have

$$||G|| = ||\mathbb{A} + \mathbb{B} + \mathbb{C}|| \le ||\mathbb{A}|| + ||\mathbb{B}|| + ||\mathbb{C}|| = O_{\mathbb{P}}(\frac{1}{\sqrt{\mathbf{n}}} + \frac{\log \mathbf{n}}{\mathbf{n}} + \frac{\log^2 \mathbf{n}}{\mathbf{n}}) = O_{\mathbb{P}}(\frac{1}{\sqrt{\mathbf{n}}})$$
(35)

Now we denote subscript ij for summary.

2387 Proof of Theorem 3.1. Using $||G_{ij} - \mathbb{A}_{ij}|| = ||\mathbb{B}_{ij} + \mathbb{C}_{ij}|| \le ||\mathbb{B}_{ij}|| + ||\mathbb{C}_{ij}||$ and Lemma I.5

$$\mathbb{P}\bigg(\|G_{ij} - \mathbb{A}_{ij}\| \ge C \frac{\log^2 \mathbf{n}}{\mathbf{n}}\bigg) \le \mathbb{P}\bigg(\|G_{ij} - \mathbb{A}_{ij}\| \ge C(\frac{\log n}{n} + \frac{\log^2 \mathbf{n}}{\mathbf{n}})\bigg) \le Cne^{-c\log^2 \mathbf{n}}.$$
(36)

23912392Therefore, almost surely, in the proportional limit,

$$\|G_{ij} - \mathbb{A}_{ij}\| \le C \frac{\log^2 \mathbf{n}}{\mathbf{n}} = \frac{\kappa}{\sqrt{\mathbf{n}}} \frac{C}{\kappa} \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} \le \|\mathbb{A}_{ij}\| \frac{C}{\kappa} \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} \le \kappa' \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} \left(\|G_{ij}\| + \|G_{ij} - \mathbb{A}_{ij}\| \right).$$
(37)

We get $(1 - \kappa' \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}}) ||G_{ij} - \mathbb{A}|| \le \kappa' \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} ||G_{ij}||$. For large enough \mathbf{n} for $1 - \kappa' \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} \ge \frac{1}{2}$,

$$||G_{ij} - \mathbb{A}_{ij}|| \le \kappa' \frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}} ||G_{ij}|| \le C \frac{\log^2 \mathbf{n}}{\mathbf{n}}$$

Sum up for $\forall ij$, 2402

$$||G - \sum_{i < j} \mathbb{A}_{ij}|| = ||\sum_{i < j} G_{ij} - \mathbb{A}_{ij}|| \le \sum_{i < j} ||G_{ij} - \mathbb{A}_{ij}|| \le C \frac{\log^2 \mathbf{n}}{\mathbf{n}}$$

J. Proof of Theorem 3.3

Lemma J.1. The following facts will be used in subsequent proofs. Remark $\beta_{ij} \triangleq \frac{1}{n} X_{ij}^T y$ in Theorem 3.2.

- ²⁴¹¹ ₂₄₁₂ A. $||X_{ij}|| = O_{\mathbb{P}}(\sqrt{\mathbf{n}}), ||y|| = O_{\mathbb{P}}(\sqrt{\mathbf{n}}), ||\beta_{ij}|| = O_{\mathbb{P}}(1)$
- 2413 B. $||X_{ij}\beta_{ij}a_{ij}|| = ||X\beta_{ij}||_2 ||a_{ij}||_2 = O_{\mathbb{P}}(\sqrt{n})$
- 2415 C. $||W_0|| = O_{\mathbb{P}}(1), ||W|| = ||W_0 + G|| \le ||W_0|| + ||G|| = O_{\mathbb{P}}(1)$

2416 D.
$$||X_{ij}G|| = O_{\mathbb{P}}(\sqrt{n})$$

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

 $F. \ M_b \triangleq \|X\beta\|_{\infty} = \max_{1 \le i \le \mathbf{n}} | < \tilde{X}[i], \beta > | \le C \log^{1/2} \mathbf{n}, \text{ w.p. } 1 - 2\mathbf{n}e^{-c\log \mathbf{n}}$ G. $M_{W_0} \triangleq \sup_{k>1} ||(W_0 W_0^{\top})^{\circ k}|| \le C \text{ w.p. } 1 - o(1)$ *H.* $||A^{\circ k}|| \le ||A||^k$ Proof. It is evident from Lemma L.3, equation 15 in the proportional regime, that A, B, C, and D hold. To proof E, F, and G, we employ proof techniques adapted from Moniri et al. (2024). For E, by Lemma I.3, with $t = \log^{\frac{1}{2}} n$, $M_a \leq \frac{C \log^{\frac{1}{2}} n}{\sqrt{n}}$, w.p. 1 - o(1). For F, $\mathbb{P}(C|x^T\beta| > t) = \mathbb{P}(C|x^T\beta - Ex^T\beta + Ex^T\beta| > t)$ (38) $<\mathbb{P}(C|x^T\beta - Ex^T\beta| > t - C|Ex^T\beta|) < 2\exp(-ct^2).$ Then, $\mathbb{P}(|x^T\beta| > t) < 2\exp(-c(t - Ex^T\beta)^2) < 2\exp(-ct^2)$. Therefore, $M_b < C \log^{\frac{1}{2}} n$, w.p. 1 - o(1) with $t = \log^{\frac{1}{2}} n$. For G, refer Moniri et al. (2024). For H, refer Bai & Silverstein (2010) Corollary A.21. **Corollary J.2** (Corollary of Theorem 3.1). By Lemma J.1, we have w.p. 1 - o(1), $\|\tilde{X}G - c_1 \tilde{X} \sum_{i=1}^{\infty} \beta_{ij} a_{ij}^T \| = O(\frac{\log^2 \mathbf{n}}{\mathbf{n}} \cdot \sqrt{\mathbf{n}}) = O(\frac{\log^2 \mathbf{n}}{\sqrt{\mathbf{n}}})$ (39)*Remark* J.3. $W_1 = W_0 + G$, so $\tilde{X}W_1 = \tilde{X}W_0 + \tilde{X}G$. \tilde{X} is i.i.d. copy of training data X We generalize Corollary J.2 i.e. monomial approximation of data-gradient product in polynomial form as Lemma J.4. **Lemma J.4** (Polynomial Approximation of Data-Gradient Product). For any $k \in \mathbb{N}$, sufficiently large **n**, and w.p. 1 - o(1), $\|(\tilde{X}G)^{ok} - c_1^k (\tilde{X}\sum_{i \leq i} \beta_{ij} a_{ij}^T)^{ok}\| = O(\mathbf{n}^{-\frac{k}{2}} \log^{2k} \mathbf{n})$ (40)*Proof of Lemma J.4.* k = 1 is trivial Corollary J.2. We follow Moniri et al. (2024) for $k \ge 2$. We need to show $\exists C > 0$, w.p. 1-o(1) $\|(\tilde{X}G)^{ok} - c_1^k (\tilde{X}\sum_{i \neq j} \beta_{ij} a_{ij}^T)^{ok}\| \le C \mathbf{n}^{-\frac{k}{2}\log^{2k} \mathbf{n}}$ (41) $(\tilde{X}G)^{ok} = (\tilde{X}G - c_1\tilde{X}\sum_{i < i}\beta_{ij}a_{ij}^T + c_1\tilde{X}\sum_{i < j}\beta_{ij}a_{ij}^T)^{ok}$ (42) $=\sum_{i=1}^{\kappa} (k_j) (\tilde{X}G - c_1 \tilde{X} \sum_{i < i} \beta_{ij} a_{ij}^T)^{oj} \odot (c_1 \tilde{X} \sum_{i < i} \beta_{ij} a_{ij}^T)^{o(k-j)} + c_1^k (\tilde{X} \sum_{i < j} \beta_{ij} a_{ij}^T)^{ok}$ Thus, $(\tilde{X}G)^{ok} - c_1^k (\tilde{X}\sum_{i \in \mathcal{A}} \beta_{ij} a_{ij}^T)^{ok}$ (43) $= \sum_{i=1}^{k} \binom{k}{j} (\tilde{X}G - c_1 \tilde{X} \sum_{i < i} \beta_{ij} a_{ij}^T)^{oj} \odot c_1^{k-j} (\sum_{i < j} (\tilde{X}\beta_{ij} a_{ij}^T))^{o(k-j)}$ Now we will show $||(\tilde{X}G - c_1\tilde{X}\sum_{i < j}\beta_{ij}a_{ij}^T)^{oj} \odot c_1^{k-j}(\sum_{i < j}(\tilde{X}\beta_{ij}a_{ij}^T))^{o(k-j)}|| = O_{\mathbb{P}}(\log^{k+j}\mathbf{n} \cdot \mathbf{n}^{-\frac{1}{2}k}).$

2475 $\|(\tilde{X}G - c_1\tilde{X}\sum_{i < j}\beta_{ij}a_{ij}^T)^{oj} \odot c_1^{k-j}(\sum_{i < j}(\tilde{X}\beta_{ij}a_{ij}^T))^{o(k-j)}\|$ 2476 2477 2478 $\leq C \| (\tilde{X}G - c_1 \tilde{X} \sum_{i < i} \beta_{ij} a^T)^{oj} \odot (\tilde{X}\beta a^T)^{o(k-j)} \|$ 2479 2480 $\leq C \|\operatorname{diag}(\tilde{X}\beta)^{ok-j}||_{op}||(\tilde{X}G^T - c_1\tilde{X}\sum_{i < j}\beta_{ij}a^T)^{oj}||_{op}||\operatorname{diag}(a)^{ok-j}\|$ (44)2482 $\leq C(M_a M_b)^{k-j} \| (\tilde{X}G - c_1 \tilde{X} \sum_{i < j} \beta_{ij} a^T)^{oj} \|^j$ 2483 2484 2485 $\leq C(n^{-\frac{1}{2}(k-j)}\log^{k-j}\mathbf{n})\log^{2j}\mathbf{n}\cdot\mathbf{n}^{-\frac{1}{2}j}$ 2486 $= O_{\mathbb{D}}(\mathbf{n}^{-\frac{1}{2}k}\log^{k+j}\mathbf{n})$ 2487 2488 2489 Therefore, $\|(\tilde{X}G)^{ok} - c_1^k (\tilde{X}\sum_{i < i} \beta_{ij} a_{ij}^T)^{ok}\| = O_{\mathbb{P}}(\mathbf{n}^{-\frac{k}{2}} \log^{2k} \mathbf{n})$ (45)2490 2491 2492 2493 **Lemma J.5.** Following condition in section 2, Assume event $\Omega = \sup_{k\geq 1} ||(W_0 W_0^T)^{ok}||_{op} \leq C$ occur, following statement 2494 holds. 2495 $||H_i(\tilde{X}W_0)||_{op} = O_{\mathbb{P}}(\sqrt{n}\log^{\frac{3}{2}}n\sqrt{j!})$ 2496 2497 Lemma J.6. Given random matrix A, Following statement holds, 2498 $\mathbb{P}(||A||_{op} \ge t) \le \mathbb{P}(||\frac{1}{n}AA^T - EAA^T||_{op} \ge \frac{t^2}{n} - ||EAA^T||_{op})$ 2499 2500 2501 Proof of Lemma J.6. 2502 $\mathbb{P}(||A||_{op} \ge t) = \mathbb{P}(||A||_{op}^2 \ge t^2) = \mathbb{P}(||\frac{1}{n}AA^T||_{op} \ge \frac{t^2}{n})$ 2504 $=\mathbb{P}(||\frac{1}{n}AA^{T} - EAA^{T} + EAA^{T}||_{op} \geq \frac{t^{2}}{n})$ 2506 (46) $\leq \mathbb{P}(||\frac{1}{n}AA^{T} - EAA^{T}||_{op} + ||EAA^{T}||_{op} \geq \frac{t^{2}}{n})$ 2508 $= \mathbb{P}(||\frac{1}{n}AA^{T} - E(AA^{T})||_{op} \ge \frac{t^{2}}{n} - E||AA^{T}||_{op})$ 2509

2510 2511

2512 Lemma J.7. Following condition of Lemma J.5, 2513

$$E||H_j(\tilde{X}W_0)H_j(\tilde{X}W_0)^\top||_{op} \le Cj!$$

 $E(e^{(X-\mu)t}) < e^{\frac{k^2}{2}t^2}$

 $Ee^{Xt} < e^{\frac{k^2}{2}t^2 + \mu t}$

(47)

Proof of Lemma J.7. For non-centered Sub Gaussian random variable X with mean μ , 2517

2514

2515 2516

Firstly, we proof $\mu = 0$ case. For centered Sub Gaussian vector g, let $z = g^{\top}u, z' = g^{\top}v, \rho$ -correlated. s.t. $||u||^2 = g^{\top}v$ $||v||^2 = 1, u^T v = \rho$, then by equation 47 2524

$$\mathbb{E}\exp(sz+tz') \le \exp(\frac{k^2}{2}||u||^2s^2 + k^2 < \vec{u}, \vec{v} > st + \frac{k^2}{2}||v||^2t^2)$$

$$\leq \exp(\frac{\kappa}{2}(s^2 + 2\rho st + t^2))$$

Dividing by $\exp(\frac{k^2}{2}(s^2+t^2))$, then $E[\exp(sz - \frac{k^2}{2}s^2)\exp(tz' - \frac{k^2}{2}t^2)] \le \exp(\rho st) = \sum_{i=0}^{\infty} \frac{\rho^i}{j!} s^j t^j$ Using proof techniques similar to those in Lemma M.1, one can acquire $EH_i(u^Tq)H_k(v^Tq) < j!\rho^j \mathbf{1}_{i-k}$ (48)For $\mu \neq 0$ case, considering non-centered Sub Gaussian Random vector g with mean μ and centered Sub Gaussian Random vector ξ s.t. $g = \xi + \mu$. We use proof techniques similar to those in Theorem M.11. Denote $\nu = \min(j, k)$. Considering $u^{\top}g, v^{\top}g$, $\mathbb{E}[H_i(u^T\mu + u^T\xi)H_k(v^T\mu + v^T\xi)]$ $= \mathbb{E}[\{\sum_{i=1}^{j} {j \choose i} (u^{T} \mu)^{i} H_{j-i}(u^{T} \xi)\} \cdot \{\sum_{i=1}^{k} {k \choose h} (v^{T} \mu)^{h} H_{k-h}(v^{T} \xi)\}]$ $= \mathbb{E}\left[\sum_{q}^{\nu} {\binom{\nu}{q}}^2 (u^T \mu)^{j-q} (v^T \mu)^{k-q} H_q(u^T \xi) H_q(v^T \xi)\right] \text{ by equation 48}$ (49) $\leq \sum_{n=0}^{\nu} {\binom{\nu}{q}}^2 (u^T \mu)^{j-q} (v^T \mu)^{k-q} \cdot \nu! \rho^{\nu}$ $\leq C \min(j, k)!$ Proof of Lemma J.5. Let $A = H_j(\tilde{X}W_0)$, then $\mathbb{P}(||A||_{op} \ge t) \le \mathbb{P}\left(||\frac{1}{n}AA^{T} - EAA^{T}||_{op} \ge \frac{t^{2}}{n} - ||EAA^{T}||_{op}\right) \quad \text{(by Lemma J.6)}$ $\leq \frac{1}{\frac{t^2}{2} - ||EAA^T||_{op}} E\left[||\frac{1}{n}AA^T - EAA^T||_{op} \right] \quad \text{(by Markov's inequality)}$ $\leq \left[\frac{t^2}{n} - E\left[||AA^T||_{op}\right]\right]^{-1} \delta \max\left(\sqrt{||EAA^T||_{op}}, \delta\right) \quad \text{(by Theorem 5.48 in Vershynin (2010))}$ $\leq \left[\frac{t^2}{n} - E\left[||AA^T||_{op}\right]\right]^{-1} \delta \max\left(\sqrt{E\left[||AA^T||_{op}\right]}, \delta\right) \quad \text{(by Jensen's inequality)}.$ Let $M = E \max_i ||H_j(W_0 \tilde{x}_i)||^2$ and $\delta = C \sqrt{\frac{M \log n}{N}}$. Moreover, we note that $\frac{||\tilde{x}_i||^{2j}}{N}$ is sub-weibull random variable and bound of (Kuchibhotla & Chakrabortty, 2022) proposition A.6 can be applied. Use property of $\frac{||\tilde{x}_i||^{2j}}{N}$, W_0 and hermite polynomials, we have $M \le c_j E \max_i ||(W_0 \tilde{x}_i)^{\circ j}||_2^2 \le c_j E \max_i ||x||^{2j} \le c_j N (\log n)^{\frac{1}{2}}.$ Therefore, $\delta \leq C \log n$. Let $t^2 = n \cdot Q_n E ||AA^T||_{op}$ s.t. Q_n is positive and increasing. Building on the result derived

2585 above, we can continue expanding the expression as follows: 2586 $\left[\frac{t^2}{n} - E\left[||AA^T||_{op}\right]\right]^{-1} \delta \max\left(\sqrt{E\left[||AA^T||_{op}\right]}, \delta\right)$ 2587 2588 $\leq [\frac{t^2}{n} - E||AA^T||_{op}]^{-1}C\log n \max(\sqrt{E||AA^T||_{op}}, \log n)$ 2589 2590(50) $= [E||AA^{T}||_{op}(Q_{n}-1)]^{-1}C\log n\max(\sqrt{E||AA^{T}||_{op}},\log n)$ $\leq C \frac{\log n \max(\sqrt{E||AA^T||_{op}}, \log n)}{E||AA^T||_{op}Q_n}$ 2594 2595 2596 Choosing $Q_n = \log^3 n$, and using Lemma J.7, we conclude the proof. 2597 2598 *Fact* J.8. For any vector u, v and any matrix A, B2599 A. $||uv^T||_{op} = ||u||_2 ||v||_2$ 2600 2601 B. $||u||_{\infty} < ||u||_{2} < \sqrt{n} ||u||_{\infty}$ 2602 2603 C. $||u^{\circ k}|| < ||u||^k$ 2604 D. $||u^{\circ k}||_2 \leq \sqrt{n} ||u^{\circ k}||_{\infty} \leq \sqrt{n} \max_i (|u_i^k|) = \sqrt{n} (\max_i |u_i|)^k = \sqrt{n} ||u||_{\infty}^k$ 2605 2606 E. Schur product theorem 2607 $||A \circ B||_{op} = \sup_{||x||=1} tr(A^T \operatorname{diag}(x)B\operatorname{diag}(x)) \le ||A||_{op} \cdot ||B||_{op}$ 2609 2610 Next, let $L = O(\log n)$. 2611 Denote $\sigma_L(z) = \sum_{k=1}^L c_k H_k(z), F^L = \sigma_L(\tilde{X}W)$ and $F_0^L = \sigma_L(\tilde{X}W_0)$. 2612 2613 Then, $F = F^L + (\sigma - \sigma_L)(\tilde{X}W)$. 2614 2615 Using Lemma J.5, w in assumption 2.1, w.p. 1 - o(1)2616 $||E[(\sigma - \sigma_L)(W_0X)(\sigma - \sigma_L)(W_0X)^T]||$ 2617 $\leq C \sum_{k=L+1}^{\infty} k! c_k^2 \leq C \sum_{k=L+1}^{\infty} k^{-3-w} \leq C \int_L^{\infty} k^{-\frac{3}{2}-w} dk \leq C L^{-2-w}.$ 2618 (51)2619 2621 Therefore, following same proof technique as Lemma J.5, J.6, J.7, 2622 2623 $||(\sigma - \sigma_L)(\tilde{X}W_0)||_{op} = o_{\mathbb{P}}(\sqrt{n\log^3 n \cdot L^{-2-w}}) = o_{\mathbb{P}}(\sqrt{n})$ (52)2624 2625 Also, because $||W||_{op} = O(1)$, 2626 $||(\sigma - \sigma_L)(\tilde{X}W)||_{op} = o(\sqrt{n\log^3 n} \cdot L^{-2-w}) = o_{\mathbb{P}}(\sqrt{n})$ (53)2628 2629 Finally, we proof Theorem 3.3. 2630 Proof of Theorem 3.3. We write $F^L + F_0^L = F^L + F_0^L$, then $F^L = F_0^L + \sum_{k=1}^L c_k (H_k(\tilde{X}W) - H_k(\tilde{X}W_0))$. We have to 2631 study $H_k(\tilde{X}W) - H_k(\tilde{X}W_0)$ term. 2633 $H_k(\tilde{X}W) - H_k(\tilde{X}W_0)$ 2634 2635 $=H_k(\tilde{X}W_0^T+\tilde{X}G^T)-H_k(\tilde{X}W_0)$ (54) $= (\tilde{X}G)^{ok} + \sum_{i=1}^{k-1} \binom{k}{j} H_{k-j}(\tilde{X}W_0) \circ (\tilde{X}G)^{oj}$ 2637 2639

2640	Thus,	
2641		
2642	I_{k-1}	
2643	$F^L = F^L + \sum_{n=1}^{L} a_n (\tilde{Y}C)^{\circ k} + \sum_{n=1}^{L} \sum_{n=1}^{n-1} a_n {k \choose k} H_{i-1} (YW_i) \circ (\tilde{Y}C)^{\circ j}$	
2644	$I' = I'_0 + \sum_{k=1}^{\infty} c_k(AG) + \sum_{k=1}^{\infty} \sum_{j=1}^{\infty} c_k(j) II_{k-j}(AW_0) \circ (AG)$	
2645	$k=1$ $k=1$ $j=1$ $\langle 0 \rangle$	
2646	$\frac{L}{\sum k} \frac{1}{\sqrt{2}} \sum c \frac{T}{2} \frac{1}{2} \frac{1}{\sqrt{2}} \frac$	
2647	$=F_0^{\scriptscriptstyle D}+\sum c_1^\kappa c_k (X\sum eta_{ij}a_{ij}^{\scriptscriptstyle I})^{ m ok}$	
2047	$k{=}1$ $i{<}j$	
2048		
2049	$ -\sum c_1^k c_k (X \sum \beta_{ij} a_{ij}^T)^{ok}$	
2650	\wedge $k=1$ $i < j$	
2651	Δ_1 L	
2652	$+\sum c_k (\tilde{X}G)^{ok}$	(55)
2653	$\begin{bmatrix} & & \\ & & \\ & & \\ & & \\ & & \end{bmatrix}$	
2654	$\begin{bmatrix} L & k-1 & \\ & & \end{pmatrix}$	
2655	$+\sum \sum c_k \binom{\kappa}{k} H_{k-i}(\tilde{X}W_0) \circ (\tilde{X}G)^{\circ j}$	
2656	$\left(\sum_{k=1}^{j} \sum_{j=1}^{j} n^{k} \left(j \right) \right)$	
2657	$\Delta_2 \mid L_k = 1$	
2658	$\int \sum_{k=1}^{\infty} \sum_{j=1}^{n-1} c^{j} c_{i} \left(k \right) H_{i-1} \left(\tilde{X} W_{c} \right) \circ \left[\tilde{X} \sum_{k=1}^{\infty} \beta_{i+1} a^{T} \right] \circ j$	
2659	$\sum_{k=1}^{n} \sum_{i=1}^{n} c_1 c_k (j)^{m_{k-j}(\mathcal{M}, \mathcal{V}_0)} \circ [\mathcal{M} \sum_{i \leq i} \rho_{ij} u_{ij}]$	
2660	$L \kappa = 1 j - 1 \qquad i < j$	
2661	$\mathbf{A} = \sum_{k=1}^{L} \sum_{j=1}^{k-1} j_{j} \left(k \right) \mathbf{H} = (\tilde{\mathbf{Y}} \mathbf{H}) = [\tilde{\mathbf{Y}} \sum_{k=1}^{K} \rho_{k} + T_{k} \rho_{k}]$	
2662	$\Delta_3 + \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} c_1^{i} c_k (j)^{H_{k-j}(X,W_0)} \circ [X \sum_{i=1}^{\infty} \beta_{ij} a_{ij}]^{j}$	
2663		
2664		
2665		
2666	$ F_0^L = \Theta(\sqrt{n})$ by Moniri et al. (2024).	
2667		
2668	$\left\ \sum_{k=1}^{L} c_1^{\kappa} c_k (X \sum_{i < j} \beta_{ij} a_{ij}^{I})^{\delta \kappa}\right\ \text{ is bigger than } \sqrt{n}.$	
2669	For $\Delta_1, \Delta_2, \Delta_3$, it is derived as follows	
2670		
2671		
2672		
2673		
2674	$\ \Delta_1\ \le \sum c_k \ (XG)^{o\kappa} - c_1^{\kappa} (X\sum \beta_{ij} a_{ij}^1)^{o\kappa}\ $	
2675	$k{=}1$ $i{<}j$	(56)
2676	$\frac{L}{2}$ or $\log^2 n$	(50)
2677	$\leq C \sum \log^{2k} n \cdot n^{-\frac{\kappa}{2}} = O(\frac{\log n}{\sqrt{2}}) = o(1)$	
2678	$\overline{K=1}$ \sqrt{n}	
2679		
2680		
2681		
2682		
2683	$\frac{L}{k-1}$ $\binom{k}{k}$	
2684	$\ \Delta_2\ \leq \sum \sum c_k \binom{n}{i} \ H_{k-j}(XW_0^T) \circ [(XG^T)^{\circ j} - c_1^j[X\sum \beta_{ij}a_{ij}^T]^{\circ j}]\ $	
2685	$\overline{k=1} \overline{j=1}$ (J) $\overline{i$	
2005	L $k-1$	
2000	$\leq C \sum \sum \ H_{k-i}(\tilde{X}W_0^T)\ \ (\tilde{X}G^T)^{\circ j} - c_1^j (\tilde{X}\sum \beta_{ij}a_{ij}^T)^{\circ j} \ $	
200/	k=1 $j=1$ $i < j$	(57)
2000	L $k-1$	(57)
2089	$\leq C \sum \sum \sqrt{n} \log^{\frac{3}{2}} n \sqrt{j!} \cdot n^{-\frac{j}{2}} \log^{2j} n$	
2090	$ \sum_{k=1}^{-}$ $\sum_{j=1}^{-}$ \cdots $\sum_{j=1}^{-}$ \cdots $\sum_{k=1}^{-}$	
2691	L $k-1$ $ \frac{3}{2}+2i$	
2692	$\leq C \sum \sum \frac{\sqrt{n\sqrt{j! \log^2 n}}}{2} = O(\log^{\frac{7}{2}} n)$	
2693	$ = \sum_{k=1}^{\infty} \sum_{i=1}^{\infty} \sqrt{n^i} $	
2694		

$$\|\Delta_3\| \leq$$

$$\|\Delta_3\| \le C \sum_{k=1}^{L} \sum_{j=1}^{k-1} \|H_{k-j}(\tilde{X}W_0) \circ [\tilde{X}\sum_{i < j} \beta_{ij} a_{ij}^T]^{\circ j}\|$$

$$\leq C \sum_{k=1}^{L} \sum_{j=1}^{k-1} \|\operatorname{diag}(\tilde{X}\beta)^{\circ j}\| \|H_{k-j}(\tilde{X}W_{0})\| \|\operatorname{diag}(a)^{\circ j}\| \\ \leq C \sum_{k=1}^{L} \sum_{j=1}^{k-1} (M_{a}M_{b})^{j} \|H_{k-j}(\tilde{X}W_{0})\| \\ \leq C \sum_{k=1}^{L} \sum_{j=1}^{k-1} n^{-\frac{1}{2}j} \log^{j} n \sqrt{n} \log^{\frac{3}{2}} = O(\log^{\frac{5}{2}} n)$$
(58)

Therefore, we conclude the proof.

K. Proof of Clustering Risk Analysis in two-classes case

Definition K.1. Given N, d, let

$$S_{d,k}^{(1)} = \mathbb{E}_{w \sim Unif(\mathbb{S}^{d-1})}[(w^{T}e_{1})^{k}] \in \mathbb{R}_{+}$$

$$S_{d,k,k'}^{(2)} = E_{w}[(w^{T}\hat{\mu}_{1})^{k}(w^{T}\hat{\mu}_{2})^{k'}]$$

$$\rho_{k,k'}^{(1)} = NS_{d,k+k'}^{(1)}\mathbf{1}_{k+k' \text{ is even}} \in \mathbb{R}_{+}$$

$$\rho_{k,k'}^{(2)}(\cos(\mu_{1},\mu_{2})) = NS_{d,k,k'}^{(2)}\mathbf{1}_{k+k' \text{ is even}} \in \mathbb{R}_{+}$$

$$\rho_{k,k',r}^{(3)} = \frac{c_{1}^{k}S_{d,k'}^{(1)}}{N^{\frac{k}{2}-1}} \binom{k}{r}(r-1)!!(k-1)!!\mathbf{1}_{k,k',r \text{ is even}} \in \mathbb{R}_{+}$$

$$\rho_{k,k',r,r'}^{(4)} = \frac{2c_{1}^{k+k'}S_{d,k}^{(1)}}{N^{\frac{k'}{2}-1}} \binom{k'}{r'}(r'-1)!!(k'-1)!!\mathbf{1}_{k,k',r' \text{ is even}} \in \mathbb{R}_{+}$$
(59)

For $S_{d,k,k'}^{(2)}$, it depends on $\cos(\mu_1,\mu_2)$. As $\cos(\mu_1,\mu_2)$ increases, $S_{d,k,k'}^{(2)}$ grows, while it decreases as $\cos(\mu_1,\mu_2)$ decreases. e.g. when $\mu_1 = \mu_2$, $S_{d,k,k'}^{(2)} = S_{d,k+k'}^{(1)}$, and when $\mu_1 = -\mu_2 = -S_{d,k+k'}^{(1)}$. **Lemma K.2.** Let $C_{d,k} \triangleq \mathbb{E}_{\omega}[(\omega^{\top} e_1)^k]$ s.t. $\omega \sim Unif(\mathbb{S}^{d-1})$, then

$$\mathbb{E}_{\omega}[(\omega^{\top}\mu)^{k}] = \|\mu\|^{k} S_{d,k}^{(1)} \mathbf{1}_{\mathbf{k} \text{ is even}}$$

$$\tag{60}$$

Proof of K.2. The uniform distribution on the sphere is origin-symmetric. Therefore, when k is odd, Expectation is zero. In the other case, also use isotropic property of uniform sphere,

$$E_{\omega}[(\omega^{\top}\mu)^{k}] = \|\mu\|^{k} E_{\omega}[(\omega^{\top}e_{1})^{k}] = \|\mu\|^{k} S_{d,k}^{(1)}$$

In the proof below, we utilize the results of Corollary M.12, Corollary M.13, and Lemma K.2.

Lemma K.3. Given vector $a \in \mathbb{R}^{\mathbf{N}} \beta \in \mathbb{R}^{\mathbf{d}}$ and Gaussian Random vector $x \sim \mathcal{N}(\mu, I)$. Let $b = x^{\top}\beta \sim \mathcal{N}(\mu^{\top}\beta, \|\beta\|^2)$, then

$$\mathbb{E}_x (x^\top \beta a^\top)^{\circ k} = \sum_{r=0}^k \binom{k}{r} (\mu^\top \beta)^{k-r} \|\beta\|^r (r-1)!! \mathbf{1}_{\mathrm{r} \text{ is even}} a^{\circ k^\top}$$
(61)

$$\mathbb{E}_a a^{\circ k} = \frac{(k-1)!! \mathbf{1}_{k \text{ is even}}}{N^{\frac{k}{2}}} \mathbb{1}$$
(62)

2747
2748
$$\mathbb{E}_{a}a^{\circ k^{\top}}a^{\circ k'} = \frac{(k+k'-1)!!\mathbf{1}_{k+k' \text{ is even}}}{N^{\frac{k+k'}{2}}-1}$$
(63)

2750	Proof. This follows directly from Corollary M.12.	
2752	Proof of Proposition 4.5. Let cohesion of initialized feature as	
2753 2754	$coh_0 = \mathbb{E}_{W_0}[\mathbb{E}_{x \sim c_1} F_0^L(x)^T \mathbb{E}_{x' \sim c_1} F_0^L(x')]$	(64)
2755 2756	Let <i>cohesion</i> of feature after training as	
2757 2758	$coh_1 = \mathbb{E}_{W_0,a}[\mathbb{E}_{x \sim c_1} F_L(x)^T \mathbb{E}_{x' \sim c_1} F_L(x')]$	(65)
2759	Calculate coh_0 By Lemma K.2,	
2760 2761 2762	$\sum_{n=1}^{L} \sum_{m=1}^{L} \sum_{m$	
2763	$con_0 = \mathbb{E}_{W_0}[\mathbb{E}_{x \sim c_1}[\sum_{k=1}^{k} c_k H_k(W_0 x)] \mathbb{E}_{x' \sim c_1}[\sum_{k'=1}^{k} c_{k'} H_{k'}(W_0 x)]]$	
2764 2765	$= \sum_{k=1}^{L} c_k c_{k'} \mathbb{E}_{W_0} [\sum_{k=1}^{N} (W_0[q]^T \mu_1)^{k+k'}]$	
2766 2767	$\begin{array}{c} k=1,k'=1 \\ L \end{array} \qquad \qquad$	(66)
2768 2769	$= N \sum_{k=1,k'=1} c_k c_{k'} (\ \mu_1\ ^{k+k'} S_{d,k+k'}^{(1)}) 1_{(k+k')\text{even}}$	
2770	$\sum_{k=1,k=1}^{L} \sum_{k=1,\dots,k=k} (1) _{k=1} _{k+k'}$	
2772	$= \sum_{k=1,k'=1}^{k} c_k c_{k'} \rho_{k,k'} \ \mu\ $	
2773	Calculate coh_1	
2775 2776	$coh_{\star} = \mathbb{E}_{W} [\mathbb{E} [\sum_{k=1}^{L} (c_{k}H_{\star}(W^{T}r) + c_{k}c^{k}(r^{T}\beta a)^{ok}]^{T}\mathbb{E} [\sum_{k=1}^{L} (c_{k}H_{\star}(W^{T}r) + c^{k}(r^{T}\beta a)^{ok}]]$	
2777 2778	$\sum_{k=1}^{2} \sum_{k=1}^{2} \sum_{k=1}^{2} \left[\sum_{k=1}^{2} \left(c_k n_k (w_0 x) + c_k c_1 (x \beta u) - \right) \sum_{k=1}^{2} \sum_{k=1}^{2} \left(c_k n_k (w_0 x) + c_1 (x \beta u) - \right) \right]$	
2779 2780	$= \mathbb{E}_{W_0,a} \left[\sum_{k=1}^{L} c_k c_{k\prime} \left[\mathbb{E}_x H_k (W_0^T x)^T \mathbb{E}_{x\prime} H_{k\prime} (W_0^T x\prime) \right] \right]$	
2781 2782	$+ 2\mathbb{E}_{r}H_{k}(W_{0}^{T}x)^{T}\mathbb{E}_{r'}c_{1}^{k'}(x^{T}\beta a)^{ok'} + c_{1}^{k+k'}\mathbb{E}_{r}(x^{T}\beta a)^{ok'}\mathbb{E}_{r'}(x^{T}\beta a)^{ok'}]$	
2783	$L = \frac{L}{k' \pi} = \frac{1}{\pi} \frac{L}{k' \pi} = \frac{1}{\pi} \frac{L}{k' \pi} \frac{L}{k'$	
2785	$= coh_0 + 2 \sum_{k,k'=1} c_k c_{k'} c_1 \mathbb{E}_{W_0} \mathbb{E}_x H_k(W_0 x) \mathbb{E}_a \mathbb{E}_{x'}(x \beta a)$	
2786 2787	+ $\sum_{k=1}^{L} c_k c_{k\prime} c_1^{k+k\prime} \mathbb{E}_a[\mathbb{E}_x (x^\top \beta a)^{ok^T} \mathbb{E}_{x\prime} (x^{\prime \top} \beta a)^{ok}]$	
2788 2789	k, k'=1	
2790 2791	$= coh_0 + 2N \sum_{\substack{k,k'=1\\ k',k''}} c_k c_{k'} c_{1'}^{k'} (\ \mu_1\ ^k S_{d,k}^{(1)}) (\frac{1}{N^{\frac{k'}{2}}} \sum_{r'=0}^{r'} {\binom{k'}{r'}} (\mu_1^T \beta)^{k'-r'} \ \beta\ ^{r'} (r'-1)!! (k'-1)!! 1_{k,k',r' \text{ is even}}$	
2792 2793	$\sum_{k=0}^{L} c_k c_{k'} c_1^{k+k'} \sum_{k=0}^{k} \sum_{k=0}^{k'} (k) (k') (T_{k'}) (T_{k'}) (k+k'-r-r') (k$	
2794	$+\sum_{k,k'=1} \frac{n}{N^{\frac{k+k'}{2}} - 1} \sum_{r=0} \sum_{r'=0} {\binom{r}{r'}} {\binom{r'}{r'}} {\binom{\mu_1^{r}\beta}{r'+r'}} \ \beta\ ^{r+r'} (r-1)!!(r'-1)!!1_{k+k',r,r' \text{is even}}$	
2796		

2798 Proof of Proposition 4.6. Let separability of initialized feature as

 $sep_{0} = -\mathbb{E}_{W_{0}}[\mathbb{E}_{x \sim c_{1}}F_{0}^{L}(x)^{T}\mathbb{E}_{x' \sim c_{2}}F_{0}^{L}(x')]$ (67)

2802 Let *separability* of feature after training as

$$sep_1 = -\mathbb{E}_{W_0,a}[\mathbb{E}_{x \sim c_1} F_L(x)^T \mathbb{E}_{x' \sim c_2} F_L(x')]$$
(68)

2805 Calculate sep_0 By Lemma K.2,

$$sep_{0} = -\sum_{k=1,k'=1}^{L} c_{k}c_{k'}\mathbb{E}_{W_{0}}\left[\sum_{q=1}^{N} (W_{0}[q]^{T}\mu_{1})^{k}(W_{0}[q]^{T}\mu_{2})^{k'}\right]$$

$$= -N\sum_{k=1,k'=1}^{L} c_{k}c_{k'}\mathbb{E}_{w\sim Unif(\mathbb{S}^{d-1})}\left[(w^{T}\mu_{1})^{k}(w^{T}\mu_{2})^{k'}\right]$$

$$= -N\sum_{k=1,k'=1}^{L} c_{k}c_{k'}\|\mu_{1}\|^{k}\|\mu_{2}\|^{k'}E_{w}\left[(w^{T}\hat{\mu_{1}})^{k}(w^{T}\hat{\mu_{2}})^{k'}\right]$$

$$= -N\sum_{k=1,k'=1}^{L} c_{k}c_{k'}\|\mu_{1}\|^{k}\|\mu_{2}\|^{k'}S_{d,k,k'}^{(2)}\mathbf{1}_{k+k'\text{ is even}}$$

$$= -\sum_{k=1,k'=1}^{L} c_{k}c_{k'}\|\mu_{1}\|^{k}\|\mu_{2}\|^{k'}\rho_{k,k'}^{(1)}$$
(69)

2823 **Calculate** *sep*₁

$$\begin{aligned} & 2824 \\ 2825 \\ 2826 \\ 2827 \\ 2828 \\ 2829 \\ 2829 \\ 2829 \\ 2830 \\ 2830 \\ 2830 \\ 2831 \\ 2832 \\ 2833 \\ 2834 \\ 2835 \\ 2836 \\ 2836 \\ 2837 \\ 2836 \\ 2837 \\ 2838 \\ 2839 \\ 2840 \\ 2840 \end{aligned} = sep_0 - \sum_{k,k\prime=1}^{L} c_k c_{k\prime} \begin{bmatrix} \mathbb{E}_{x \sim c_1} H_k (W_0^T x)^T \mathbb{E}_{x\prime \sim c_2} C_1^{k\prime} (x^T \beta a)^{ok\prime} \\ + \mathbb{E}_{x \sim c_1} c_1^k (x^T \beta a)^{ok'} \mathbb{E}_{x\prime \sim c_2} (x^T \beta a)^{ok\prime} \\ + \mathbb{E}_{x \sim c_1} c_1^k (x^T \beta a)^{ok'} \mathbb{E}_{x\prime \sim c_2} (x^T \beta a)^{ok\prime} \end{bmatrix} \\ = sep_0 - \sum_{k,k\prime=1}^{L} c_k c_{k\prime} \begin{bmatrix} c_1^{k\prime} (\|\mu_1\|^k S_{d,k}^{(1)}) \frac{1}{N^{\frac{k\prime}{2}-1}} \sum_{r\prime=0}^{k\prime} {k\prime} (k') (\mu_2^T \beta)^{k\prime-r\prime} \|\beta\|^{r\prime} (r\prime - 1)!! (k\prime - 1)!! \mathbf{1}_{k,k\prime,r\prime} \text{ is even} \\ + c_1^k (\|\mu_2\|^{k\prime} S_{d,k\prime}^{(1)}) \frac{1}{N^{\frac{k}{2}-1}} \sum_{r=0}^{k} {k\prime} (k') (\mu_1^T \beta)^{k-r} \|\beta\|^{r\prime} (r-1)!! (k-1)!! \mathbf{1}_{k,r,k\prime} \text{ is even} \\ + c_1^{k+k\prime} \sum_{r=0}^{k} \sum_{r\prime=0}^{k} {k\prime} (k') (\mu_1^T \beta)^{k-r} (\mu_2^T \beta)^{k\prime-r\prime} \|\beta\|^{r+r\prime} (r-1)!! (k-1)!! \mathbf{1}_{k,r,k\prime} \text{ is even} \\ + c_1^{k+k\prime} \sum_{r=0}^{k} \sum_{r\prime=0}^{k} {k\prime} (k') (\mu_1^T \beta)^{k-r} (\mu_2^T \beta)^{k\prime-r\prime} \|\beta\|^{r+r\prime} (r-1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ + c_1^{k+k\prime} \sum_{r=0}^{k} \sum_{r\prime=0}^{k} {k\prime} (k') (\mu_1^T \beta)^{k-r} (\mu_2^T \beta)^{k\prime-r\prime} \|\beta\|^{r+r\prime} (r-1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)! \mathbf{1}_{k+k\prime,r,r\prime} \text{ is even} \\ \frac{1}{N^{\frac{k+k\prime}-1}} (k+k\prime - 1)!$$

28442845L. Additional Lemmas of Sub-Gaussian Distribution

For more detailed explanation and well known results of Sub-Gaussian we used, please refer to Vershynin (2010; 2018).
We show below that the truncated Gaussian distribution, utilized in our synthetic data experiments, is a sub-Gaussian distribution.

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

Proof. Denote $\mathcal{N}_{(a,b)}(0,\sigma^2)$ is Truncated Gaussian distribution which have support on (a,b) s.t. $a, b \in (-\infty,\infty)$. 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.

28552856 L.1. Generalization of centered Sub-Gaussian results toward non-centered

We verify below that the results on centered sub-Gaussian distributions from Vershynin (2018) can be extended to non-centered sub-Gaussian distributions.

Lemma L.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.

$$\begin{aligned} \|\sum X_{i}\|_{\psi_{2}} &\leq \|\sum (X_{i} - \mathbb{E}X_{i})\|_{\psi_{2}} + \|\sum \mathbb{E}X_{i}\|_{\psi_{2}} \\ &\leq \|\sum (X_{i} - \mathbb{E}X_{i})\|_{\psi_{2}} + \sum \|\mathbb{E}X_{i}\|_{\psi_{2}} \\ &\leq \|\sum (X_{i} - \mathbb{E}X_{i})\|_{\psi_{2}} + C \sum \|X_{i}\|_{\psi_{2}} < \infty \end{aligned}$$
(70)

Lemma L.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, $||A|| \le CK(\sqrt{m + n} + t) \text{ w.p. } 1 - \exp(-t^2)$
(71)

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

Lemma L.4. (Expectation of operator norm for non-centered Sub-Gaussian matrix generalization of 4.4.6 in Vershynin (2018)) $\mathbb{E}[|A|| < OV(\sqrt{m})$

$$\mathbb{E}||A|| \le CK(\sqrt{m} + \sqrt{n})$$

Alternatively, $\mathbb{E}||A|| \le CK(\sqrt{m+n})$, and, $\mathbb{E}||A||^2 \le C(m+n)$ (72)

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

M. Additional Results of Expectation of Hermite Polynomials

The non standard gaussian expectation of the product of two Hermite polynomials is computed as follows. It is an generalization of results of standard Gaussian distributions in O'Donnell (2021); Moniri et al. (2024) into a generalized multivariate Gaussian. We start with previously known facts, and derive our generalized results. These findings provide a useful analysis tool for Hermite polynomials, and may offer a foundation for broader applications in future works involving nonlinear activations decomposable into Hermite polynomials under the assumption of a multivariate Gaussian distribution.

M.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 orthogonality of Hermite polynomials.

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

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 $\mathbb{E}[Z_1Z_2] = \rho$. Then for any $k_1, k_2 \in \{0, 1, \cdots, \}$

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

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

$$\mathbb{E}_{Z \sim \mathcal{H}(0,1)}[H_j(W^\top Z)H_k(W^\top Z)^\top] = k!(W^\top W)^{\circ j}\mathbf{1}_{\mathbf{j}=\mathbf{k}}$$
(73)

 $\mathbb{E}_{Z \sim n(0,1)}[H_j(W^{\top}Z)^{\top}H_k(W^{\top}Z)] = k! \sum ||W[i]||^{2j} \mathbf{1}_{j=k} = k! N \mathbf{1}_{j=k}$ (74)

2915 *Proof.* We apply H_j element-wise. By Lemma M.1, we can acquire the above result.

2916

2926

2929

2934

2937

2938

The following remark presents a modified condition of Lemma M.1 for the case where $a, b \notin S^{d-1}$ in Lemma M.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 M.1.

2921 *Remark* M.3 (the modified condition of Lemma M.1). For $d \in \mathbb{N}$, $u, v \in \mathbb{R}^d$, $Z \sim \mathcal{N}(0, I_d)$,

$$Z_{2923} \quad Z_1 = \langle u, Z \rangle \sim \mathcal{N}(0, ||u||_2^2), Z_2 = \langle v, Z \rangle \sim \mathcal{N}(0, ||v||_2^2).$$

²⁹²⁴ Then, Z_1, Z_2 is $\rho \triangleq \langle \frac{u}{||u||}, \frac{v}{||v||} \rangle$ - correlated

$$corr(Z_1, Z_2) = \frac{\mathbb{E}[Z_1 Z_2]}{\sqrt{V(Z_1)}\sqrt{V(Z_2)}} = \frac{\mathbb{E}_Z \langle u, Z \rangle \langle v, Z \rangle}{||u|| ||v||}$$
$$= \frac{\mathbb{E}_g \sum_i \sum_j u_i v_j Z_i Z_j}{||u|| ||v||} = \frac{\sum_i \sum_j u_i v_j \mathbb{E}_Z[Z_i Z_j]}{||u|| ||v||}$$
$$= \frac{\langle u, v \rangle}{||u|| ||v||}$$
(75)

2935 Additionally,

$$\begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} ||u||^2 & \langle u, v \rangle \\ \langle v, u \rangle & ||v||^2 \end{pmatrix} \right)$$
(76)

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 M.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.

$$\mathbb{E}[X_A] = \sum_{\sigma \in \prod(A)} \prod_{(i,j) \in \sigma} \mathbb{E}[X_{\alpha_i} X_{\alpha_j}] \mathbf{1}_{\mathrm{d} \text{ is even}}.$$
(77)

Now, we generalize the assumptions from the previous works so that Lemma M.1 holds for arbitrary vectors as Remark M.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 M.5 (Generalization of Lemma M.1 for centered Gaussian distribution). For $d \in \mathbb{N}$, $u, v \in \mathbb{R}^d$, $g \sim \mathcal{N}(0, I_d)$, $2955 \quad \langle u, g \rangle \sim \mathcal{N}(0, ||u||_2^2), \langle v, g \rangle \sim \mathcal{N}(0, ||v||_2^2).$

2956

2959

2960

2961 2962 2963 $\mathbb{E}_{g}[H_{j}(u^{\top}g)H_{k}(v^{\top}g)] = \frac{j!\langle u, v \rangle^{j}}{||u||^{2}||v||^{2}}\mathbf{1}_{j=k} - \frac{(||u||^{2}-1)(||v||^{2}-1)}{||u||^{2}||v||^{2}}\mathbb{E}_{g}[(v^{\top}g)^{k}(u^{\top}g)^{j}] + \frac{(||v||^{2}-1)}{||v||^{2}}\mathbb{E}_{g}[H_{j}(u^{\top}g)(v^{\top}g)^{k}] + \frac{(||u||^{2}-1)}{||u||^{2}}\mathbb{E}_{g}[H_{k}(v^{\top}g)(u^{\top}g)^{j}]$ (78)

Remark M.6. The same results can be derived as in Lemma M.1 when the variance is 1 in Thm. M.5.
2965
2966

2967 Proof of Theorem M.5. (Generalize Chapter 11.2 O'Donnell (2021)'s derivation to non unit variance)

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

2970 First, we study about $\mathbb{E}_{g \sim \mathcal{N}(0,\sigma^2)}[e^{tg}]$ in order to analysis non unit variance case.

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

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

$$= e^{\frac{1}{2}t^2}$$
(79)

2979 $\mathbb{E}_{Z,Z'}[\exp(sZ+tZ')]$ study

2980 Studying $\mathbb{E}_{Z,Z'}[\exp(sZ + tZ')]$, we can derive what we need to show.

$$\mathbb{E}_{Z,Z'}[\exp(sZ + tZ')] = \mathbb{E}_{g \sim n(0,I)}[\exp(s\langle u, g \rangle) + \exp(t\langle v, g \rangle)] \\ = \prod_{i} \mathbb{E}_{g \sim n(0,1)}[\exp((su_{i} + tv_{i})g_{i})] \qquad \text{Use equation 79} \\ = \prod_{i} \exp(\frac{1}{2}(su_{i} + tv_{i})^{2}) = \prod_{i} \exp(\frac{1}{2}s^{2}||u||^{2} + \langle u, v \rangle st + \frac{1}{2}t^{2}||v||^{2})$$
(80)

2989 Therefore,

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

Fact M.7. One can verify below propositions with simple calculations.

Let $P_j(z) + z^j = H_j(z), C_u = ||u||^2 - 1, a > 0.$ 2994 Let $f(s) = \exp(sz - \frac{1}{2}s^2), \bar{f}(s) = \exp(sz - \frac{1}{2}as^2)$, then

2995
2996 A. By Taylor expansion,
$$\exp(\langle u, v \rangle st) = \sum_{j=0}^{\infty} \frac{1}{j!} \langle u, v \rangle^j s^j t^j$$
.

2997
2998 B. By Taylor expansion,
$$\bar{f}(s) = \sum_{j=0}^{\infty} \frac{1}{j!} \bar{f}^{(n)}(0) s^j$$

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

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.

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

$$(81)$$

3011 Therefore,

$$\mathbb{E}_{g}[H_{j}(u^{\top}g)H_{j}(v^{\top}g)] = \frac{j!\langle u, v \rangle^{j}}{||u||^{2}||v||^{2}}\mathbf{1}_{j=k} - \frac{(||u||^{2}-1)(||v||^{2}-1)}{||u||^{2}||v||^{2}}\mathbb{E}_{g}[(v^{\top}g)^{j}(u^{\top}g)^{j}] + \frac{(||v||^{2}-1)}{||v||^{2}}\mathbb{E}_{g}[H_{j}(u^{\top}g)(v^{\top}g)^{j}] + \frac{(||u||^{2}-1)}{||u||^{2}}\mathbb{E}_{g}[H_{j}(v^{\top}g)(u^{\top}g)^{j}]$$
(82)

3019 Note that the result of Lemma M.8 can be applied for concrete calculation, and conclude the proof. 3020 Lemma M.8. For $d \in \mathbb{N}$ $u, v \in \mathbb{R}^d$, $a \in \mathcal{H}(0, L)$, $\bar{Z}_1 = \langle u, a \rangle$, $\bar{Z}_2 = \langle u, a \rangle$

Lemma M.8. For $d \in \mathbb{N}$, $u, v \in \mathbb{R}^d$, $g \sim \mathcal{N}(0, I_d)$, $\overline{Z}_1 = \langle u, g \rangle$, $\overline{Z}_2 = \langle v, g \rangle$.

$$\begin{cases} 3022\\ 3023\\ 3024 \end{cases} \sim \mathcal{N}\left(\begin{pmatrix} 0\\ 0\\ Z_2 \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} 0\\ 0\\ 0 \end{pmatrix}, \begin{pmatrix} ||u||^2 & \langle u, v \rangle\\ \langle v, u \rangle & ||v||^2 \end{pmatrix}\right)$$
(83)

 X_{α_i} is defined at Thm. M.4

$$\begin{aligned} & 3028 \\ & 3029 \\ & 3030 \\ & 3030 \\ & 3031 \\ & 3032 \end{aligned} \qquad \mathbb{E}_{\bar{Z}_1, \bar{Z}_2}[\bar{Z}_1^j \bar{Z}_2^k] = j! \sum_{m=0}^{\lfloor \frac{j}{2} \rfloor} \frac{(-1)^m}{m!(j-2m)!2^m} \sum_{\sigma \in \prod(\{\{\bar{Z}_1\} \times j-2m\} \cup \{\{\bar{Z}_2\} \times k\}\})} \prod_{(p,q) \in \sigma} \mathbb{E}[X_{\alpha_p} X_{\alpha_q}] \mathbf{1}_{j+k-2m \text{ is even}} \\ & (84) \\ & \mathbb{E}_{\bar{Z}_1, \bar{Z}_2}[\bar{Z}_1^j \bar{Z}_2^k] = \sum_{m=0}^{\lfloor \frac{j}{2} \rfloor} \prod_{m=0}^{\lfloor \frac{j}{2} \rfloor} \mathbb{E}[X_{\alpha_p} X_{\alpha_q}] \mathbf{1}_{j+k \text{ is even}} \end{aligned}$$

Proof. By explicit formula of Hermite polynomials

$$\mathbb{E}_{\bar{Z}_1,\bar{Z}_2}[H_j(\bar{Z}_1)(\bar{Z}_2)^k] = j! \sum_{m=0}^{\lfloor \frac{j}{2} \rfloor} \frac{(-1)^m}{m!(j-2m)!2^m} \mathbb{E}_{\bar{Z}_1,\bar{Z}_2}[\bar{Z}_1^{j-2m}\bar{Z}_2^k]$$
(85)

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 Gaussian, so we can apply Thm. M.4 with $A = \{\{\bar{Z}_1\} \times p\} \cup \{\{\bar{Z}_2\} \times q\}\}, \mathbb{E}[\bar{Z}_1^p\bar{Z}_2^q] = \sum_{\sigma \in \prod(A)} \prod_{(\tau,\upsilon) \in \sigma} \mathbb{E}[X_{\alpha_\tau}X_{\alpha_\upsilon}].\mathbf{1}_{p+q \text{ is even}}$

Corollary M.9 (Corollary of Lemma M.8). Remark $Z_1 \sim \mathcal{N}(0, ||u||^2)$ For the case k = 0,

 $\sigma \in \prod \left(\left\{ \left\{ \bar{Z}_1 \right\} \times j \right\} \cup \left\{ \left\{ \bar{Z}_2 \right\} \times k \right\} \right\} \right) (p,q) \in \sigma$

$$\mathbb{E}_{\bar{Z}_1}[\bar{Z}_1^j] = \|u\|^j (j-1)!! \mathbf{1}_{j \text{ is even}}$$
(86)

3050 Proof.

$$\mathbb{E}_{\bar{Z}_{1},\bar{Z}_{2}}[\bar{Z}_{1}^{j}\bar{Z}_{2}^{k}] = \mathbb{E}_{\bar{Z}_{1}}[\bar{Z}_{1}^{j}] = \sum_{\sigma \in \prod(\{\bar{Z}_{1}\} \times j\})} \prod_{(p,q) \in \sigma} \mathbb{E}[X_{\alpha_{p}}X_{\alpha_{q}}]\mathbf{1}_{j \text{ is even}}$$

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

$$(87)$$

M.2. Expectation of a product of two Hermite polynomials—Generalization toward non-centered Gaussian

3060 We will change Theorem M.5 and Lemma M.8 to adopt a generalized Gaussian assumption with a mean of zero.

Lemma M.10 (Taylor expansion of Hermite polynomials from Lemma C.2 Moniri et al. (2024)). For any $k_1, k_2 \in \{0, 1, \dots, \}$ and $x, y \in \mathbb{R}$,

$$H_k(x+y) = \sum_{j=0}^k \binom{k}{j} x^j H_{k-j}(y).$$
(88)

Theorem M.11 (Generalization of Thm. M.5 for any Gaussian distribution). For $d \in \mathbb{N}$, $u, v \in \mathbb{R}^d$, $\xi \sim \mathcal{N}(0, 1)$, 3068 $g \sim \mathcal{N}(\mu, \Sigma)$, $Z_1 = \langle u, g \rangle \sim \mathcal{N}(\mu^\top u, u^\top \Sigma u)$, $Z_2 = \langle v, g \rangle \sim \mathcal{N}(\mu^\top v, v^\top \Sigma v)$.

Proof of Theorem M.11. 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 M.10,

$$H_j(\sqrt{u^{\top}\Sigma u}\xi + u^{\top}\mu) = \sum_{\alpha=0}^j \binom{j}{\alpha} (u^{\top}\mu)^{\alpha} H_{j-\alpha}(\sqrt{\mu^{\top}\Sigma u}\xi).$$
(90)

$$\mathbb{E}_{g}[H_{j}(u^{\top}g)H_{k}(v^{\top}g)] = \mathbb{E}_{\xi}[H_{j}(\sqrt{u^{\top}\Sigma u}\xi + u^{\top}\mu)H_{k}(\sqrt{v^{\top}\Sigma v}\xi + v^{\top}\mu)]$$

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

$$= \sum_{\alpha=0}^{j}\sum_{\beta=0}^{k} \binom{j}{\alpha}\binom{k}{\beta}(u^{\top}\mu)^{\alpha}(v^{\top}\mu)^{\beta}\mathbb{E}_{\xi}[H_{j-\alpha}(\sqrt{\mu^{\top}\Sigma u}\xi)H_{k-\beta}(\sqrt{\mu^{\top}\Sigma v}\xi)]$$
(91)

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

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

In summary,

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 M.9. It corresponds to the case k = 0 in Lemma M.8.

Corollary M.12 (Corollary of Lemma M.8). Given $\omega \in \mathbb{R}^d$, let Gaussian Random Variable $Z \sim \mathcal{N}(\mu^\top \omega, \|\omega\|^2)$, then

$$\mathbb{E}_{Z}(Z)^{k} = \sum_{t=0}^{k} {k \choose t} (\mu^{\top} \omega)^{k-t} \mathbb{E}_{\bar{Z} \sim \mathcal{N}(0, \|\omega\|^{2})}[\bar{Z}^{t}]$$

$$= \sum_{t=0}^{k} {k \choose t} (\mu^{\top} \omega)^{k-t} (t-1)!! \cdot \|\omega\|^{t} \mathbf{1}_{\text{t is even}}.$$
(94)

The following corollary, which computes the Gaussian expectation of Hermite polynomials, is derived from the explicit form of Hermite polynomials and Corollary M.9. It corresponds to the case k = 0 in Theorem M.11.

3135	Corollary M.13.	Given $\omega \in \mathbb{S}^{d-1}$, let Gaussian Random Variable $Z \sim \mathcal{N}(\mu^{\top}\omega, 1)$, then	
3136			

$$\mathbb{E}_{x}[H_{k}(\omega^{\top}x)] = \mathbb{E}_{\xi \sim \mathcal{H}(0,1)}[H_{k}(\omega^{\top}\mu + \xi)]$$
$$= \sum_{i=0}^{k} \binom{k}{j} (\omega^{\top}\mu)^{\circ j} E[H_{k}(\xi)H_{0}(\xi)] = (\omega^{\top}\mu)^{k}$$
(95)

N. Information of ImageNet subset used in Experiments 3142

3144 3145	Table 5: Configuration of Expr. V					Table 6: Configuration of Expr. VI				
3146		Vehicle	Bird	Product	Clothing		Step 0	Step 1	Step 2	Step 3
3147	D	98	100	11316	3985	Vehicle	25	50	75	98
3148	D+I(Sub)	138	159	11568	4031	Bird	25	50	75	100
3149 3150	D+I	1098	1100	12316	4985	Clothing	2829 996	5658 1992	8487 2989	3985

In this section, we present the criteria used to select classes for constructing the ImageNet subsets. We manually verified 3152 3153 the label information to select the classes. The ImageNet subsets corresponding to the base fine-grained datasets were constructed as follows: I(V), I(B), I(P), and I(C), representing the Vehicle, Bird, Product, and Clothing subsets, respectively. 3154 These subsets consist of 59, 40, 353, and 46 classes, respectively. To balance the number of samples per class with those in 3155 the base fine-grained datasets, we extracted 82, 58, 5, and 6 samples per class for I(V), I(B), I(P), and I(C), respectively. 3156

3158 N.1. I(V): The Vehicle classes chosen in ImageNet

3159 Total 40 classes. 3160

3143

3151

3157

3165

3161 ambulance, cab, convertible, fire engine, forklift, freight car, garbage truck, go-kart, golfcart, half track, harvester, horse cart, jeep, jinrikisha, limousine, minibus, minivan, Model T, moped, motor scooter, mountain bike, moving van, oxcart, passenger 3163 car, pickup, police van, racer, recreational vehicle, school bus, snowmobile, snowplow, sports car, streetcar, tank, tow truck, 3164 tractor, trailer truck, tricycle, trolleybus, unicycle

3166 N.2. I(B): The bird classes chosen in ImageNet

3167 Total 59 classes. 3168

3169 cock, hen, ostrich, brambling, goldfinch, house finch, junco, indigo bunting, robin, bulbul, jay, magpie, chickadee, water 3170 ouzel, bald eagle, vulture, great grey owl, black grouse, ptarmigan, ruffed grouse, prairie chicken, peacock, quail, partridge, 3171 African grey, macaw, sulphur-crested cockatoo, lorikeet, coucal, bee eater, hornbill, hummingbird, jacamar, toucan, drake, 3172 red-breasted merganser, goose, black swan, tusker, white stork, black stork, spoonbill, flamingo, little blue heron, American 3173 egret, bittern, crane, limpkin, European gallinule, American coot, bustard, ruddy turnstone, red-backed sandpiper, redshank, 3174 dowitcher, oystercatcher, pelican, king penguin, albatross 3175

3176 N.3. I(P): The Product classes chosen in ImageNet

Total 353 classes. 3178

3179 abacus, accordion, acoustic guitar, altar, analog clock, apiary, ashcan, assault rifle, backpack, balance beam, balloon, 3180 ballpoint, Band Aid, banjo, barbell, barber chair, barometer, barrel, barrow, baseball, basketball, bassinet, bassoon, bathing 3181 cap, bath towel, bathtub, beach wagon, beacon, beaker, bearskin, beer bottle, beer glass, bell cote, bib, bicycle-built-for-two, 3182 binder, binoculars, bobsled, bolo tie, bonnet, bookcase, bottlecap, bow tie, brass, breakwater, broom, bucket, buckle, 3183 bulletproof vest, caldron, candle, cannon, canoe, can opener, car mirror, carousel, carpenter's kit, carton, car wheel, cash 3184 machine, cassette, cassette player, CD player, cello, cellular telephone, chain, chain saw, chest, chiffonier, chime, china 3185 cabinet, cleaver, clog, cocktail shaker, coffee mug, coffeepot, coil, combination lock, computer keyboard, confectionery, 3186 corkscrew, cornet, cradle, crash helmet, crate, crib, Crock Pot, croquet ball, crutch, dam, desk, desktop computer, dial 3187 telephone, digital clock, digital watch, dining table, dishrag, dishwasher, disk brake, dogsled, doormat, drum, drumstick, 3188 dumbbell, Dutch oven, electric fan, electric guitar, electric locomotive, envelope, espresso maker, face powder, feather boa, 3189

3190 file, fire screen, flagpole, flute, folding chair, football helmet, fountain pen, four-poster, French horn, frying pan, gasmask, gas pump, goblet, golf ball, gondola, gong, grand piano, grille, guillotine, hair slide, hair spray, hammer, hamper, hand blower, 3192 hand-held computer, handkerchief, hard disc, harmonica, harp, hatchet, holster, honeycomb, hook, horizontal bar, hourglass, 3193 iPod, iron, jack-o'-lantern, jigsaw puzzle, joystick, knot, ladle, lampshade, laptop, lawn mower, lens cap, letter opener, 3194 lighter, lipstick, lotion, loudspeaker, loupe, magnetic compass, mailbox, maraca, marimba, matchstick, maypole, measuring 3195 cup, medicine chest, microphone, microwave, milk can, mixing bowl, modem, monitor, mountain tent, mousetrap, muzzle, 3196 nail, neck brace, necklace, nipple, notebook, oboe, ocarina, odometer, oil filter, organ, oscilloscope, oxygen mask, packet, 3197 paddle, paddlewheel, padlock, paintbrush, paper towel, parachute, parallel bars, park bench, parking meter, pay-phone, 3198 pedestal, pencil box, pencil sharpener, perfume, Petri dish, photocopier, pick, picket fence, piggy bank, pill bottle, pillow, 3199 ping-pong ball, plastic bag, plate rack, plow, plunger, Polaroid camera, pole, pool table, pop bottle, pot, potter's wheel, 3200 power drill, prayer rug, printer, prison, projectile, projector, puck, punching bag, purse, quill, quilt, racket, radiator, radio, radio telescope, rain barrel, reel, reflex camera, refrigerator, remote control, revolver, rifle, rocking chair, rotisserie, rubber eraser, rugby ball, rule, safe, safety pin, saltshaker, sax, scabbard, scale, scoreboard, screen, screw, screwdriver, seat belt, sewing machine, shield, shopping basket, shopping cart, shovel, shower cap, shower curtain, ski, sleeping bag, sliding door, 3204 slot, snorkel, soap dispenser, soccer ball, sock, solar dish, soup bowl, space bar, space heater, spatula, spider web, spindle, 3205 spotlight, steel drum, stethoscope, stole, stopwatch, stove, strainer, stretcher, studio couch, sunscreen, swab, switch, syringe, 3206 table lamp, tape player, teapot, teddy, television, tennis ball, theater curtain, thimble, thresher, throne, tile roof, toaster, tobacco shop, toilet seat, torch, totem pole, tray, tripod, trombone, tub, turnstile, typewriter keyboard, umbrella, vacuum, 3208 vase, vault, velvet, vending machine, violin, volleyball, waffle iron, wall clock, wallet, wardrobe, washbasin, washer, water 3209 bottle, water jug, water tower, whiskey jug, whistle, window screen, window shade, wine bottle, wing, wok, wooden spoon, 3210 comic book, crossword puzzle, street sign, traffic light, book jacket, menu, plate 3211

3212 N.4. I(C): The Clothing classes chosen in ImageNet

3213 Total 46 classes.

abaya, academic gown, apron, bikini, brassiere, breastplate, cardigan, chain mail, Christmas stocking, cloak, cowboy boot,
cowboy hat, cuirass, diaper, fur coat, gown, hoopskirt, jean, jersey, kimono, knee pad, lab coat, Loafer, mailbag, mask,
military uniform, miniskirt, mitten, overskirt, pajama, poncho, running shoe, sandal, sarong, ski mask, sombrero, suit,
sunglass, sunglasses, sweatshirt, swimming trunks, trench coat, vestment, wig, Windsor tie, wool

O. 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 = 300random 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.

3227 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. Please refer Algorithm 3

59

3228

3245	
3246	
3247	
3248	
3249	
3250	
3251	
32.52	
3253	
3254	
3255	
3256	
3257	
3258	
3259	
3260	
3261	
3262	
3263	Algorithm 3 Gaussian-Sampled Random Rotation Matrix Generation
3264	Insuite Number of dimensions of number of motions le
3265	Input: Number of dimensions n , number of matrices κ
3266	Litital source list O
3267	Set of C
3268	Set $\epsilon \leftarrow 0.5$
3260	$\mathbf{IOF} \ i \leftarrow \mathbf{U} \ \mathbf{IO} \ k - \mathbf{I} \ \mathbf{IO}$
3209	If $i \mod \left(\frac{\pi}{16}\right) = 0$ and $i \neq 0$ then
2271	$\epsilon \leftarrow \epsilon \times 0.22360679775$
2272	end if $(0, 1)$ is $(0, 1)$ if $(0, 1)$ if $(0, 1)$ is $(0, 1)$ if $(0, 1)$ if $(0, 1)$ is $(0, 1)$ if $(0, 1)$ if $(0, 1)$ is $(0, 1)$ if $(0, 1)$ i
2072	Generate random matrix $M: M \sim \mathcal{H}(0, 1)^{n \times n}$
3273	Compute perturbed matrix: $A \leftarrow I_n + \epsilon \times M$
3274	Compute QR decomposition: $Q, R \leftarrow QR(A)$
3213	if $det(Q) < 0$ then
3270	Flip first column of $Q: Q[:, 0] \leftarrow -Q[:, 0]$
3211	end if
3278	Add Q to Q
3279	end for
3280	return 2
3281	
3282	
2202	
3284	
3283	
3280	
3287	
3288	
3289	
3290	
3291	
3292	
3293	
3294	
3295	
3296	
3297	
3298	
3299	