

A. Related Work

Theory of CL. While there has been much progress in theoretically understanding CL, most prior work (Wang & Isola, 2020; Graf et al., 2021; Lee et al., 2021; Tosh et al., 2021a;b; Arora et al., 2019b; Tsai et al., 2020; HaoChen et al., 2021) are focused on understanding how CL clusters examples using semantically meaningful information or providing generalization guarantees on downstream tasks. Feature learning has only been studied by (Wen & Li, 2021; Ji et al., 2021) which show that CL learns semantically meaningful features from the data. In contrast, we show that CL may not learn *all* semantically relevant features. Other important recent work (Saunshi et al., 2022; HaoChen & Ma, 2022) studied the role of inductive bias of the function class in the success of CL. Our analysis, however, is focused on understanding failure modes of CL i.e. class collapse and feature suppression.

Class Collapse in Supervised CL. Chen et al. (2022) empirically demonstrates *class collapse* on test data, but does not offer any rigorous theoretical explanation. Graf et al. (2021) proves that optimizing the supervised contrastive loss leads to class-collapsed training set representations. However, we show that there exist many minimizers with such class-collapsed training set representations and not all of them suffer from class collapse at *test time*. We also present the first theoretical characterization of class collapse at test time.

Feature Suppression in Unsupervised CL. Feature suppression has been empirically observed by Tian et al. (2020); Chen et al. (2021); Robinson et al. (2021) but we lack a theoretical formulation of this phenomenon. Li et al. (2023) shows that InfoNCE has local minimums that exhibit feature suppression, thus attributing this phenomenon to failure of optimizing the loss. However, Robinson et al. (2021) shows that the InfoNCE loss can be minimized by many models, some of which learn all task-relevant features, while others do not. We put forth the only theoretical characterization of feature suppression and consequently, use this understanding to suggest practical solutions to remedy this problem.

Joint Supervised and Unsupervised Contrastive Loss. Recently, several versions of loss functions that combine supervised and unsupervised contrastive losses have been proposed. For example, Chen et al. (2022) put forth a weighted sum of supervised CL loss and class-conditional InfoNCE (which has similar effect as \mathcal{L}_{UCL} in our setting) to avoid class collapse. Islam et al. (2021) empirically observed that the joint objective of supervised and unsupervised contrastive loss leads to better transferability of the learned models than their supervised counterparts. We provide the first theoretically rigorous analysis of which features the *minimum norm* global minimizer of the joint loss learns, provably demonstrating that it can avoid class collapse and feature suppression. To the best of our knowledge, this is the only theoretical result that can be used to understand the empirical success of joint losses.

B. Problem Formulation

B.1. Data distribution

We define data distribution $\mathcal{D}_{\text{orig}}$ below. Each example $(\mathbf{x}, y, y_{\text{sub}}) \in \mathcal{D}_{\text{orig}}$ is generated as follows:

$$\mathbf{x} = \mathbf{u} + \boldsymbol{\xi}, \quad \text{where}$$

$$\mathbf{u} = (y\phi_1 + \mu_1)\mathbf{v}_1 + (y_{\text{sub}}\phi_2 + \mu_2)\mathbf{v}_2 + (\rho_k\phi_k + \mu_k)\mathbf{v}_k,$$

and k is uniformly selected from $3, \dots, K$; and $y, y_{\text{sub}}, \rho_k$ are uniformly sampled from $\{-1, 1\}$.

Features and Noise. We assume features and noise form an orthonormal basis of \mathbb{R}^d , i.e., a set of unit orthogonal vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_d\}$ in \mathbb{R}^d . W.l.o.g., one can let \mathbf{v} 's be the standard basis, where the first K basis are feature vectors. $\{\phi_1, \dots, \phi_K\}$ are constants indicating the strength of each feature, and $\{\mu_1, \dots, \mu_K\}$ are the means of the corresponding entries in the feature vectors. In particular:

- Class Feature: \mathbf{v}_1 .
- Subclass Feature: \mathbf{v}_2 .
- (Class and subclass) irrelevant features:¹ $\mathbf{v}_3, \dots, \mathbf{v}_K$.
- Noise $\boldsymbol{\xi} \sim \mathcal{D}_{\boldsymbol{\xi}}$: $\mathcal{D}_{\boldsymbol{\xi}}$ is a uniform distribution over features $\sigma_{\boldsymbol{\xi}}\mathbf{v}_1, \dots, \sigma_{\boldsymbol{\xi}}\mathbf{v}_d$, where $\sigma_{\boldsymbol{\xi}}$ indicates the variance of the noise.²

¹In the rest of the paper, we use irrelevant features to refer to features that may have semantic meaning but are irrelevant to class and subclass.

²This definition of noise is nearly identical to Gaussian noise $\mathcal{N}(0, \frac{\sigma_{\boldsymbol{\xi}}^2}{d}\mathbf{I}_d)$ in the high-dimensional regime but keeps the analysis clear.

We sample n examples from $\mathcal{D}_{\text{orig}}$ to form the original dataset $\hat{\mathcal{D}}_{\text{orig}}$.

Assumption B.1 (Balanced Dataset). All combinations of $(y_i, y_{\text{sub},i}, k_i, \rho_i)$ are equally represented in $\hat{\mathcal{D}}_{\text{orig}}$.³

A Concrete Example of the Above Data Distribution. Let $y = 1$ be dogs and $y = -1$ be cats, $y_{\text{sub}} = 1$ if they are fluffy and $y_{\text{sub}} = -1$ if they are not-fluffy. Then $(\phi_1 + \mu_1)\mathbf{v}_1 + (\phi_2 + \mu_2)\mathbf{v}_2$ denotes a fluffy dog. Here, the background can be interpreted as an irrelevant feature: let $\rho_3 = 1$ for grass and $\rho_3 = -1$ for forest. Then $(\phi_1 + \mu_1)\mathbf{v}_1 + (\phi_2 + \mu_2)\mathbf{v}_2 + (\phi_3 + \mu_3)\mathbf{v}_3$ represents a fluffy dog on grass. Note that each example only selects one irrelevant feature, which mimics the real world, where examples do not necessarily have all types of objects in the background i.e. many examples have neither grass or forests as their background.

Rationale for Including Feature Means μ_i . In general, it is unreasonable to expect all features to have 0 expectation over entire data, thus we introduce μ to further generalize our analysis. We find that considering a non-zero mean for the subclass feature is sufficient to provide novel insights into class collapse (Theorem C.5). Therefore, for clarity, we set all the μ 's except μ_2 to zero.

Relation to Sparse Coding Model. This data distribution is a variant of the sparse coding model which is usually considered as a provision model for studying the feature learning process in machine learning (e.g., (Zou et al., 2021; Wen & Li, 2021; Liu et al., 2021)). It naturally fits into many settings in machine learning, and in general mimics the outputs of intermediate layers of neural networks which have been shown to be sparse (Papayan et al., 2017). It is also used to model the sparse occurrences of objects in image tasks (Olshausen & Field, 1997; Vinje & Gallant, 2000; Foldiak, 2003; Protter & Elad, 2008; Yang et al., 2009; Mairal et al., 2014) and polysemy of words in language tasks (Arora et al., 2018).

B.2. Data Augmentation $\mathcal{A}(\cdot)$

For each example in $\hat{\mathcal{D}}_{\text{orig}}$, we generate m augmentations to form $\hat{\mathcal{D}}_{\text{aug}}$. We consider the following augmentation strategy: given an example $\mathbf{x} = \mathbf{u} + \boldsymbol{\xi}$, its augmentation is given by $\mathcal{A}(\mathbf{x}) = \mathbf{u} + \boldsymbol{\xi}'$, where $\boldsymbol{\xi}'$ is a new random variable from \mathcal{D}_{ξ} independent of $\boldsymbol{\xi}$. This is an abstract of augmentations used in practice where two augmentations from the same example share certain parts of the features and have the correlation between their noise parts removed or weakened.

Assumption B.2 (High dimensional regime). d is at least $\omega(n^2m^2)$.

Assumption B.3 (Sufficient sample size). The noise-to-sample-size ratio is not too large $\frac{\sigma_{\xi}^2}{mn} = o(1)$.

B.3. Linear Model

We consider a linear model with p outputs. The model has weights $\mathbf{W} \in \mathbb{R}^{p \times d}$ and bias $\mathbf{b} \in \mathbb{R}^p$ where $p \geq 3$. The function represented by the model is $f_{\Theta}(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$, where we define $\Theta \in \mathbb{R}^{p \times (d+1)}$ as the concatenated parameter $[\mathbf{W} \ \mathbf{b}]$. We establish theoretical proofs of class collapse and feature suppression for linear model, and also empirically verified them for (non-linear) deep neural networks.

B.4. Loss function

For unsupervised contrastive learning, we use the unsupervised spectral contrastive loss popular in prior theoretical and empirical work (HaoChen et al., 2021; Saunshi et al., 2022; HaoChen & Ma, 2022) and for supervised contrastive learning, we consider the natural generalization of this loss to incorporate supervision. Let \mathcal{A}_i denote the set of augmentations in $\hat{\mathcal{D}}_{\text{aug}}$ generated from the i -th original example with $\mathcal{A}(\cdot)$. Let \mathcal{S}_{+1} and \mathcal{S}_{-1} denote the set of augmentations in $\hat{\mathcal{D}}_{\text{aug}}$ with class labels $+1$ and -1 , respectively. Let $\hat{\mathbb{E}}$ denote the empirical expectation. Then we have the following loss functions:

$$\begin{aligned} \mathcal{L}_{\text{UCL}}(\Theta) = & -2\hat{\mathbb{E}}_{i \in [n], \mathbf{x} \in \mathcal{A}_i, \mathbf{x}^+ \in \mathcal{A}_i} [f_{\Theta}(\mathbf{x})^{\top} f_{\Theta}(\mathbf{x}^+)] \\ & + \hat{\mathbb{E}}_{\mathbf{x} \in \hat{\mathcal{D}}_{\text{aug}}, \mathbf{x}^- \in \hat{\mathcal{D}}_{\text{aug}}} [(f_{\Theta}(\mathbf{x})^{\top} f_{\Theta}(\mathbf{x}^-))^2] \end{aligned} \quad (1)$$

$$\begin{aligned} \mathcal{L}_{\text{SCL}}(\Theta) = & -2\hat{\mathbb{E}}_{c \in \{-1, 1\}, \mathbf{x} \in \mathcal{S}_c, \mathbf{x}^+ \in \mathcal{S}_c} [f_{\Theta}(\mathbf{x})^{\top} f_{\Theta}(\mathbf{x}^+)] \\ & + \hat{\mathbb{E}}_{\mathbf{x} \in \hat{\mathcal{D}}_{\text{aug}}, \mathbf{x}^- \in \hat{\mathcal{D}}_{\text{aug}}} [(f_{\Theta}(\mathbf{x})^{\top} f_{\Theta}(\mathbf{x}^-))^2]. \end{aligned} \quad (2)$$

Our results can be extended to the Gaussian noise setting.

³This can be approximately achieved when n is sufficiently larger than K . While our analysis can be generalized to consider imbalanced data, this is outside the scope of this work.

C. Main Results

C.1. Simplicity Bias Contributes to Class Collapse in Supervised CL

We make two key observations through our theoretical analysis and experiments (henceforth we refer to class collapse at *test time* simply as ‘class collapse’):

1. Theoretically, not all global minimizers exhibit class collapse, but the *minimum norm* minimizer does.
2. Theoretically and empirically, when the model is trained using (S)GD, some subclasses are *provably* learned early in training. Empirically, however, those subclasses will eventually be unlearned i.e. S(GD) converges to minimizers that exhibit class collapse.

Altogether, these observations suggest that class collapse, which has been observed in practice when certain gradient-based algorithms are used to minimize the loss, cannot be explained by simply analyzing the loss function. This highlights the importance of studying the dynamics and inductive bias of training algorithms in contrastive learning.

C.1.1. WHAT MINIMIZERS HAVE CLASS COLLAPSE?

We first define class collapse in terms of the alignment between the model weights and the subclass feature.

Definition C.1 (Exact class collapse). We say exact class collapse happens at test time when:

$$\forall \beta \in \mathbb{R}^p, \Pr_{(\mathbf{x}, y, y_{\text{sub}}) \sim \mathcal{D}_{\text{orig}}} (y_{\text{sub}} \beta^\top f_{\Theta}(\mathbf{x}) > 0) = 1/2.$$

The definition means that no linear classifier on the embeddings of examples drawn from $\mathcal{D}_{\text{orig}}$ can predict the subclass label with accuracy beyond random guess.⁴

This is different from class collapse on the training set which is not defined on the population set $\mathcal{D}_{\text{orig}}$ but on the training samples $\hat{\mathcal{D}}_{\text{orig}}$.

Proposition C.2. For any $\Theta^* \in \min_{\Theta} \mathcal{L}_{\text{SCL}}(\Theta)$, we have $f_{\Theta^*}(\mathbf{x}_i) = f_{\Theta^*}(\mathbf{x}_j)$ for all $\mathbf{x}_i, \mathbf{x}_j$ in the training set $\hat{\mathcal{D}}_{\text{aug}}$ such that $y_i = y_j$.

This directly implies that minimizing the loss results in class collapse on the training set. However, the following theorem C.3 shows that minimizing the loss does not necessarily lead to class collapse on the test set. To determine whether class collapse occurs, we need to determine whether the model learns the subclass feature. With a linear model, this exactly corresponds to constant alignment between weights and the subclass feature.

Theorem C.3 (Minimizing $\mathcal{L}_{\text{SCL}} \not\Rightarrow$ Class Collapse). With high probability i.e. at least $1 - O(\frac{m^2 n^2}{d}) = 1 - o(1)$, there exists $\Theta^* = [\mathbf{W}^* \ \mathbf{b}^*]$ such that $\Theta^* \in \min_{\Theta} \mathcal{L}_{\text{SCL}}(\Theta)$ \mathbf{W}^* has constant alignment with subclass feature \mathbf{v}_2 i.e.

$$\|\mathbf{W}^* \mathbf{v}_2\| = \Omega(1).$$

Hence, there exists a linear classifier in the embedding space that can predict subclass labels almost perfectly. I.e.,

$$\exists \beta, \text{ s.t. } \Pr_{(\mathbf{x}, y, y_{\text{sub}}) \sim \mathcal{D}_{\text{orig}}} (y_{\text{sub}} \beta^\top \mathbf{W}^* \mathbf{x} > 0 | y) = 1 - o(1).$$

We prove the theorem in Appendix G. The proof utilizes Lemma F.1 which implies that, due to the high-dimensionality, the noise vectors have non-trivial effects on the empirical covariance matrix by rotating its kernel space. This results in the kernel space to have a $\Theta(\frac{\sigma_{\xi}}{\sqrt{mn}})$ alignment with the subclass feature. Since minimizers of the loss can behave arbitrarily on this kernel space, without any additional restriction, they can have any alignment with the subclass feature.

Next, we show that, the *minimum norm* minimizer exhibits class collapse.

⁴Actually we are able to analyze a stronger version of class collapse: $\Pr_{(\mathbf{x}, y, y_{\text{sub}}) \sim \mathcal{D}_{\text{orig}}} (f_{\Theta}(\mathbf{x}) | y_{\text{sub}}) = \Pr_{(\mathbf{x}, y, y_{\text{sub}}) \sim \mathcal{D}_{\text{orig}}} (f_{\Theta}(\mathbf{x}))$, which means the distributions of embeddings given and not given the subclass label are exactly the same. Nonetheless, we present this simpler formulation for clarity.

Theorem C.4 (Minimizing \mathcal{L}_{SCL} + Minimum Norm \implies Class Collapse). Assume $\mu_2 = 0$. Let $\Theta^{**} = [\mathbf{W}^{**} \ \mathbf{b}^{**}]$ be the minimum norm minimizer of \mathcal{L}_{SCL} , i.e.,

$$\Theta^{**} = \arg \min_{\Theta} \|\Theta^*\|_F \text{ s.t. } \Theta^* \in \arg \min_{\Theta} \mathcal{L}_{\text{SCL}}(\Theta).$$

Then with high probability i.e. at least $1 - O(\frac{m^2 n^2}{d}) = 1 - o(1)$, \mathbf{W}^{**} has no alignment with subclass feature \mathbf{v}_2 i.e.

$$\|\mathbf{W}^{**} \mathbf{v}_2\| = 0.$$

This means class collapse occurs at test time (Definition C.1), and no linear classifier does better than random guess for predicting subclass labels.

Theorems C.3 and C.4 show that minimizing the training loss does not necessarily lead to class collapse on test data, but does with additional constraint on the weights of the model. This is not due to a degenerate solution, as we show that both minimizers learn the class feature \mathbf{v}_1 (see corollary F.5).

C.1.2. INTRIGUING PROPERTIES OF GD

We now further our theoretical characterization of class collapse by investigating the setting where \mathcal{L}_{SCL} is minimized by GD. This is an important step toward understanding class collapse in practice, where similar optimization algorithms are used to minimize the loss. Our findings indicate that it is likely the simplicity bias of commonly used optimization algorithms that eventually leads to class collapse.

We consider GD with a constant learning rate η . The weights are initialized from a Gaussian distribution, i.e., the initial weight Θ_0 has each of its element drawn from $\mathcal{N}(0, \frac{\sigma_0^2}{d})$. And the weights at training epoch t are given by:

$$\Theta_t = \Theta_{t-1} - \eta \nabla_{\Theta} \mathcal{L}_{\text{SCL}}(\Theta_{t-1}).$$

Early in Training Some Subclasses are Provably Learned. By analyzing the training dynamics of GD, we find that subclasses are learned early in training.

Theorem C.5 (Early in training subclass features are learned). Assume $\sigma_0 \sqrt{\frac{p}{d}} = o(1)$ and $\sigma_{\xi} = o(1)$. If the subclass feature has a constant non-zero mean such that $1 + \mu^2 > \phi_1^2$, then with probability at least $1 - O(\frac{m^2 n^2}{d} + \frac{1}{\text{poly}(p)}) = 1 - o(1)$ the following holds:

- $\|\mathbf{W}_0 \mathbf{v}_2\| = o(1)$.
- $\exists t = O(\ln(\frac{1}{\sigma_0 \sqrt{\frac{p}{d}}}))$, s.t. $\|\mathbf{W}_t \mathbf{v}_2\| = \Omega(1)$, and
- $\exists \beta$, s.t. $\Pr_{(\mathbf{x}, y, y_{\text{sub}}) \sim \mathcal{D}_{\text{orig}}}(y_{\text{sub}} \beta^{\top} \mathbf{W}_t \mathbf{x} > 0 | y) = 1 - o(1)$.

The above theorem shows that there exists an epoch where the weights have constant alignment with the subclass feature and produce distinguishable subclass embeddings (proof in Appendix J).

The key step of our analysis is showing that early in training, GD aligns the weights with the first eigenvector of the covariance matrix of class centers. This alignment grows exponentially faster than alignments with any other directions. When $1 + \mu^2 > \phi_1^2$, the subclass feature has a constant projection onto the first eigenvector and is therefore learned by the model.

More importantly, the same phenomenon can be observed in *neural networks*. We use SGD to train a ResNet18 (He et al., 2016) on CIFAR-100 (Krizhevsky et al., 2009) with supervised CL loss (Khosla et al., 2020) with 20 class (superclass) labels, and perform linear evaluation on embeddings of test data with 100 subclass (class) labels (see details in Appendix K). We observe that the subclass accuracy increases during the first 200 epochs before it starts to drop (Figure 3(a)). Some subclasses can even achieve a high accuracy around 80% (Figure 3(b)). This is surprising as it confirms that models trained with commonly used loss functions *do* learn subclass features early in training.

Empirical Evidence Showing that Class Collapse Eventually Happens in (S)GD. We simulate our theoretical analysis using numerical experiments to show that gradient descent converges to a minimizer that exhibits class collapse, despite learning subclasses early in training. We visualize the embeddings of test data at different epochs in Figure 1, and plot the

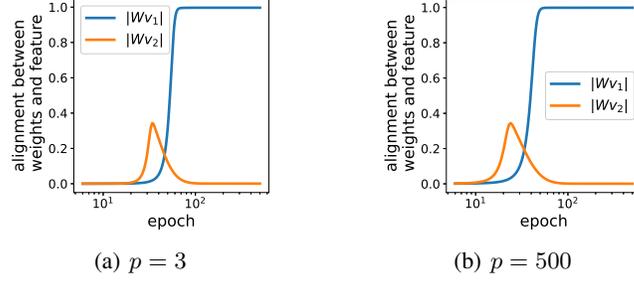


Figure 2. $\|\mathbf{W}_t \mathbf{v}_1\|$ and $\|\mathbf{W}_t \mathbf{v}_2\|$ at different epochs. Both features are learned early in training, but \mathbf{v}_2 is unlearned later.

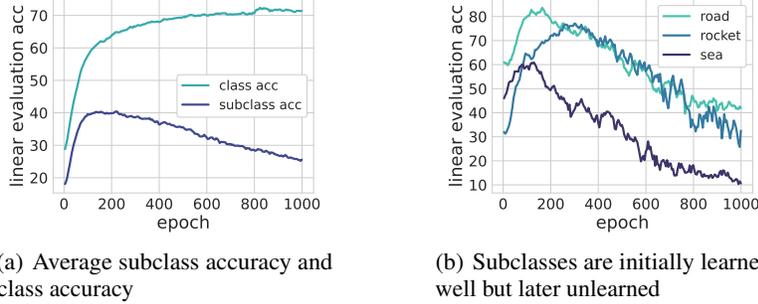


Figure 3. (a) Average subclass accuracy and class accuracy. (b) Accuracy in subclasses ‘road’, ‘rocket’ and ‘sea’. In both plots, the subclass accuracy increases and then decreases, which confirms that subclasses are learned early in training before class collapse happens. The class accuracy only increases during training.

alignment between weights and class/subclass features in Figure 2. Subclasses are perfectly separated and the weights align with both \mathbf{v}_1 and \mathbf{v}_2 after around 100 epochs of training. The model then starts unlearning \mathbf{v}_2 which causes the alignment to drop, thus subclasses are merged in the embedding space. We also confirm that same conclusion holds for neural networks in realistic settings. In Figure 3, we see that the subclass accuracy drops after around 200 epochs of training and eventually reaches a low value. In contrast, the class accuracy does not drop during training.

Minimum Norm Minimizer Exhibits Class Collapse. Note that in Theorem C.5, assuming $\mu \neq 0$ leads us to discovering that subclasses are learned early in training. Here, we extend Theorem C.4 to this setting under asymptotic class collapse.

Definition C.6 (Asymptotic Class Collapse). We say asymptotic class collapse happens when $\|\mathbf{W} \mathbf{v}_2\| = O(\frac{\sigma_\xi}{\sqrt{mn}}) = o(1)$.

This definition implies that: (1) representations of subclasses are not well separated, hence it is nearly impossible to distinguish between them, and (2) the distinguishability of subclasses is at odds with generalization, which improves as number of augmented views per example m and size of training data n increase. Thus, while this definition is a relaxation of Definition C.1, practically, this results in equally severe class collapse.

Theorem C.7 (Extension of Theorem C.4 for $\mu_2 \neq 0$). Let $\Theta^{**} = [\mathbf{W}^{**} \mathbf{b}^{**}]$ be the minimum norm minimizer of \mathcal{L}_{SCL} :

$$\Theta^{**} = \arg \min_{\Theta^*} \|\Theta^*\|_F \text{ s.t. } \Theta^* \in \arg \min_{\Theta} \mathcal{L}_{SCL}(\Theta).$$

Then with probability at least $1 - O(\frac{m^2 n^2}{d}) = 1 - o(1)$, asymptotic class collapse happens, i.e.,

$$\|\mathbf{W}^{**} \mathbf{v}_2\| = O(\frac{\sigma_\xi}{\sqrt{mn}}) = o(1).$$

C.1.3. SIMPLICITY BIAS OF (S)GD

We reiterate our main findings:

1. Minimizing the supervised contrastive loss *does not* necessarily lead to class collapse.
2. However, *simpler* minimizers of the supervised contrastive loss (e.g. *minimum norm*) do suffer from class collapse.

3. Optimizing with (S)GD does learn the subclass features early in training, but eventually unlearns them, resulting in class collapse.

These coupled with the fact that (S)GD is known to have a bias towards simpler solutions (Kalimeris et al., 2019) prompt us to conjecture:

The simplicity bias of (S)GD leads it to unlearn subclass features, thus causing class collapse.

The simplicity bias of (S)GD has not been rigorously studied for CL, and our results indicate the surprising role it may play in class collapse. Note that, the supervised contrastive loss is different than common supervised objectives, where the role of such bias of (S)GD is understood better (Gunasekar et al., 2018; Soudry et al., 2018; Ji & Telgarsky, 2019; Wu et al., 2019; Lyu et al., 2021). Rather, the supervised CL objective can be reformulated as a matrix factorization objective, where the debate on the bias of (S)GD (e.g., minimum norm (Gunasekar et al., 2017) or rank (Arora et al., 2019a; Razin & Cohen, 2020)) is still ongoing.

C.2. Understanding Feature Suppression in Unsupervised CL

Empirically, feature suppression can be observed due to a variety of reasons (Li et al., 2023; Chen et al., 2021; Robinson et al., 2021). Easy features for unsupervised CL are those that allow the model to discriminate between examples (highly discriminative). Here, we consider different ways irrelevant features can be easy (highly discriminative) and characterize how this can lead to feature suppression. We show that the types of feature suppression we consider can be largely attributed to insufficient embedding dimensionality and/or poor data augmentations. Surprisingly, we find again that the minimum norm simplicity bias is critical in explaining this phenomenon.

C.2.1. FEATURE SUPPRESSION DUE TO EASY IRRELEVANT FEATURES AND LIMITED EMBEDDING SPACE

In Theorem C.8, we show that easy (discriminative) irrelevant features can suppress the class feature when the embedding dimensionality is limited. For clarity, we let $\mu_2 = 0$.

Theorem C.8 (Feature Suppression 1). *Assume $p \leq K$. Let L be the $(K + 1)$ -element tuple $[1, \phi_1^2, \phi_2^2, \frac{\phi_3^2}{K-2}, \dots, \frac{\phi_K^2}{K-2}]$ whose last K elements are the variances of features. If ϕ_1^2 is not among the p largest elements in L , then with probability at least $1 - O(\frac{m^2 n^2}{d}) = 1 - o(1)$: (1) there exists a global minimizer Θ^* of \mathcal{L}_{UCL} such that $\|\mathbf{W}^* \mathbf{v}_1\| = \Omega(1)$, (2) However, the minimum norm minimizer Θ^{**} satisfies $\|\mathbf{W}^{**} \mathbf{v}_1\| = 0$.*

We prove the theorem in Appendix H. The elements except the first one in tuple L can be interpreted as the variance of examples at each coordinate $v_k, k = 1, 2, \dots, K$, which indicates how much the examples are discriminated by each feature. The theorem shows that when the embedding space is not large enough to represent all the K features (which requires $K + 1$ dimensions), the minimum norm minimizer only picks the most discriminative ones. In practice, the embedding space in unsupervised CL is relatively low-dimensional (compared to input dimensionality) and thus the model cannot fit all the information about inputs into the embedding space. As is suggested by Theorem C.8, if the training algorithm prefers functions with certain simple structures, only the easiest (most discriminative) features that can be mapped into the embedding space by less complex functions (e.g., smaller norm) are learned. The class features are suppressed if they are not amongst the easiest ones.

Remark C.9. Following the same analysis we can also show that when ϕ_1 is among the p largest elements in L , i.e., the class feature is among the easiest (most discriminative) ones, the class feature \mathbf{v}_1 is learned by the minimum norm minimizer; when ϕ_1 is exactly on par with some other element as the p -th largest, there exist both minimum norm minimizers that learn and do not learn the class feature \mathbf{v}_1 .

Numerical Experiments with GD. Our theory for the minimum norm minimizer matches the experimental results for models trained with GD. We let $p = K$ and let $1 \geq \phi_2^2 \geq \frac{\phi_3^2}{K-2} \geq \dots \geq \frac{\phi_{K-1}^2}{K-2} > \phi_1^2$ so that ϕ_1^2 must be among the smallest two variances i.e. \mathbf{v}_1 is among the two most difficult features. Then we vary ϕ_K and see how the trained weights align with \mathbf{v}_1 . Consistent with Theorem 1, Figure 4 shows that \mathbf{v}_1 is suppressed when $\frac{\phi_K^2}{K-2} > \phi_1^2$. Interestingly, we also see that the result at $\frac{\phi_K^2}{K-2} = \phi_1^2$ diverges, indicating that GD can find both minimizers that learn and do not learn \mathbf{v}_1 when the variances at \mathbf{v}_1 and \mathbf{v}_K are the same.

Empirically Verifying Benefits of Larger Embedding Size. Theorem C.8 also provides one practical solution for feature suppression due to limited embedding size: increasing the embedding size so that every feature can be learned by the model.

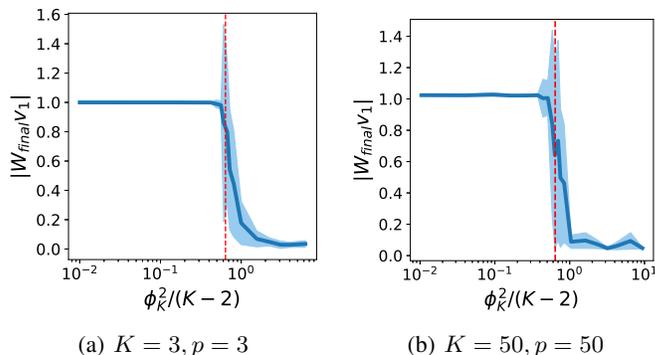


Figure 4. The irrelevant feature suppresses the class feature when its variance is beyond the variance of the class feature (the red vertical line). We let $d = 2000, p = K, \phi_1 = 0.8, \phi_2 = 1, \mu = 0, \frac{\phi_k^2}{K-2} > \phi_1, \forall k \in [K-1]$ and vary ϕ_K . Thus whether ϕ_1^2 is among the p largest variances only depends on ϕ_K . We train the linear model to convergence. Plots show that the alignment between the trained weights and v_1 drops when ϕ_K increases. We report the average of 10 runs. The result diverges at $\frac{\phi_K^2}{K-2} = \phi_1^2$ indicating that the model can learn either v_1 or v_K in this case.

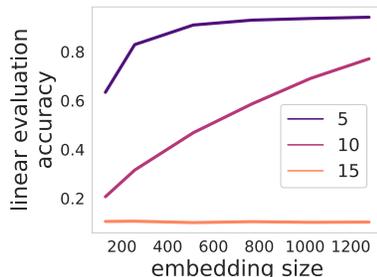


Figure 5. Effect of embedding size on feature suppression in MNIST RandBit(Chen et al., 2021). Legends show the number of bits in the extra channel which indicates how easy (discriminative) the irrelevant features are. We observe that (1) increasing the easiness of irrelevant features exacerbates feature suppression; (2) increasing the embedding size alleviates feature suppression.

Table 2. Effect of embedding size on feature suppression in CIFAR-10/100 RandBit. ‘Acc’ refers to class accuracy and ‘Sub Acc’ refers to subclass accuracy. We see that increasing embedding size alleviates feature suppression, improving class/subclass accuracy.

| w | CIFAR-10 RandBit | | CIFAR-100 RandBit | |
|-----|------------------|-------|-------------------|-------|
| | Sub Acc | Acc | Sub Acc | Acc |
| 4 | 34.38 | 86.73 | 11.67 | 23.53 |
| 64 | 71.96 | 96.82 | 34.11 | 52.32 |
| 128 | 76.69 | 97.65 | 38.51 | 57.40 |

To provide empirical evidence for this, we conduct two sets of experiments:

First, we train 5-layer convolutional networks on the RandomBit dataset with the same setup as in (Chen et al., 2021), but we vary the embedding size (see details in Appendix K). Varying the # bits in the extra channel intuitively controls how discriminative the irrelevant feature are, i.e., how easy-to-learn it is for CL. In this setting, the random bit can suppress the MNIST digits. We make two observations in Figure C.2.1: (1) with a fixed embedding size, increasing easiness (number of random bits) of the irrelevant features exacerbates feature suppression; (2) with a fixed easiness of irrelevant features, increasing the embedding size alleviates feature suppression.

Second, we train ResNet18 (He et al., 2016) on the CIFAR-10/100 RandBit Dataset, constructed similarly to the MNIST RandBit dataset but with images from CIFAR-10/100 (Krizhevsky et al., 2009) (see Appendix K.1). For CIFAR-10, we use 2 random bits, and for CIFAR-100, we use one random bit as the class irrelevant features. Table 2 presents the test performance for different values of the model width w , where a larger w indicates a larger embedding size (see Appendix K.3 for details). On both datasets, increasing the embedding size alleviates feature suppression, leading to improvements in both class and subclass accuracies. We also provide additional experiments and discussion in Appendix K.3. Both experimental results confirm the conclusion drawn from the theoretical analysis.

C.2.2. FEATURE SUPPRESSION DUE TO HIGH-DIMENSIONAL IRRELEVANT FEATURES AND IMPERFECT AUGMENTATION

Empirically, another form of feature suppression has been observed that cannot be remedied by larger embedding dimensionality (Li et al., 2023). We characterize this form of feature suppression by defining easy irrelevant features as being: (1) drawn from a high dimensional space so that the collection of irrelevant features is large and discriminating based on irrelevant features is easier, (2) less altered by data augmentation compared to the class feature.

For (1), formally we assume $K = \omega(n^2)$, as opposed to assumption D.1 which implies that K is smaller than n . A consequence of this assumption is that with high probability the n original examples each have a unique irrelevant feature. For (2) we consider the following imperfect data augmentation:

Definition C.10 (Imperfect data augmentation $\mathcal{A}'(\cdot)$). For a given example $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\xi} \in \hat{\mathcal{D}}_{\text{orig}}$,

$$\mathcal{A}'(\mathbf{x}) = \mathbf{u} + \zeta' \mathbf{v}_1 + \zeta'' \mathbf{v}_2 + \boldsymbol{\xi}',$$

where $\zeta' \sim \mathcal{N}(0, \sigma_{\zeta'}^2)$, $\zeta'' \sim \mathcal{N}(0, \sigma_{\zeta''}^2)$, $\sigma_{\zeta'}^2, \sigma_{\zeta''}^2 \neq 0$ and $\boldsymbol{\xi}'$ is a new random variable drawn from $\mathcal{N}(\boldsymbol{\xi}, \boldsymbol{\Sigma}_{\boldsymbol{\xi}})$ with $\text{rank}(\boldsymbol{\Sigma}_{\boldsymbol{\xi}}) \leq \frac{m}{2}$.

In the definition, the data augmentation adds small perturbations (ζ' and ζ'') to class and subclass features, weakly alters the noise, but preserves the irrelevant features. For example, on Colorful-Moving-MNIST (Tian et al., 2020) constructed by assigning each MNIST digit a background object image selected randomly from STL-10, the colorful background objects are high-dimensional and the colors are invariant to data augmentations without color distortion.

Theorem C.11 (Feature Suppression 2). *If $K = \omega(n^2)$ and augmentation is $\mathcal{A}'(\cdot)$, with probability $\geq 1 - o(\frac{n^2 m^2}{d} + \frac{1}{n}) = 1 - o(1)$, the minimum norm minimizer $\boldsymbol{\Theta}^* = [\mathbf{W}^*, \mathbf{b}^*]$ satisfies $\|\mathbf{W}^* \mathbf{v}_1\| = 0$.*

This theorem shows that feature suppression can happen even when embedding dimensionality p is arbitrarily large and helps understand empirical observations made both in our work (Figure C.2.1, the line with 15 bits) and previous work. For example Li et al. (2023) showed that on Colorful-Moving-MNIST, the colorful background can suppress learning the digits especially when color distortion is not used in augmentation, and increasing embedding size does not address the issue.

In conclusion, Theorem C.11 highlights that designing data augmentations that disrupt the highly-discriminative irrelevant features is a key to addressing feature suppression.

C.3. Combining Supervised and Unsupervised CL Losses Can Avoid Both Class Collapse and Feature Suppression

We now consider the following loss which is a weighted sum of the supervised and unsupervised CL loss functions:

$$\mathcal{L}_{\text{joint}, \beta}(\boldsymbol{\Theta}) = \beta \mathcal{L}_{\text{SCL}}(\boldsymbol{\Theta}) + (1 - \beta) \mathcal{L}_{\text{UCL}}(\boldsymbol{\Theta}).$$

Similar loss functions have been proposed recently with notable empirical success. For example, Chen et al. (2022) put forth a weighted sum of supervised CL loss and class-conditional InfoNCE (which has similar effect as \mathcal{L}_{UCL} in our setting) to avoid class collapse. Islam et al. (2021) empirically observed that the joint objective of supervised and unsupervised contrastive loss leads to better transferability of the learned models than their supervised counterparts. However, we still lack a theoretical understanding of why this weighted sum of losses can outperform both losses.

From our investigation of class collapse and feature suppression, the benefit of the joint objective $\mathcal{L}_{\text{joint}}$ becomes evident: the unsupervised term in $\mathcal{L}_{\text{joint}}$ increases the chance of learning features that do not appear relevant to the labels but might be useful for downstream tasks, while the supervised term in $\mathcal{L}_{\text{joint}}$ ensures that even hard-to-learn class features are learnt. Thus, $\mathcal{L}_{\text{joint}}$ can learn rich representations capturing more task relevant information than either $\mathcal{L}_{\text{UCL}}(\boldsymbol{\Theta})$ or $\mathcal{L}_{\text{SCL}}(\boldsymbol{\Theta})$. We show below that with an appropriate choice of β , $\mathcal{L}_{\text{joint}}$ can provably succeed where \mathcal{L}_{SCL} fails due to collapse and \mathcal{L}_{UCL} fails due to feature suppression (for clarity, we let $\mu = 0$).

Theorem C.12. *W.L.O.G., assume $\phi_3 \geq \phi_4 \geq \dots \geq \phi_K$. If $p \leq K$, $\phi_2^2 > \frac{\phi_{p-2}^2}{K-2}$ and $\phi_1^2 < \frac{\phi_{p-1}^2}{K-2}$, then by Theorem C.4 the minimum norm minimizer of \mathcal{L}_{SCL} suffers from class collapse and by Theorem C.8 the minimum norm minimizer of \mathcal{L}_{UCL} suffers from feature suppression. However, for constant $\beta \in (0, 1)$, the minimum norm minimizer of $\mathcal{L}_{\text{joint}, \beta}$, denoted by $\boldsymbol{\Theta}^* = [\mathbf{W}^*, \mathbf{b}^*]$, satisfies $\|\mathbf{W}^* \mathbf{v}_1\| = \Omega(1)$ and $\|\mathbf{W}^* \mathbf{v}_2\| = \Omega(1)$.*

Empirically Verifying Benefits of the Joint Loss. We empirically examine the impact of the joint loss on MNIST RandBit, CIFAR-100, and CIFAR-100 RandBit. The training details are in Appendix K.2. The results indicate that the

Table 3. Joint loss alleviates class collapse on CIFAR-100.

| Loss | Subclass Acc |
|------------------------------|--------------|
| SCL | 26.11 |
| Joint loss ($\beta = 0.8$) | 41.37 |

Table 4. Joint loss alleviates feature suppression on MNIST RandBit.

| Loss | Class Acc |
|------------------------------|-----------|
| UCL | 61.21 |
| Joint loss ($\beta = 0.5$) | 79.37 |

Table 5. Joint loss alleviates both class collapse and feature suppression on CIFAR-100 RandBit.

| Loss | Subclass Acc | Class Acc |
|------------------------------|--------------|-----------|
| SCL | 28.13 | 61.10 |
| UCL | 34.11 | 52.32 |
| Joint loss ($\beta = 0.8$) | 35.72 | 63.94 |

joint loss significantly improves performance in scenarios where SCL suffers from class collapse (Table 3) and UCL suffers from feature suppression (Table 4). Furthermore, on CIFAR-100 RandBit dataset, where both phenomena can occur simultaneously, the joint loss effectively alleviates both issues (Table 5).

D. Preliminaries for The Proofs

D.1. Additional Assumptions, Propositions and Definitions

We assume the dataset is balanced. This can be approximately achieved when n is sufficiently larger than K . While our analysis can be generalized to consider imbalanced data, this is outside the scope of this work.

Assumption D.1 (Balanced Dataset). All combinations of $(y_i, y_{\text{sub},i}, k_i, \rho_i)$ are equally represented in $\hat{\mathcal{D}}_{\text{orig}}$.

The following proposition shows that class collapse on training set is directly implied by minimizing the training loss.

Proposition D.2. For any $\Theta^* \in \min_{\Theta} \mathcal{L}_{\text{SCL}}(\Theta)$, we have $f_{\Theta^*}(\mathbf{x}_i) = f_{\Theta^*}(\mathbf{x}_j)$ for all $\mathbf{x}_i, \mathbf{x}_j$ in the training set $\hat{\mathcal{D}}_{\text{aug}}$ such that $y_i = y_j$.

The following definition defines asymptotic class collapse, which we will demonstrate in Theorem C.7.

Definition D.3 (Asymptotic Class Collapse). We say asymptotic class collapse happens when $\|\mathbf{W}\mathbf{v}_2\| = O(\frac{\sigma_{\xi}}{\sqrt{mn}}) = o(1)$.

D.2. Effective dataset

Analyzing training a linear model with bias on the data is equivalent to analyzing training a linear model without bias on:

$\left\{ \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix} : \mathbf{x}_i \in \mathcal{D}_{\text{aug}} \right\}$. Equivalently we can consider a dataset distribution where

$$\mathbf{x} = \mathbf{u} + \boldsymbol{\xi},$$

where $\mathbf{u} = \mathbf{v}_0 + y\phi_1\mathbf{v}_1 + (y_{\text{sub}}\phi_2 + \mu)\mathbf{v}_2 + \rho\phi_k\mathbf{v}_k$.

The definitions are identical to the one in Section B.1 except that each data now is in \mathbb{R}^{d+1} and has one constant feature \mathbf{v}_0 orthogonal to other \mathbf{v} 's. We train a linear model $f(\mathbf{x}) = \mathbf{W}\mathbf{x}$ on such data. The definition of other notations such as $\hat{\mathcal{D}}_{\text{aug}}$ in the following analysis are also adapted to this dataset accordingly. Other notations such as $\hat{\mathcal{D}}_{\text{aug}}$ in the subsequent analysis are adjusted accordingly to accommodate this dataset.

D.3. Loss functions

The loss functions can be rewritten as follows

$$\begin{aligned} \mathcal{L}_{\text{SCL}} &= -2\hat{\mathbb{E}}_{i \in [n], \mathbf{x} \in \mathcal{A}_i, \mathbf{x}^+ \in \mathcal{A}_i} [\mathbf{x}^\top \mathbf{W}^\top \mathbf{W} \mathbf{x}^+] + \hat{\mathbb{E}}_{\mathbf{x} \in \hat{\mathcal{D}}_{\text{aug}}, \mathbf{x}^- \in \hat{\mathcal{D}}_{\text{aug}}} [(\mathbf{x}^\top \mathbf{W}^\top \mathbf{W} \mathbf{x}^-)^2] \\ &= -\text{Tr}(2\mathbf{M}^+ \mathbf{W} \mathbf{W}^\top) + \text{Tr}(\mathbf{M} \mathbf{W}^\top \mathbf{W} \mathbf{M} \mathbf{W}^\top \mathbf{W}) \end{aligned} \quad (3)$$

$$\begin{aligned} \mathcal{L}_{\text{UCL}} &= -2\hat{\mathbb{E}}_{c \in \{-1, 1\}, \mathbf{x} \in \mathcal{S}_c, \mathbf{x}^+ \in \mathcal{S}_c} [\mathbf{x}^\top \mathbf{W}^\top \mathbf{W} \mathbf{x}^+] + \hat{\mathbb{E}}_{\mathbf{x} \in \hat{\mathcal{D}}_{\text{aug}}, \mathbf{x}^- \in \hat{\mathcal{D}}_{\text{aug}}} [(\mathbf{x}^\top \mathbf{W}^\top \mathbf{W} \mathbf{x}^-)^2] \\ &= -\text{Tr}(2\tilde{\mathbf{M}} \mathbf{W} \mathbf{W}^\top) + \text{Tr}(\mathbf{M} \mathbf{W}^\top \mathbf{W} \mathbf{M} \mathbf{W}^\top \mathbf{W}) \end{aligned} \quad (4)$$

$$\begin{aligned} \mathcal{L}_{\text{joint}} &= (1 - \beta)\mathcal{L}_{\text{SCL}} + \beta\mathcal{L}_{\text{UCL}} \\ &= -\text{Tr}(2\bar{\mathbf{M}} \mathbf{W} \mathbf{W}^\top) + \text{Tr}(\mathbf{M} \mathbf{W}^\top \mathbf{W} \mathbf{M} \mathbf{W}^\top \mathbf{W}), \end{aligned} \quad (5)$$

where we define the following

Definition D.4. $\mathbf{M}, \mathbf{M}^+, \tilde{\mathbf{M}}$ are the covariance matrices of training examples, class centers and augmentation centers, respectively

$$\begin{aligned} \mathbf{M} &= \frac{1}{mn} \sum_{i=1}^{mn} \mathbf{x}_i \mathbf{x}_i^\top \\ \mathbf{M}^+ &= \frac{1}{2} \sum_{c \in \{-1, 1\}} \left(\frac{2}{mn} \sum_{\mathbf{x} \in \mathcal{S}_c} \mathbf{x} \right) \left(\frac{2}{mn} \sum_{\mathbf{x} \in \mathcal{S}_c} \mathbf{x} \right)^\top \\ \tilde{\mathbf{M}} &= \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{m} \sum_{\mathbf{x} \in \mathcal{A}_i} \mathbf{x} \right) \left(\frac{1}{m} \sum_{\mathbf{x} \in \mathcal{A}_i} \mathbf{x} \right)^\top, \end{aligned}$$

and

$$\bar{\mathbf{M}} = (1 - \beta)\mathbf{M}^+ + \beta\tilde{\mathbf{M}}.$$

E. Minimizers of Loss Functions

We start with a technical lemma which we will need:

Lemma E.1. *The product of two positive semidefinite matrices is diagonalizable.*

Next, we present a lemma that facilitates the analysis of minimizers for various contrastive loss functions. To apply the lemma, simply substitute the respective covariance matrices ($\mathbf{M}, \mathbf{M}^+, \tilde{\mathbf{M}}$) into \mathbf{P} and \mathbf{Q} as indicated.

Lemma E.2. *Let $\mathbf{P}, \mathbf{Q} \in \mathbb{R}^{d+1}$ be positive semidefinite matrices such that $\text{colsp}(\mathbf{P}) \subset \text{colsp}(\mathbf{Q})$. Consider the function $\mathcal{L} : \mathbb{R}^{p \times (d+1)} \rightarrow \mathbb{R}$ given by*

$$\mathcal{L}(\mathbf{W}) = \text{Tr}[-2\mathbf{W}^\top \mathbf{W} \mathbf{P} + \mathbf{W}^\top \mathbf{W} \mathbf{Q} \mathbf{W}^\top \mathbf{W} \mathbf{Q}] \quad (6)$$

Then \mathbf{W} is a global minimizer of \mathcal{L} if and only if

$$\mathbf{W}^\top \mathbf{W} \mathbf{Q} = [\mathbf{Q}^\dagger \mathbf{P}]_p$$

where notation $[\mathbf{A}]_p$ represents the matrix composed of the first p eigenvalues and eigenvectors of a positive semidefinite \mathbf{A} (if $p \geq \text{rank} \mathbf{A}$ then $[\mathbf{A}]_p = \mathbf{A}$).

Moreover, if $p \geq \text{rank}(\mathbf{P})$, then \mathbf{W}^{**} is a minimum norm global minimizer if and only if

$$\mathbf{W}^{**\top} \mathbf{W}^{**} = \mathbf{Q}^\dagger \mathbf{P} \mathbf{Q}^\dagger$$

Proof. First consider points that satisfy the first order condition

$$\nabla_{\mathbf{W}}(\mathcal{L}) = -4\mathbf{W} \mathbf{P} + 4\mathbf{W} \mathbf{Q} \mathbf{W}^\top \mathbf{W} \mathbf{Q} = 0 \quad (7)$$

with $\lambda_1, \dots, \lambda_r, \dots, \lambda_q \neq 0$ for some $r \leq q \leq s$, where $r = \mathbf{rank} \mathbf{W} \leq p, q = \mathbf{rank}(P)$.

Then for all such \mathbf{W} ,

$$\begin{aligned} \mathcal{L} &= \text{Tr}[-2\mathbf{W}^\top \mathbf{W} \mathbf{P} + \mathbf{W}^\top \mathbf{W} \mathbf{Q} \mathbf{W}^\top \mathbf{W} \mathbf{Q}] \\ &= -2 \sum_{i=1}^r \lambda_i^2 + \sum_{i=1}^r \lambda_i^2 \\ &= - \sum_{i=1}^r \lambda_i^2 \end{aligned}$$

It is clear from the above expression that the minimum among critical points is achieved if and only if

$$\mathbf{W}^\top \mathbf{W} \mathbf{Q} = [\mathbf{Q}^\dagger \mathbf{P}]_p$$

(note that if matching anything beyond the q th eigenvalue is trivial since all such eigenvalues are zero).

It remains to check the behavior as $\|\mathbf{W}\|_F$ grows large. Equivalently, $\mathbf{W}^\top \mathbf{W}$ has a large eigenvalue λ . Let \mathbf{w} be a corresponding eigenvector. If $\mathbf{w} \in \ker \mathbf{Q}$, then $\mathbf{Q} \mathbf{w} = \mathbf{P} \mathbf{w} = 0$, so we see that the loss is unchanged. Otherwise, \mathbf{w} has some nonzero alignment with $\mathbf{colsp}(\mathbf{W})$. But then $\text{Tr}[\mathbf{W}^\top \mathbf{W} \mathbf{Q} \mathbf{W}^\top \mathbf{W} \mathbf{Q}]$ grows quadratically in λ , but $\text{Tr}[-2\mathbf{W}^\top \mathbf{W} \mathbf{P}]$ grows at most linearly in λ , hence the loss is large. We conclude that the previously found condition in fact specifies the global minimizers of \mathcal{L} .

From now on, assume that $p \geq q$. Then the global minimum is achieved if and only if

$$\mathbf{W}^\top \mathbf{W} \mathbf{Q} = \mathbf{Q}^\dagger \mathbf{P} \tag{9}$$

Let us now consider the minimum norm solution, i.e. the one that minimizes $\text{Tr}(\mathbf{W}^\top \mathbf{W})$. Note that $\mathbf{W}^\top \mathbf{W}$ and $\mathbf{Q}^\dagger \mathbf{P} \mathbf{Q}^\dagger$ are positive semidefinite. Let \mathcal{B} be an orthonormal basis of eigenvectors for $\mathbf{colsp}(\mathbf{Q})$, \mathcal{C} an orthonormal basis for $\ker \mathbf{Q}$. Then in the orthonormal basis $\mathcal{B} \cup \mathcal{C}$, we have the following block form of $\mathbf{Q}^\dagger \mathbf{P} \mathbf{Q}^\dagger$

$$\mathbf{Q}^\dagger \mathbf{P} \mathbf{Q}^\dagger = \begin{pmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \tag{10}$$

where \mathbf{A} is positive semidefinite.

Now equation 9 implies that $\mathbf{W} \mathbf{W}^\top$ has the form

$$\mathbf{W}^\top \mathbf{W} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{pmatrix} \tag{11}$$

where \mathbf{C} is also positive semidefinite matrix. Then $\|\mathbf{W}\|_F = \text{Tr}[\mathbf{W}^\top \mathbf{W}]$ is minimized exactly when $\text{Tr}[\mathbf{C}] = 0$. But this holds if and only if $\mathbf{C} = \mathbf{0}$. Now suppose for the sake of contradiction $\mathbf{B} \neq \mathbf{0}$, say $b_{ij} \neq 0$ for some i, j . Then $\mathbf{W}^\top \mathbf{W}$ contains a submatrix

$$\begin{pmatrix} a_{ii} & b_{ij} \\ b_{ij} & 0 \end{pmatrix} \tag{12}$$

which has negative determinant. But this implies that $\mathbf{W}^\top \mathbf{W}$ is not positive semidefinite, a contradiction. We conclude that $\mathbf{B} = \mathbf{0}$ so that the minimum norm solution is precisely

$$\mathbf{W}^{**\top} \mathbf{W}^{**} = \mathbf{Q}^\dagger \mathbf{P} \mathbf{Q}^\dagger.$$

This completes the proof.

□

F. Some Properties of The Covariance Matrices

We assume $\frac{\sigma_\xi^2}{mn} = o(1)$.

With probability $\geq 1 - O(\frac{m^2 n^2}{d})$, we have that $\xi_i^\top \mathbf{v}_k = 0, \forall k, i$ and $\xi_i^\top \xi_j = 0, \forall i, j$. The following discussion focuses on the properties of M, M^+ , and \tilde{M} when this condition is met.

Write $\mathbf{X} = \mathbf{V} \begin{bmatrix} \mathbf{S} \\ \sigma_\xi \mathbf{I}_{mn} \end{bmatrix}$ where $\mathbf{V} = [\mathbf{v}_0, \mathbf{v}_1 \dots \mathbf{v}_K \dots \mathbf{v}_{K+1} \dots \mathbf{v}_{mn+K}]$ where \mathbf{v}_{K+i} is the noise vector selected by example \mathbf{x}_i , and

$$\mathbf{S} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ y_1 \phi_1 & y_2 \phi_1 & \dots & y_{mn} \phi_1 \\ \mu + y_{\text{sub},1} \phi_2 & \mu + y_{\text{sub},2} \phi_2 & \dots & y_{\text{sub},mn} \phi_2 \\ \rho_1 \mathbb{1}_{k_1=3} \phi_3 & \rho_2 \mathbb{1}_{k_2=3} \phi_3 & \dots & \rho_{mn} \mathbb{1}_{k_{mn}=3} \phi_3 \\ \rho_1 \mathbb{1}_{k_1=4} \phi_4 & \rho_2 \mathbb{1}_{k_2=4} \phi_4 & \dots & \rho_{mn} \mathbb{1}_{k_{mn}=4} \phi_4 \\ \vdots & \vdots & \ddots & \vdots \\ \rho_1 \mathbb{1}_{k_1=K} \phi_K & \rho_2 \mathbb{1}_{k_2=K} \phi_K & \dots & \rho_{mn} \mathbb{1}_{k_{mn}=K} \phi_K \end{bmatrix} \\ = \mathbf{S}' \tilde{\mathbf{Y}}, \quad (13)$$

where

$$\mathbf{S}' := \begin{bmatrix} \sqrt{mn} & 0 & 0 & 0 & \dots & 0 \\ 0 & \sqrt{mn} \phi_1 & 0 & 0 & \dots & 0 \\ \sqrt{mn} \mu & 0 & \sqrt{mn} \phi_2 & 0 & \dots & 0 \\ 0 & 0 & 0 & \sqrt{\frac{mn}{K-2}} \phi_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \sqrt{\frac{mn}{K-2}} \phi_K \end{bmatrix},$$

and

$$\tilde{\mathbf{Y}} := \begin{bmatrix} \frac{1}{\sqrt{mn}} & \frac{1}{\sqrt{mn}} & \dots & \frac{1}{\sqrt{mn}} \\ y_1 \frac{1}{\sqrt{mn}} & y_2 \frac{1}{\sqrt{mn}} & \dots & y_{mn} \frac{1}{\sqrt{mn}} \\ y_{\text{sub},1} \frac{1}{\sqrt{mn}} & y_{\text{sub},2} \frac{1}{\sqrt{mn}} & \dots & y_{\text{sub},mn} \frac{1}{\sqrt{mn}} \\ \rho_1 \mathbb{1}_{k_1=3} \sqrt{\frac{K-2}{mn}} & \rho_2 \mathbb{1}_{k_2=3} \sqrt{\frac{K-2}{mn}} & \dots & \rho_{mn} \mathbb{1}_{k_{mn}=3} \sqrt{\frac{K-2}{mn}} \\ \rho_1 \mathbb{1}_{k_1=4} \sqrt{\frac{K-2}{mn}} & \rho_2 \mathbb{1}_{k_2=4} \sqrt{\frac{K-2}{mn}} & \dots & \rho_{mn} \mathbb{1}_{k_{mn}=4} \sqrt{\frac{K-2}{mn}} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_1 \mathbb{1}_{k_1=K} \sqrt{\frac{K-2}{mn}} & \rho_2 \mathbb{1}_{k_2=K} \sqrt{\frac{K-2}{mn}} & \dots & \rho_{mn} \mathbb{1}_{k_{mn}=K} \sqrt{\frac{K-2}{mn}} \end{bmatrix}$$

It should be noted that the rows of $\tilde{\mathbf{Y}}$ are orthonormal due to the assumption of a balanced dataset. Consequently, to obtain the singular value decomposition (SVD) of \mathbf{S} , it suffices to find the SVD of $\mathbf{S}' = \mathbf{P}' \mathbf{\Lambda}' \mathbf{Q}'^\top$. Moreover, the right singular vectors of \mathbf{S} with non-zero singular values are given by the rows of $\mathbf{Q}'^\top \tilde{\mathbf{Y}}$.

We write M as $\mathbf{V} \mathbf{G} \mathbf{V}^\top$ where \mathbf{G} is given by

$$\begin{bmatrix} \frac{1}{mn} \mathbf{S} \mathbf{S}^\top & \frac{\sigma_\xi}{mn} \mathbf{S} \\ \frac{\sigma_\xi}{mn} \mathbf{S}^\top & \frac{\sigma_\xi^2}{mn} \mathbf{I}_{mn} \end{bmatrix}.$$

Now we are ready to show the following lemma which describes the SVD of \mathbf{G} .

Lemma F.1. *Let $\mathbf{S} \in \mathbb{R}^{K \times nm}$ be a rank- K matrix with SVD $\mathbf{P} \mathbf{\Lambda} \mathbf{Q}^\top$, where $\mathbf{P} \in \mathbb{R}^{K \times K}$, $\mathbf{\Lambda} \in \mathbb{R}^{K \times mn}$ and $\mathbf{Q} \in \mathbb{R}^{mn \times mn}$. The mn non-zero eigenvalues of the following matrix \mathbf{G}*

$$\begin{bmatrix} \frac{1}{mn} \mathbf{S} \mathbf{S}^\top & \frac{\sigma_\xi}{mn} \mathbf{S} \\ \frac{\sigma_\xi}{mn} \mathbf{S}^\top & \frac{\sigma_\xi^2}{mn} \mathbf{I}_{mn} \end{bmatrix}$$

are given by $\frac{\sigma_\xi^2}{mn} + \frac{\lambda_1^2}{mn}, \frac{\sigma_\xi^2}{mn} + \frac{\lambda_2^2}{mn}, \dots, \frac{\sigma_\xi^2}{mn} + \frac{\lambda_K^2}{mn}, \frac{\sigma_\xi^2}{mn}, \dots, \frac{\sigma_\xi^2}{mn}$, with the corresponding eigenvectors $\begin{bmatrix} \frac{1}{\sqrt{1+r_1^2}} \mathbf{p}_1 \\ \frac{r_1}{\sqrt{1+r_1^2}} \mathbf{q}_1 \end{bmatrix}, \begin{bmatrix} \frac{1}{\sqrt{1+r_2^2}} \mathbf{p}_2 \\ \frac{r_2}{\sqrt{1+r_2^2}} \mathbf{q}_2 \end{bmatrix}, \dots, \begin{bmatrix} \frac{1}{\sqrt{1+r_K^2}} \mathbf{p}_K \\ \frac{r_K}{\sqrt{1+r_K^2}} \mathbf{q}_K \end{bmatrix}, \begin{bmatrix} \mathbf{0}_K \\ \mathbf{q}_{K+1} \end{bmatrix}, \dots, \begin{bmatrix} \mathbf{0}_K \\ \mathbf{q}_{mn} \end{bmatrix}$, where $r_k = \frac{\sigma_\xi}{\lambda_k}$.

Proof. Let $\begin{bmatrix} \mathbf{P}\mathbf{a} \\ \mathbf{Q}\mathbf{b} \end{bmatrix}$ where $\mathbf{a} \in \mathbb{R}^K$ and $\mathbf{b} \in \mathbb{R}^{mn}$ be an eigenvector of \mathbf{G} . By the definition of eigenvector there should exist α such that $\mathbf{G} \begin{bmatrix} \mathbf{P}\mathbf{a} \\ \mathbf{Q}\mathbf{b} \end{bmatrix} = \alpha \begin{bmatrix} \mathbf{P}\mathbf{a} \\ \mathbf{Q}\mathbf{b} \end{bmatrix}$, i.e.,

$$\begin{cases} \frac{1}{mn} \mathbf{P}\mathbf{\Lambda}\mathbf{\Lambda}^\top \mathbf{a} + \frac{\sigma_\xi}{mn} \mathbf{P}\mathbf{\Lambda}\mathbf{b} = \alpha \mathbf{P}\mathbf{a} \\ \frac{\sigma_\xi}{mn} \mathbf{Q}\mathbf{\Lambda}^\top \mathbf{a} + \frac{\sigma_\xi^2}{mn} \mathbf{Q}\mathbf{b} = \alpha \mathbf{Q}\mathbf{b}, \end{cases}$$

which reduces to

$$\begin{cases} (\alpha \mathbf{I}_K - \frac{1}{mn} \mathbf{\Lambda}\mathbf{\Lambda}^\top) \mathbf{a} = \frac{\sigma_\xi}{mn} \mathbf{\Lambda}\mathbf{b} \\ \frac{\sigma_\xi}{mn} \mathbf{\Lambda}^\top \mathbf{a} = (\alpha - \frac{\sigma_\xi^2}{mn}) \mathbf{b}. \end{cases}$$

Firstly, we observe that the rank of \mathbf{G} is at most mn because $\mathbf{G} = \frac{1}{mn} \begin{bmatrix} \mathbf{S} \\ \sigma_\xi \mathbf{I}_{mn} \end{bmatrix} \begin{bmatrix} \mathbf{S} \\ \sigma_\xi \mathbf{I}_{mn} \end{bmatrix}^\top$. Then it is easy to check that the eigenvalues and eigenvectors in Lemma F.1 satisfy the above conditions and the eigenvectors are indeed orthonormal, which completes the proof. \square

Corollary F.2. *The projection of \mathbf{v}_2 onto $\ker \mathbf{M}$ has magnitude $\Theta(\frac{\sigma_\xi}{\sqrt{mn}})$.*

Corollary F.3. *Assuming the dataset is balanced, then*

$$\sqrt{\mathbf{v}_2^\top \mathbf{M}^\dagger \mathbf{M} \mathbf{v}_2} = \begin{cases} 0, & \text{if } \mu = 0 \\ O(\frac{\sigma_\xi}{\sqrt{mn}}), & \text{if } \mu \neq 0 \text{ and } \mu = \Theta(1). \end{cases}$$

Proof. Let $\mathbf{L}\mathbf{A}\mathbf{L}^\top$ be the eigendecomposition of \mathbf{G} . Then

$$\mathbf{M}^\dagger \mathbf{v}_2 = \mathbf{V}\mathbf{L}\mathbf{A}^\dagger \mathbf{L}^\top \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

When $\mu = 0$, we can express the SVD of \mathbf{S} (equation 13) and apply Lemma F.1 to obtain the following result.

$$\lambda_3 = \sqrt{mn}\phi_2, \quad a_3 = \frac{\sigma_\xi^2}{mn} + \phi_2^2, \quad r_3 = \frac{\sigma_\xi}{\sqrt{mn}\phi_2}$$

$$\mathbf{p}_k = \mathbf{e}_k, \quad \forall k \in [K] \quad \text{and} \quad \mathbf{q}_3 = \begin{bmatrix} \frac{1}{\sqrt{mn}} y_{\text{sub},1} \\ \frac{1}{\sqrt{mn}} y_{\text{sub},2} \\ \vdots \\ \frac{1}{\sqrt{mn}} y_{\text{sub},mn} \end{bmatrix}, \quad \text{and} \quad \mathbf{l}_3 = \begin{bmatrix} \frac{1}{\sqrt{1+r_3^2}} \mathbf{p}_3 \\ \frac{r_3}{\sqrt{1+r_3^2}} \mathbf{q}_3 \end{bmatrix}.$$

Thus

$$\begin{aligned} \mathbf{M}^\dagger \mathbf{v}_2 &= \frac{1}{a_3 \sqrt{1+r_3^2}} \mathbf{V}\mathbf{l}_3 \\ &= \frac{1}{a_3 \sqrt{1+r_3^2}} \left(\frac{1}{\sqrt{1+r_3^2}} \mathbf{v}_2 + \frac{r_3}{\sqrt{1+r_3^2}} \sum_{i=1}^{mn} \frac{1}{\sqrt{mn}} y_{\text{sub},i} \mathbf{v}_{K+i} \right). \end{aligned}$$

Let \bar{x}_y be the average of examples with label y and let \mathcal{S}_y collect indices of examples with label y . Then

$$\bar{x}_y = \mathbf{v}_0 + \mathbf{v}_1 + \frac{2\sigma_\xi}{mn} \sum_{i \in \mathcal{S}_y} \mathbf{v}_{K+i}, \quad (14)$$

and

$$\bar{x}_y^\top \mathbf{M}^\dagger \mathbf{v}_2 = \frac{r_3}{a_3(1+r_3^2)} \frac{2\sigma_\xi}{mn} \sum_{i \in \mathcal{S}_y} \frac{1}{\sqrt{mn}} y_{\text{sub},i} = 0.$$

Write \mathbf{M}^+ as

$$\mathbf{M}^+ = \frac{1}{2}(\bar{\mathbf{x}}_{+1}\bar{\mathbf{x}}_{+1}^\top + \bar{\mathbf{x}}_{-1}\bar{\mathbf{x}}_{-1}^\top).$$

Then

$$\mathbf{v}_2^\top \mathbf{M}^\dagger \mathbf{M}^+ \mathbf{M}^\dagger \mathbf{v}_2 = \frac{1}{2}((\mathbf{v}_2^\top \mathbf{M}^\dagger \bar{\mathbf{x}}_{+1})^2 + (\mathbf{v}_2^\top \mathbf{M}^\dagger \bar{\mathbf{x}}_{-1})^2) = 0.$$

When $\mu \neq 0$, then there are at most two of \mathbf{p}_k 's that are not orthogonal to \mathbf{e}_3 (say \mathbf{p}_1 and \mathbf{p}_3). Additionally, all of their elements, except for the first one, are zero. The remaining corresponding quantities satisfy

$$\begin{aligned} \lambda_1, \lambda_3 &= \Theta(\sqrt{mn}), \\ a_1 &= \frac{\lambda_1^2}{mn} + \frac{\sigma_\xi^2}{mn}, \quad a_3 = \frac{\lambda_3^2}{mn} + \frac{\sigma_\xi^2}{mn} \\ r_1 &= \frac{\sigma_\xi}{\lambda_1}, \quad r_3 = \frac{\sigma_\xi}{\lambda_3}, \end{aligned}$$

and \mathbf{q}_1 and \mathbf{q}_3 are just linear combinations of $\bar{\mathbf{y}}_{\text{sub}}$ and $\frac{1}{\sqrt{mn}}\mathbf{1}$, where $\bar{\mathbf{y}}_{\text{sub}}$ is a vector whose i -th element is $\frac{1}{\sqrt{mn}}y_{\text{sub},i}$. Then

$$\mathbf{M}^\dagger \mathbf{v}_2 = \mathbf{V} \begin{bmatrix} \frac{1}{a_1} \frac{1}{\sqrt{1+r_1^2}} c_{3,1} \mathbf{l}_1 + \frac{1}{a_3} \frac{1}{\sqrt{1+r_3^2}} c_{3,3} \mathbf{l}_3 \end{bmatrix}$$

where $c_{i,j} = \mathbf{p}_j^\top \mathbf{e}_i$ are constants. For $i = 0, 2$

$$\begin{aligned} \mathbf{v}_0^\top \mathbf{M}^\dagger \mathbf{v}_2 &= \mathbf{e}_1^\top \begin{bmatrix} \frac{1}{a_1} \frac{1}{\sqrt{1+r_1^2}} c_{3,1} \mathbf{l}_1 + \frac{1}{a_3} \frac{1}{\sqrt{1+r_3^2}} c_{3,3} \mathbf{l}_3 \end{bmatrix} \\ &= \frac{1}{a_1} \frac{1}{1+r_1^2} c_{3,1} c_{1,1} + \frac{1}{a_3} \frac{1}{1+r_3^2} c_{3,3} c_{1,3} \\ &= \left(\frac{mn}{\lambda_1^2} - \Theta\left(\frac{\sigma_\xi^2}{mn}\right)\right) (1 - \Theta\left(\frac{\sigma_\xi}{\lambda_1}\right)) c_{3,1} c_{1,1} + \left(\frac{mn}{\lambda_3^2} - \Theta\left(\frac{\sigma_\xi^2}{mn}\right)\right) (1 - \Theta\left(\frac{\sigma_\xi}{\lambda_3}\right)) c_{3,3} c_{1,3} \\ &= \end{aligned}$$

where $|\epsilon_1| = O\left(\frac{\sigma_\xi}{\sqrt{mn}}\right)$. Similarly,

$$\mathbf{v}_2^\top \mathbf{M}^\dagger \mathbf{v}_2 = \frac{mn}{\lambda_1^2} c_{3,1} c_{3,1} + \frac{mn}{\lambda_3^2} c_{3,3} c_{3,3} + \epsilon_2,$$

where $|\epsilon_2| = O\left(\frac{\sigma_\xi}{\sqrt{mn}}\right)$. For $i > K$

$$\begin{aligned} \mathbf{v}_i^\top \mathbf{M}^\dagger \mathbf{v}_2 &= \mathbf{v}_i \mathbf{V} \begin{bmatrix} \frac{1}{a_1} \frac{1}{1+r_1^2} c_{3,1} \mathbf{l}_1 + \frac{1}{a_3} \frac{1}{1+r_3^2} c_{3,3} \mathbf{l}_3 \end{bmatrix} \\ &= \mathbf{e}_i^\top \begin{bmatrix} \frac{1}{a_1} \frac{1}{1+r_1^2} c_{3,1} \mathbf{l}_1 + \frac{1}{a_3} \frac{1}{1+r_3^2} c_{3,3} \mathbf{l}_3 \end{bmatrix} \\ &= \epsilon_3, \end{aligned}$$

1045 where $|\epsilon_3| = O(\frac{\sigma_\xi}{mn})$. Additionally,

$$1046 \quad \bar{\mathbf{x}}_y = \mathbf{v}_0 + \mathbf{v}_1 + \mu \mathbf{v}_2 + \frac{2\sigma_\xi}{mn} \sum_{i \in \mathcal{S}_y} \mathbf{v}_{K+i}.$$

1050 Then

$$1051 \quad \bar{\mathbf{x}}_y^\top \mathbf{M}^\dagger \mathbf{v}_2 = \frac{mn}{\lambda_1^2} c_{3,1} c_{1,1} + \frac{mn}{\lambda_3^2} c_{3,3} c_{1,3} + \frac{mn}{\lambda_1^2} c_{3,1} c_{3,1} + \frac{mn}{\lambda_3^2} c_{3,3} c_{3,3} + O(\frac{\sigma_\xi}{\sqrt{mn}}).$$

1054 By straightforward calculation, we can verify that $\frac{mn}{\lambda_1^2} c_{3,1} c_{1,1} + \frac{mn}{\lambda_3^2} c_{3,3} c_{1,3} + \frac{mn}{\lambda_1^2} c_{3,1} c_{3,1} + \frac{mn}{\lambda_3^2} c_{3,3} c_{3,3} = 0$. This equation
1055 can be equivalently examined as the satisfaction of the following condition:

$$1056 \quad [1 \ \mu] \begin{bmatrix} 1 & \mu \\ \mu & \mu^2 + \phi_2^2 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = 0.$$

1060 Therefore $|\bar{\mathbf{x}}_y^\top \mathbf{M}^\dagger \mathbf{v}_2| = O(\frac{\sigma_\xi}{\sqrt{mn}})$, and consequently $\sqrt{\mathbf{v}_2^\top \mathbf{M}^\dagger \mathbf{M} + \mathbf{M}^\dagger \mathbf{v}_2} = O(\frac{\sigma_\xi}{\sqrt{mn}})$. \square

1062 **Corollary F.4.** *Similar to Corollary F.3, we also have $\sqrt{\mathbf{v}_k^\top \mathbf{M}^\dagger \mathbf{M} + \mathbf{M}^\dagger \mathbf{v}_k} = 0$, $k = 3, 4, \dots, K$.*

1064 **Corollary F.5.** $\sqrt{\mathbf{v}_1^\top \mathbf{M}^\dagger \mathbf{M} + \mathbf{M}^\dagger \mathbf{v}_1} = \Theta(1)$. *It can be proved using the same strategy as in Corollary F.3.*

1066 **Lemma F.6.** (1) *The first K eigenvectors/eigenvalues of $\tilde{\mathbf{M}}$ match those of \mathbf{M} . (2) $\mathbf{M}^\dagger \tilde{\mathbf{M}}$ is identity on $\text{colsp}(\tilde{\mathbf{M}})$ and
1067 null on $\ker(\tilde{\mathbf{M}})$, i.e., $\mathbf{M}^\dagger \tilde{\mathbf{M}} = \tilde{\mathbf{M}}^\dagger \mathbf{M}$.*

1068 *Proof.* We assign indices to the training examples such that the augmented examples from the same original example are
1069 indexed from $(l-1) \times m + 1$ to $l \times m$, where l ranges from 1 to n . Next, we define matrix $\tilde{\mathbf{V}} = [\tilde{\mathbf{v}}_1, \tilde{\mathbf{v}}_2, \dots, \tilde{\mathbf{v}}_n] \in \mathbb{R}^{d \times n}$,
1070 where

$$1071 \quad \tilde{\mathbf{v}}_i = \mathbf{v}_i, \quad \forall 1 \leq i \leq K,$$

$$1072 \quad \tilde{\mathbf{v}}_i = \frac{1}{\sqrt{m}} \sum_{j=1}^m \mathbf{v}_{K+(i-1) \times m + j}, \quad \forall K+1 \leq i \leq n.$$

1077 In other words, $\tilde{\mathbf{V}}$ can be written as

$$1078 \quad \tilde{\mathbf{V}} = \mathbf{V} \mathbf{T},$$

1081 where

$$1082 \quad \mathbf{T} = \begin{bmatrix} \mathbf{I}_K & & & & & \\ & \frac{1}{\sqrt{m}} \mathbf{1}_{m \times 1} & 0 & 0 & \dots & 0 \\ & 0 & \frac{1}{\sqrt{m}} \mathbf{1}_{m \times 1} & 0 & \dots & 0 \\ & 0 & 0 & \frac{1}{\sqrt{m}} \mathbf{1}_{m \times 1} & \dots & 0 \end{bmatrix}$$

1088 Note that, by the definition of our augmentation, the center of augmentations of the i -th original example, i.e., $\tilde{\mathbf{x}}_i =$
1089 $\frac{1}{m} \sum_{j=1}^m \mathbf{x}_{K+(i-1) \times m + j}$, can be considered as an example with the same features as \mathbf{x}_i but with an added noise term of
1090 $\frac{\sigma_\xi}{\sqrt{m}} \tilde{\mathbf{v}}_i$. Therefore we can change the basis to $\tilde{\mathbf{V}}$ and express $\tilde{\mathbf{M}}$ as

$$1091 \quad \tilde{\mathbf{M}} = \tilde{\mathbf{V}} \tilde{\mathbf{G}} \tilde{\mathbf{V}}^\top,$$

1094 where

$$1095 \quad \tilde{\mathbf{G}} = \begin{bmatrix} \frac{1}{n} \tilde{\mathbf{S}} \tilde{\mathbf{S}}^\top & \frac{1}{n} \frac{\sigma_\xi}{\sqrt{m}} \tilde{\mathbf{S}} \\ \frac{1}{n} \frac{\sigma_\xi}{\sqrt{m}} \tilde{\mathbf{S}}^\top & \frac{1}{n} \frac{\sigma_\xi^2}{m} \mathbf{I}_n \end{bmatrix}$$

1100 and

1101

1102

1103

1104 where

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115 and

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

$$\tilde{\mathbf{S}} = \tilde{\mathbf{S}}' \bar{\mathbf{Y}}_{\text{orig}}, \quad (15)$$

$$\tilde{\mathbf{S}}' := \begin{bmatrix} \sqrt{n} & 0 & 0 & 0 & \dots & 0 \\ 0 & \sqrt{n}\phi_1 & 0 & 0 & \dots & 0 \\ \sqrt{n}\mu & 0 & \sqrt{n}\phi_2 & 0 & \dots & 0 \\ 0 & 0 & 0 & \sqrt{\frac{n}{K-2}}\phi_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \sqrt{\frac{n}{K-2}}\phi_K \end{bmatrix},$$

$$\bar{\mathbf{Y}}_{\text{orig}} := \begin{bmatrix} \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \dots & \frac{1}{\sqrt{n}} \\ y_1 \frac{1}{\sqrt{n}} & y_2 \frac{1}{\sqrt{n}} & \dots & y_n \frac{1}{\sqrt{n}} \\ y_{\text{sub},1} \frac{1}{\sqrt{n}} & y_{\text{sub},2} \frac{1}{\sqrt{n}} & \dots & y_{\text{sub},n} \frac{1}{\sqrt{n}} \\ \rho_1 \mathbb{1}_{k_1=3} \sqrt{\frac{K-2}{n}} & \rho_2 \mathbb{1}_{k_2=3} \sqrt{\frac{K-2}{n}} & \dots & \rho_n \mathbb{1}_{k_n=3} \sqrt{\frac{K-2}{n}} \\ \rho_1 \mathbb{1}_{k_1=4} \sqrt{\frac{K-2}{n}} & \rho_2 \mathbb{1}_{k_2=4} \sqrt{\frac{K-2}{n}} & \dots & \rho_n \mathbb{1}_{k_n=4} \sqrt{\frac{K-2}{n}} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_1 \mathbb{1}_{k_1=K} \sqrt{\frac{K-2}{n}} & \rho_2 \mathbb{1}_{k_2=K} \sqrt{\frac{K-2}{n}} & \dots & \rho_n \mathbb{1}_{k_n=K} \sqrt{\frac{K-2}{n}} \end{bmatrix}_{\text{orig}}. \quad (16)$$

We note that we use the subscript ‘orig’ of a matrix to indicate that its elements represent the corresponding quantities on the original dataset (e.g., y_i is the label of the i -th original example). Let $\tilde{\mathbf{P}}' \tilde{\mathbf{\Lambda}}' \tilde{\mathbf{Q}}'^{\top}$ be the SVD of $\tilde{\mathbf{S}}$. Similar to equation 13, we observe that $\tilde{\mathbf{P}}' \tilde{\mathbf{\Lambda}}' (\tilde{\mathbf{Q}}'^{\top} \bar{\mathbf{Y}}_{\text{orig}})$ serves as an eigendecomposition of $\tilde{\mathbf{S}}$.

Now we make the following observations:

1. By Lemma F.1 (with \mathbf{G} replaced by $\tilde{\mathbf{G}}$) and the fact that $\tilde{\mathbf{\Lambda}}'$ collects the eigenvalues of $\tilde{\mathbf{S}}$, the eigenvalues of $\tilde{\mathbf{G}}$ are $\frac{\sigma_{\xi}^2}{mn} + \frac{\lambda_1^2}{n}, \frac{\sigma_{\xi}^2}{mn} + \frac{\lambda_2^2}{n}, \dots, \frac{\sigma_{\xi}^2}{mn} + \frac{\lambda_K^2}{n}, \frac{\sigma_{\xi}^2}{n}, \dots, \frac{\sigma_{\xi}^2}{n}$, which are also the eigenvalues of $\tilde{\mathbf{M}}$ because $\tilde{\mathbf{V}}$ has orthonormal columns. With the observation that $\tilde{\mathbf{S}}' = \frac{1}{\sqrt{m}} \mathbf{S}'$ (\mathbf{S}' is defined in equation 13), we further conclude that the above eigenvalues equal eigenvalues of \mathbf{G} and therefore \mathbf{M} .

2. Let $\tilde{\mathbf{q}}_i$ be the i -th column of $\bar{\mathbf{Y}}_{\text{orig}}^{\top} \tilde{\mathbf{Q}}'$. By Lemma F.1 (substitute \mathbf{G} with $\tilde{\mathbf{G}}$), the i -th ($i \leq K$) eigenvector of $\tilde{\mathbf{G}}$ is given by $\begin{bmatrix} \frac{1}{\sqrt{1+\tilde{r}_i^2}} \tilde{\mathbf{p}}'_i \\ \frac{\tilde{r}_i}{1+\tilde{r}_i^2} \tilde{\mathbf{q}}_i \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{1+\tilde{r}_i^2}} \mathbf{p}_i \\ \frac{\tilde{r}_i}{1+\tilde{r}_i^2} \tilde{\mathbf{q}}_i \end{bmatrix}$, where $\tilde{r}_i = \frac{\sigma_{\xi}}{\sqrt{m}\lambda'_i} = \frac{\sigma_{\xi}}{\lambda_i}$. The corresponding eigenvector of $\tilde{\mathbf{M}}$ is $\mathbf{V}\mathbf{T} \begin{bmatrix} \frac{1}{\sqrt{1+\tilde{r}_i^2}} \mathbf{p}_i \\ \frac{\tilde{r}_i}{1+\tilde{r}_i^2} \tilde{\mathbf{q}}_i \end{bmatrix}$.

Observe that $\mathbf{T}\bar{\mathbf{Y}}_{\text{orig}}^{\top} = \bar{\mathbf{Y}}^{\top}$, therefore $\mathbf{V}\mathbf{T} \begin{bmatrix} \frac{1}{\sqrt{1+\tilde{r}_i^2}} \mathbf{p}_i \\ \frac{\tilde{r}_i}{1+\tilde{r}_i^2} \tilde{\mathbf{q}}_i \end{bmatrix} = \mathbf{V} \begin{bmatrix} \frac{1}{\sqrt{1+\tilde{r}_i^2}} \mathbf{p}_i \\ \frac{\tilde{r}_i}{1+\tilde{r}_i^2} \mathbf{q}_i \end{bmatrix}$ which is the i -th eigenvector of \mathbf{M} .

Combining the above two leads to the conclusion that the first K eigenvectors/eigenvalues of $\tilde{\mathbf{M}}$ and \mathbf{M} match. Additionally, we observe that $\text{colsp}(\tilde{\mathbf{M}}) \subseteq \text{colsp}(\mathbf{M})$. Therefore the span of the last $n - K$ eigenvectors of $\tilde{\mathbf{M}}$ is a subspace of the span of the last $mn - K$ eigenvectors of \mathbf{M} . Since Lemma F.1 tells us that the remaining $mn - K$ eigenvalues of \mathbf{M} are equal, \mathbf{M} is identity on the span of the last $mn - K$ eigenvectors. Thus $\tilde{\mathbf{M}}$ is identity on the span of the last $n - K$ eigenvectors of $\tilde{\mathbf{M}}$. Now we can conclude that $\mathbf{M}^{\dagger} \tilde{\mathbf{M}} = \tilde{\mathbf{M}}^{\dagger} \mathbf{M}$. \square

Lemma F.7. Suppose that the first $\frac{mn}{2}$ examples have class label +1 and the others have class label -1. Let $\mathbf{L}^+ \mathbf{A}^+ \mathbf{L}^{+\top}$

(where $\mathbf{L}^+ \in \mathbb{R}^{d \times 2}$) be the eigendecomposition of \mathbf{M}^+ , then

$$\begin{aligned}
 \mathbf{l}_1^+ &= \mathbf{V} \begin{bmatrix} \frac{1}{\sqrt{1+\mu^2+\frac{\sigma_\xi^2}{mn}}} \\ 0 \\ \frac{\mu}{\sqrt{1+\mu^2+\frac{\sigma_\xi^2}{mn}}} \\ \mathbf{0}_{(K-2) \times 1} \\ \frac{\sigma_\xi}{mn\sqrt{1+\mu^2+\frac{\sigma_\xi^2}{mn}}} \mathbf{1}_{\frac{mn}{2} \times 1} \\ \frac{\sigma_\xi}{mn\sqrt{1+\mu^2+\frac{\sigma_\xi^2}{mn}}} \mathbf{1}_{\frac{mn}{2} \times 1} \end{bmatrix}, \quad \mathbf{l}_2^+ = \mathbf{V} \begin{bmatrix} 0 \\ \frac{\phi_1}{\sqrt{\phi_1^2+\frac{\sigma_\xi^2}{mn}}} \\ 0 \\ \mathbf{0}_{(K-2) \times 1} \\ \frac{\sigma_\xi}{mn\sqrt{\phi_1^2+\frac{\sigma_\xi^2}{mn}}} \mathbf{1}_{\frac{mn}{2} \times 1} \\ \frac{-\sigma_\xi}{mn\sqrt{\phi_1^2+\frac{\sigma_\xi^2}{mn}}} \mathbf{1}_{\frac{mn}{2} \times 1} \end{bmatrix}, \\
 a_1 &= 1 + \mu^2 + \frac{\sigma_\xi^2}{mn}, \quad a_2 = \phi_1^2 + \frac{\sigma_\xi^2}{mn}
 \end{aligned} \tag{17}$$

G. Class Collapse in Supervised CL

G.1. Proof of Theorem C.3

Let \mathbf{l}_\perp be the projection of \mathbf{v}_2 onto $\ker \mathbf{M}$. By Corollary F.2, $\|\mathbf{l}_\perp\| = \Theta(\frac{\sigma_\xi}{\sqrt{mn}})$. Let $\mathbf{a} = \frac{mn}{\sigma_\xi^2} \mathbf{l}_\perp$. We can construct a \mathbf{W}^* that satisfies the following

$$\mathbf{W}^{*\top} \mathbf{W}^* = \mathbf{M}^\dagger \mathbf{M}^+ \mathbf{M}^\dagger + \mathbf{a} \mathbf{a}^\top,$$

which, by Lemma E.2, satisfies the condition for being a minimizer of the loss. In the meantime, \mathbf{W}^* also satisfies $\|\mathbf{W}^* \mathbf{v}_2\| = \Theta(1)$ by Corollary F.2. Note that both \mathbf{v}_2 and the projection of \mathbf{v}_2 onto $\mathbf{colsp}(\mathbf{M})$ is orthogonal to \mathbf{v}_k ($1 \leq k \leq K, k \neq 2$) as well as \mathbf{v}_k ($k > mn$) by Lemma F.1, therefore

$$\mathbf{l}_\perp \text{ is also orthogonal to } \mathbf{v}_k, \text{ for any } k \text{ s.t. } 1 \leq k \leq K, k \neq 2 \text{ and } k > mn. \tag{18}$$

Then, for \mathbf{x} from $\mathcal{D}_{\text{orig}}$ the following holds true

$$\mathbf{W}^* \mathbf{x} = c_0 \mathbf{v}_0 + c_1 y \mathbf{v}_1 + y_{\text{sub}} c_2 \mathbf{v}_2 + \mathbf{h}_x + \mathbf{W}^* \boldsymbol{\xi},$$

where c_1, c_2 are $\Theta(1)$, and \mathbf{h}_x is orthogonal to $\mathbf{v}_k, k = 0, \dots, K$ and $\mathbf{h}_x \in \mathbf{colsp}(\mathbf{M})$ (by Lemmas F.1, F.5, F.4, equation 18 and that $\|\mathbf{W}^* \mathbf{v}_2\| = \Theta(1)$). Let $\boldsymbol{\beta} = c_2 \mathbf{v}_2$, then

$$\boldsymbol{\beta}^\top \mathbf{W}^* \mathbf{x} = y_{\text{sub}} c_2^2 + \boldsymbol{\beta}^\top \mathbf{W}^* \boldsymbol{\xi}.$$

With probability $\geq 1 - \frac{mn}{d}$, $\boldsymbol{\xi} \notin \{\mathbf{v}_k\}_{k=1}^{mn}$, which indicates that $\mathbf{W}^* \boldsymbol{\xi} = 0$ by Lemma F.1 and equation 18. Therefore we can conclude

$$\Pr(y_{\text{sub}} \boldsymbol{\beta}^\top \mathbf{W}^* \mathbf{x} > 0 | y) \geq 1 - \frac{mn}{d} = 1 - o(1).$$

G.2. Proof of Theorems C.4 and C.7

Theorems C.4 and C.7 and can be proved by invoking Lemma E.2 and Corollary F.3.

H. Feature Suppression in Unsupervised CL

H.1. Feature Suppression 1

By Lemmas E.2 and F.6, when $p < K$, any global minimizer of \mathcal{L}_{UCL} satisfies

$$\mathbf{W}^\top \mathbf{W} \mathbf{M} = \sum_{i=1}^p \mathbf{r}_i \mathbf{r}_i^\top, \tag{19}$$

1210 where $\{\mathbf{r}_i\}_{i=1}^p$ can be an orthonormal basis of any p -dimensional subspace of $\text{colsp}(\tilde{M})$. By equation 13 and Lemmas F.1
 1211 and F.6, M and \tilde{M} each have an eigenvector \mathbf{c}_1 with eigenvalue $\frac{\sigma_\xi^2}{mn} + \phi_1^2$ and a $\frac{1}{\sqrt{1 + \frac{\sigma_\xi^2}{mn\phi_1^2}}}$ alignment with \mathbf{v}_1 , with the
 1212
 1213 other eigenvectors having no alignment with \mathbf{v}_1 . Thus if we include \mathbf{c}_1 in $\{\mathbf{r}_i\}_{i=1}^p$ and let $\mathbf{W}^\top \mathbf{W}$ be null on $\ker M$, then
 1214 the constructed \mathbf{W} is a minimizer of \mathcal{L}_{UCL} with $\Theta(1)$ alignment with \mathbf{v}_1 . Now let's look at the minimum norm minimizer,
 1215 which should satisfy
 1216

$$1217 \mathbf{W}^\top \mathbf{W} = \sum_{i=1}^p \mathbf{r}_i \mathbf{r}_i^\top M^\dagger,$$

1221 where $\{\mathbf{r}_i\}_{i=1}^p$ is selected such that \mathbf{W} has the smallest norm. By Lemma F.6, $\{\mathbf{r}_i\}_{i=1}^p$ should be the p -
 1222 eigenvectors of M with largest eigenvalues (so that the inverse of the eigenvalues are among the smallest). If among
 1223 $\frac{(1+\mu^2+\phi_2^2)+\sqrt{(1+\mu^2+\phi_2^2)^2-4\phi_2^2}}{2}$, $\frac{(1+\mu^2+\phi_2^2)-\sqrt{(1+\mu^2+\phi_2^2)^2-4\phi_2^2}}{2}$, $\frac{\phi_3}{\sqrt{K-2}}$, \dots , $\frac{\phi_K}{\sqrt{K-2}}$ there are p elements larger than ϕ_1 , then
 1224 $\frac{\sigma_\xi^2}{mn} + \phi_1^2$ is not among the p largest eigenvalues of M . Thus \mathbf{c}_1 is not included in $\{\mathbf{r}_i\}_{i=1}^p$ and the corresponding \mathbf{W} is
 1225 orthogonal to \mathbf{v}_1 .
 1226

1227 H.2. Feature Suppression 2

1228 We first present our result under slightly technical conditions.

1230 **Lemma H.1.** *Let $\mathbf{v}_1, \dots, \mathbf{v}_C \in \mathcal{R}^d$ be nonzero and orthogonal, U, A are subspaces that are orthogonal to each other and
 1231 all the \mathbf{v}_i . Suppose we have a data distribution $\mathcal{D} = \{(\mathbf{v}_{y_i} + \mathbf{u}_{y_i} + \mathbf{a}_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \{1, \dots, C\}$, where $\mathbf{u}_i \in U$, $\mathbf{a}_i \in A$
 1232 for all $i \in \{1, \dots, n\}$ (namely all examples in the same class c share the same \mathbf{v}_c and \mathbf{u}_c).
 1233*

1234 Denote $\mathbf{z}_{y_i} = \mathbf{v}_{y_i} + \mathbf{u}_{y_i}$, and let M, M^+ be the matrices defined for this dataset, and let \mathbf{Z}, \mathbf{Z}^+ and \mathbf{A}, \mathbf{A}^+ be the
 1235 corresponding matrices when the data is $\{(\mathbf{z}_{y_i}, y_i)\}$ and $\{(\mathbf{a}_i, y_i)\}$, respectively. Suppose that $(\mathbf{A} - \mathbf{A}^+)\mathbf{v} \neq \mathbf{0}$ for all
 1236 $\mathbf{v} \in \mathbb{R}^d$ s.t. $\mathbf{A}\mathbf{v} \neq \mathbf{0}$ and the output dimension $p \geq C$. Then $\mathbf{W}^\top \mathbf{W} = \mathbf{Z}^\dagger$ is the minimum norm solution to the contrastive
 1237 learning objective on \mathcal{D} .
 1238

1239 *Proof.* In this proof, we will use \mathbb{E} to represent the empirical expectation over the dataset \mathcal{D} . Also, let n_c denote the number
 1240 of examples in class c .
 1241

1242 We first derive the following expression for \mathbf{A}^+ :

$$1243 \mathbf{A}^+ = \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] \mathbf{Z}^\dagger \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] \quad (20)$$

1244 Define $\mathbf{B} = [\sqrt{n_1} \mathbf{z}_1, \dots, \sqrt{n_C} \mathbf{z}_C] \in \mathbb{R}^{d \times C}$, $\mathbf{C} = [\mathbf{a}_1^*, \dots, \mathbf{a}_C^*] \in \mathbb{R}^{d \times C}$, where n_c is the number of examples in class c
 1245 and $\mathbf{a}_c^* = \frac{1}{n_c} \sum_{y_i=c} \mathbf{a}_i$. Then
 1246

$$1247 \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] \mathbf{Z}^\dagger \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] = \frac{1}{n} \mathbf{C} \mathbf{B}^\top \left(\frac{1}{n} \mathbf{B} \mathbf{B}^\top \right)^\dagger \frac{1}{n} \mathbf{B} \mathbf{C}^\top \quad (21)$$

1248 Now \mathbf{B} has full column rank, so $\mathbf{B}^\top (\mathbf{B} \mathbf{B}^\top)^\dagger \mathbf{B} = \mathbf{I}$. Thus

$$1249 \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] \mathbf{Z}^\dagger \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] = \frac{1}{n} \mathbf{C} \mathbf{C}^\top \quad (22)$$

$$1250 = \sum_{c=1}^C \frac{n_c}{n} \mathbb{E}_{y_i=c}[\mathbf{a}_i] \mathbb{E}_{y_i=c}[\mathbf{a}_i]^\top \quad (23)$$

$$1251 = \mathbf{A}^+ \quad (24)$$

1252 Now we show that $\mathbf{W}^\top \mathbf{W} = \mathbf{Z}^\dagger$ is a global minimizer. It suffices to show that $M \mathbf{W}^\top \mathbf{W} M = M^+$. Note that by

assumption, we have $\langle \mathbf{z}_i, \mathbf{a}_j \rangle = 0$ for all $i \in \{1, \dots, C\}, j \in \{1, \dots, n\}$, so we have

$$\mathbf{M}\mathbf{Z}^\dagger\mathbf{M} = \mathbb{E}[(\mathbf{z}_{y_i} + \mathbf{a}_i)(\mathbf{z}_{y_i} + \mathbf{a}_i)^\top] \mathbf{M}_*^\dagger \mathbb{E}[(\mathbf{z}_{y_i} + \mathbf{a}_i)(\mathbf{z}_{y_i} + \mathbf{a}_i)^\top] \quad (25)$$

$$= (\mathbf{Z} + \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] + \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] + \mathbf{A}) \mathbf{Z}^\dagger (\mathbf{Z} + \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] + \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] + \mathbf{A}) \quad (26)$$

$$= \mathbf{Z}\mathbf{Z}^\dagger\mathbf{Z} + \mathbf{Z}\mathbf{Z}^\dagger\mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] + \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] \mathbf{Z}^\dagger\mathbf{Z} + \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] \mathbf{Z}^\dagger \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] \quad (27)$$

$$= \mathbf{Z} + \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] + \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] + \mathbf{A}^+ \quad (28)$$

$$= \mathbf{Z}^+ + \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] + \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] + \mathbf{A}^+ \quad (29)$$

$$= \mathbf{M}^+ \quad (30)$$

We now want to show that this is the minimum norm solution. It is sufficient to show that $\text{im}(\mathbf{W}^\top \mathbf{W}) = \text{im}(\mathbf{Z}^\dagger) = \text{im}(\mathbf{Z}) \subset \text{im}(\mathbf{M})$. Note that $\text{im}(\mathbf{M}) \subset \text{im}(\mathbf{A}) \oplus \text{im}(\mathbf{Z})$, so we can restrict \mathbf{M} to this subspace. We will show that \mathbf{M} is invertible on $\text{im}(\mathbf{A}) \oplus \text{im}(\mathbf{Z})$. Suppose $\mathbf{v} = \mathbf{z} + \mathbf{a}$ with $\mathbf{z} \in \text{im}(\mathbf{Z}), \mathbf{a} \in \text{im}(\mathbf{A}), \mathbf{M}\mathbf{v} = 0$. This implies that

$$\mathbf{Z}\mathbf{z} + \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] \mathbf{a} = \mathbf{0} \quad (31)$$

$$\mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] \mathbf{z} + \mathbf{A}\mathbf{a} = \mathbf{0} \quad (32)$$

Left-multiplying the first equation by $\mathbb{E}[\frac{n}{n_{y_i} \|\mathbf{v}_{y_i}\|^2} \mathbf{a}_i \mathbf{v}_{y_i}^\top]$, by orthogonality we have

$$\begin{aligned} \mathbf{0} &= \mathbb{E} \left[\frac{n}{n_{y_i} \|\mathbf{v}_{y_i}\|^2} \mathbf{a}_i \mathbf{v}_{y_i}^\top \right] (\mathbb{E}[\mathbf{z}_{y_i} \mathbf{z}_{y_i}^\top] \mathbf{z} + \mathbb{E}[\mathbf{z}_{y_i} \mathbf{a}_i^\top] \mathbf{a}) \\ &= \mathbb{E} \left[\frac{n}{n_{y_i} \|\mathbf{v}_{y_i}\|^2} \mathbf{a}_i \mathbf{v}_{y_i}^\top \right] (\mathbb{E}[(\mathbf{v}_{y_i} + \mathbf{u}_{y_i}) \mathbf{z}_{y_i}^\top] \mathbf{z} + \mathbb{E}[(\mathbf{v}_{y_i} + \mathbf{u}_{y_i}) \mathbf{a}_i^\top] \mathbf{a}) \\ &= \mathbb{E} \left[\frac{n}{n_{y_i} \|\mathbf{v}_{y_i}\|^2} \mathbf{a}_i \mathbf{v}_{y_i}^\top \right] (\mathbb{E}[\mathbf{u}_{y_i} (\mathbf{z}_{y_i}^\top \mathbf{z} + \mathbf{a}_i^\top \mathbf{a})] + \mathbb{E}[\mathbf{v}_{y_i} (\mathbf{z}_{y_i}^\top \mathbf{z} + \mathbf{a}_i^\top \mathbf{a})]) \\ &= \mathbb{E} \left[\frac{n}{n_{y_i} \|\mathbf{v}_{y_i}\|^2} \mathbf{a}_i \mathbf{v}_{y_i}^\top \right] \mathbb{E}[\mathbf{v}_{y_i} (\mathbf{z}_{y_i}^\top \mathbf{z} + \mathbf{a}_i^\top \mathbf{a})] \\ &= \sum_{c=1}^C \frac{1}{n n_c \|\mathbf{v}_c\|^2} \left(\sum_{y_i=c} \mathbf{a}_i \right) \mathbf{v}_c^\top \mathbf{v}_c \left(n_c \mathbf{z}_c^\top \mathbf{z} + \sum_{y_i=c} \mathbf{a}_i^\top \mathbf{a} \right) \\ &= \sum_{c=1}^C \frac{1}{n n_c} \left(\sum_{y_i=c} \mathbf{a}_i \right) \left(n_c \mathbf{z}_c^\top \mathbf{z} + \sum_{y_i=c} \mathbf{a}_i \mathbf{a} \right) \\ &= \frac{1}{n^2} \sum_{c=1}^C \frac{1}{n} \left(\sum_{y_i=c} \mathbf{a}_i \right) \mathbf{z}_c^\top \mathbf{z} + \frac{1}{n n_c} \left(\sum_{y_i=c} \mathbf{a}_i \right) \left(\sum_{y_i=c} \mathbf{a}_i \right)^\top \mathbf{a} \\ &= \mathbb{E}[\mathbf{a}_i \mathbf{z}_{y_i}^\top] \mathbf{z} + \mathbf{A}^+ \mathbf{a} \end{aligned}$$

Now substituting into the second equation, we find that

$$(\mathbf{A} - \mathbf{A}^+) \mathbf{a} = \mathbf{0} \quad (33)$$

But our assumptions imply that $\mathbf{a} = \mathbf{0}$. Returning to the first equation, we now have $\mathbf{Z}\mathbf{z} = \mathbf{0}$. But since \mathbf{Z} is diagonalizable, \mathbf{Z} must be invertible on its image, hence $\mathbf{z} = \mathbf{0}$. We conclude that $\mathbf{v} = \mathbf{0}$. This completes the proof. \square

We now want to show that we can simplify some of the conditions of the previous lemma to linear independence.

Lemma H.2. *Suppose $d \geq 3n - 2$ and $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ are linearly independent. Then there exists a set of nonzero orthogonal vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ s.t. $\mathbf{x}_i = \mathbf{v}_i + \mathbf{u}_i$ and $\mathbf{v}_i, \mathbf{u}_j$ are orthogonal for all $i, j \in \{1, \dots, n\}$.*

1320 *Proof.* WLOG assume the x_i are contained in the span of the first n basis vectors. The lemma amounts to finding an
 1321 orthonormal matrix $\Omega = \begin{pmatrix} A & B \\ C & D \end{pmatrix}$ s.t.
 1322

$$1324 \begin{pmatrix} A & B \\ C & D \end{pmatrix} \begin{pmatrix} X \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} AX \\ CX \end{pmatrix} = \begin{pmatrix} \Sigma \\ F \end{pmatrix} \quad (34)$$

1326 where Σ is diagonal. Since the x_i are linearly independent, X is invertible, so there exists A' s.t. $A'X$ is diagonal.

1328 We now want to construct a matrix C such that $\begin{pmatrix} A' \\ C' \end{pmatrix}$ has orthogonal columns, all with norm $l > 0$. Note that C' has at
 1329 least $2n - 2$ rows. Set $C'_{11} = 1$, and the remaining entries in the first row so that when considering A and the first row
 1330 of C' , the first column is orthogonal to every other column. Now leave $C'_{21} = 0$, set $C'_{22} = 1$, and fill out the remaining
 1331 entries in the second row so that when considering A and the first two rows of C' , the second column is orthogonal to the
 1332 remaining columns. Note that the first column remains orthogonal to all other columns. Continuing in this fashion, we can
 1333 use the first $n - 1$ rows of C' to guarantee that all n columns are orthogonal. Finally, suppose without loss of generality that
 1334 when considering the A' and the first $n - 1$ rows of C' , the first column has the largest norm l . For each of the remaining
 1335 $n - 1$ rows, set the j th row to have all zero entries except possibly in the $(j + 1)$ -th column, which is set so that the j th
 1336 column will also have norm l . Note that the columns remain orthogonal under this construction.

1337 Now $\frac{1}{l} \begin{pmatrix} A' \\ C' \end{pmatrix}$ has orthonormal columns and $\frac{1}{l} A'X$ is still diagonal. By Gram-Schmidt, we can fill out the remaining
 1338 columns of Ω to construct an orthonormal matrix. \square

1342 We now present the feature result with simplified assumptions.

1343 **Lemma H.3.** *Let Z, A be orthogonal subspaces. Suppose we have a data distribution $\mathcal{D} = \{(z_{y_i} + a_i, y_i)\}_{i=1}^n \subset$
 1344 $\mathbb{R}^d \times \{1, \dots, C\}$, where $z_i \in Z, a_i \in A$ for all $i \in \{1, \dots, n\}$, and the z_i are linearly independent.*

1345 *Let M, M^+ be the matrices defined for this dataset, and let Z, Z^+ and A, A^+ be the corresponding matrices when the
 1346 data is $\{(z_{y_i}, y_i)\}$ and $\{(a_i, y_i)\}$, respectively. Suppose that $(A - A^+)v \neq \mathbf{0}$ for all $v \in \mathbb{R}^d$ s.t. $Av \neq 0$ and the output
 1347 dimension $p \geq C$. Then $W^\top W = Z^\dagger$ is the minimum norm solution to the contrastive learning objective on \mathcal{D} .*

1350 *Proof.* Assume that $d \geq 3C - 2$, otherwise embed the distribution in a space of sufficiently large dimension. By Lemma E.2,
 1351 the minimum norm minimizer is unaffected by adding extra dimensions. Then Lemma H.2 applies, so linear independence
 1352 of the z_{y_i} is sufficient to be able to construct $v_1, \dots, v_C, y_1, \dots, y_C$ satisfying Lemma H.1, from which the conclusion
 1353 follows. \square

1358 I. Minimizer of The Joint Loss

1359 For simplicity we assume $\mu = 0$. Same strategy can be applied to prove the theorem when $\mu \neq 0$ but a more detailed
 1360 discussion on the selection of β may be required.

1362 By Lemmas F.7 and F.1 and the expression of S (equation 13), we observe that the two eigenvectors of M^+ match two of
 1363 the eigenvectors of M . By combining this with Lemma F.6, we obtain that $\beta M^\dagger M^+ + (1 - \beta)M^\dagger \tilde{M} = l_1^+ l_1^{+\top} + l_2^+ l_2^{+\top}$
 1364 on $\text{span}(\{l_1^+, l_2^+\})$ and $\beta M^\dagger M^+ + (1 - \beta)M^\dagger \tilde{M} = (1 - \beta)\tilde{M}^\dagger \tilde{M}$ on $\text{span}(\{l_1^+, l_2^+\})^\perp$. Thus the eigenvalues of
 1365 $\beta M^\dagger M^+ + (1 - \beta)M^\dagger \tilde{M}$ are $1, 1, 1 - \beta, 1 - \beta, \dots, 1 - \beta$. When $\beta \in (0, 1)$, l_1^+ and l_2^+ are the two eigenvectors
 1366 of $\beta M^\dagger M^+ + (1 - \beta)M^\dagger \tilde{M}$ with largest eigenvalues. For the remaining eigenvectors, since they have equally large
 1367 eigenvalues (same as analyzed in H), the minimum norm minimizer will select the largest $p - 2$ of them. In the setting of
 1368 Theorem C.12 $(1 - \beta)(\phi_2^2 + \frac{\sigma_\xi^2}{mn})$ is one of the $p - 2$ largest of the remaining. As a result, both components aligned with v_1
 1369 and v_2 are selected by the minimum norm minimizer of the joint loss.

1372 J. Early in Training Subclasses Are Learned

1373 We assume $\sigma_\xi = O(1)$.
 1374

1375 **J.1. Lemmas**

 1376 **Lemma J.1** (Laurent-Massart (Laurent & Massart, 2000) Lemma 1, page 1325). *Let v_1, \dots, v_d be i.i.d. Gaussian variables*
 1377 *drawn from $\mathcal{N}(0, 1)$. Let $\mathbf{a} = (a_1, \dots, a_d)$ be a vector with non-negative components. Let $Z = \sum_{i=1}^d a_i(v_i^2 - 1)$. The*
 1378 *following inequalities hold for any positive t :*
 1379

1380
$$\Pr(Z \geq 2\|\mathbf{a}\|_2\sqrt{t} + 2\|\mathbf{a}\|_\infty t) \leq e^{-t},$$
 1381
$$\Pr(Z \leq -2\|\mathbf{a}\|_2\sqrt{t}) \leq e^{-t}. \quad (35)$$

 1382 **Lemma J.2** (Mills' ratio. Exercise 6.1 in (Shorack & Shorack, 2000)). *Let v be a Gaussian random variable drawn from*
 1383 *$\mathcal{N}(0, 1)$. Then for all $\lambda > 0$,*

1384
$$\frac{\lambda}{\lambda^2 + 1} \frac{1}{\sqrt{2\pi}} e^{-\frac{\lambda^2}{2}} < \Pr(v \geq \lambda) < \frac{1}{\lambda} \frac{1}{\sqrt{2\pi}} e^{-\frac{\lambda^2}{2}}.$$

 1385 **Corollary J.3.** *Given a vector \mathbf{q} , and a random vector \mathbf{z} drawn from $\mathcal{N}(0, \frac{\sigma}{d}\mathbf{I}_d)$, w.p. $\geq 1 - O(\frac{\delta}{\sqrt{\log 1/\delta}})$, $|\mathbf{z}^\top \mathbf{q}| =$*
 1386 *$O(\frac{\|\mathbf{q}\| \sigma \sqrt{\log \frac{1}{\delta}}}{\sqrt{d}})$.*

 1387 *Proof.* This can be proven by considering the fact that $\mathbf{q}^\top \mathbf{z}$ is a Gaussian variable and applying Lemma J.2. □

 1388 **Lemma J.4.** *Let each element of $\mathbf{W}_0 \in \mathbb{R}^{p \times d}$ be randomly drawn from $\mathcal{N}(0, \frac{\sigma_0^2}{d}\mathbf{I}_d)$. Let $\mathbf{u} \in \mathbb{R}^d$ be a unit vector. With*
 1389 *probability at least $1 - \delta$, we have*

1390
$$\|\mathbf{W}_0 \mathbf{u}\| \geq \sigma_0 \sqrt{\frac{p}{d}} \sqrt{1 - 2\sqrt{\frac{\ln 2/\delta}{p}}}$$
 1391
$$\|\mathbf{W}_0 \mathbf{u}\| \leq \sigma_0 \sqrt{\frac{p}{d}} \sqrt{1 + 2\sqrt{\frac{\ln 2/\delta}{p}} + 2\frac{\ln 2/\delta}{p}}.$$

 1392 *Proof.* Firstly rewrite $\|\mathbf{W}_0 \mathbf{u}\|$ as

1393
$$\|\mathbf{W}_0 \mathbf{u}\| = \sqrt{\sum_{i=1}^p (\mathbf{w}_0^{(i)\top} \mathbf{u})^2} = \sigma_0 \sqrt{\frac{p}{d}} \sqrt{\frac{1}{p} \sum_{i=1}^p (\frac{\sqrt{d}}{\sigma_0} \mathbf{w}_0^{(i)\top} \mathbf{u})^2}.$$

 1394 By spherical symmetric, each $\frac{\sqrt{d}}{\sigma_0} \mathbf{w}_0^{(i)\top} \mathbf{u}$ is a random Gaussian variable drawn from $\mathcal{N}(0, 1)$. By lemma J.1 we have

1395
$$\Pr\left(\frac{1}{p} \sum_{i=1}^p \left(\frac{\sqrt{d}}{\sigma_0} \mathbf{w}_0^{(i)\top} \mathbf{u}\right)^2 \leq 1 - 2\sqrt{\frac{\ln 2/\delta}{p}}\right) \leq \delta/2$$
 1396
$$\Pr\left(\frac{1}{p} \sum_{i=1}^p \left(\frac{\sqrt{d}}{\sigma_0} \mathbf{w}_0^{(i)\top} \mathbf{u}\right)^2 \geq 1 + 2\sqrt{\frac{\ln 2/\delta}{p}} + 2\frac{\ln 2/\delta}{p}\right) \leq \delta/2$$

 1397 which completes the proof. □

 1423 **J.2. Proof of Theorem C.5**

 1424 We assume the dataset satisfies the condition in Section F (which holds with probability $1 - O(\frac{m^2 n^2}{d})$). Let $\mathbf{L}\mathbf{A}\mathbf{L}^\top$ (where
 1425 $\mathbf{C} \in \mathbb{R}^{d \times mn}$) be the eigendecomposition of \mathbf{M} . By equation 13 and Lemma F.7 and Lemma F.1, we observe that when
 1426 $\mu \neq 0$ all but three of \mathbf{M} 's eigenvectors are orthogonal to $\mathbf{l}_1^+, \mathbf{l}_2^+$. W.L.O.G., let $\mathbf{l}_1, \mathbf{l}_2$ and \mathbf{l}_3 be those three eigenvectors.
 1427 The corresponding three eigenvalues are all constants. Let \mathbf{l}_3^+ be a unit vector in $\text{span}(\{\mathbf{l}_1, \mathbf{l}_2, \mathbf{l}_3\}) - \text{span}(\{\mathbf{l}_1^+, \mathbf{l}_2^+\})$.
 1428
 1429

1430
 1431 Decompose \mathbf{v}_2 as $\frac{\mu}{\sqrt{1+\mu^2+\frac{\sigma_\xi^2}{mn}}}\mathbf{l}_1^+ + \frac{\sqrt{1+\frac{\sigma_\xi^2}{mn}}}{\sqrt{1+\mu^2+\frac{\sigma_\xi^2}{mn}}}\mathbf{l}_\perp$ where \mathbf{l}_\perp is a unit vector that is orthogonal to \mathbf{l}_1^+ . Since $\mathbf{v}_2 \perp \mathbf{l}_2^+$, we
 1432
 1433 have $\mathbf{l}_\perp \perp \mathbf{l}_2^+$ thus $\mathbf{M}^+\mathbf{l}_\perp = 0$.

1434 Define

$$\begin{aligned}
 1435 & \sqrt{\mathbf{M}} := \mathbf{L}\sqrt{\mathbf{A}} \\
 1436 & \Gamma_i(t) := \|\mathbf{W}_t\mathbf{l}_i^+\|, \quad i = 1, 2, 3 \\
 1437 & \Gamma_\perp(t) := \|\mathbf{W}_t\mathbf{l}_\perp\| \\
 1438 & \Gamma_{:3}(t) := \sqrt{\sum_{i=1}^3 \|\mathbf{W}_t\mathbf{l}_i^+\|^2} \\
 1439 & \mathbf{B} := [\sqrt{a_4}\mathbf{l}_4 \quad \sqrt{a_5}\mathbf{l}_5 \quad \dots \quad \sqrt{a_{mn}}\mathbf{l}_{mn}] \\
 1440 & \Gamma_B(t) := \|\mathbf{W}_t\mathbf{B}\|_F \\
 1441 & s := \|\sqrt{\mathbf{M}}\| = O(1) \\
 1442 & h := \|\sqrt{\mathbf{M}}^\top \mathbf{B}\| = \sqrt{\sum_{i=4}^{mn} a_i^2} = \sqrt{\sum_{i=3}^K \left(\frac{\sigma_\xi^2}{mn} + \frac{\phi_i^2}{(K-2)}\right)^2 + (mn-K)\frac{\sigma_\xi^4}{m^2n^2}} \\
 1443 & = O\left(\sqrt{\frac{\sigma_\xi^2}{mn} + \frac{1}{K} + \frac{\sigma_\xi^4}{mn}}\right) = O(1)
 \end{aligned}$$

1444 Then we bound $\|\mathbf{W}_t\sqrt{\mathbf{M}}\|_F$

$$\begin{aligned}
 1445 & \|\mathbf{W}_t\sqrt{\mathbf{M}}\|_F = \|\mathbf{W}_t\mathbf{L}\sqrt{\mathbf{A}}\|_F \\
 1446 & = \|\mathbf{W}_t\sqrt{a_1}\mathbf{l}_1 \quad \mathbf{W}_t\sqrt{a_2}\mathbf{l}_2 \quad \dots \quad \mathbf{W}_t\sqrt{a_{mn}}\mathbf{l}_{mn}\|_F \\
 1447 & = \sqrt{\sum_{i=1}^3 \|\mathbf{W}_t\sqrt{a_i}\mathbf{l}_i\|^2 + \sum_{i=4}^{mn} \|\mathbf{W}_t\sqrt{a_i}\mathbf{l}_i\|^2} \\
 1448 & \leq \sqrt{c\Gamma_{:3}(t)^2 + \Gamma_B(t)^2},
 \end{aligned}$$

1449 where c is a constant because a_1, a_2, a_3 are all $O(1)$ (by Lemma F.1) and each \mathbf{l}_i ($i = 1, 2, 3$) is a linear combination of
 1450 $\mathbf{l}_1^{++}, \mathbf{l}_2^{++}, \mathbf{l}_3^+$ with $O(1)$ coefficients, with $\mathbf{l}_1^{++}, \mathbf{l}_2^{++}$ representing the projections of $\mathbf{l}_1^+, \mathbf{l}_2^+$ onto $\text{span}\{\mathbf{l}_i\}_{i=1}^3$.

1451 By the rule of gradient descent we have

$$\begin{aligned}
 1452 & \mathbf{W}_{t+1} = \mathbf{W}_t + \eta(4\mathbf{W}_t\mathbf{M}^+ - 4\mathbf{W}_t\mathbf{M}\mathbf{W}_t^\top\mathbf{W}_t\mathbf{M}) \\
 1453 & = \mathbf{W}_t + 4\eta\mathbf{W}_t\mathbf{M}^+ - 4\eta\mathbf{W}_t\mathbf{M}\mathbf{W}_t^\top\mathbf{W}_t\mathbf{M}
 \end{aligned} \tag{36}$$

1454 This is followed by Lemma J.5.

1455 **Lemma J.5.** *By the update rule of GD we have the following recurrence relations*

$$\begin{aligned}
 1456 & \Gamma_1(t+1) \geq (1 + 4\eta a_1^+) \Gamma_1(t) - 4\eta(c\Gamma_{:3}^2(t) + \Gamma_B(t)^2)^{3/2}s \\
 1457 & \Gamma_1(t+1) \leq (1 + 4\eta a_1^+) \Gamma_1(t) + 4\eta(c\Gamma_{:3}^2(t) + \Gamma_B(t)^2)^{3/2}s \\
 1458 & \Gamma_2(t+1) \leq (1 + 4\eta a_2^+) \Gamma_2(t) + 4\eta(c\Gamma_{:3}^2(t) + \Gamma_B(t)^2)^{3/2}s \\
 1459 & \Gamma_3(t+1) \leq \Gamma_3(t) + 4\eta(c\Gamma_{:3}^2(t) + \Gamma_B(t)^2)^{3/2}s \\
 1460 & \Gamma_\perp(t+1) \leq \Gamma_\perp(t) + 4\eta(c\Gamma_{:3}^2(t) + \Gamma_B(t)^2)^{3/2}s \\
 1461 & \Gamma_B(t+1) \leq \Gamma_B(t) + 4\eta(c\Gamma_{:3}^2(t) + \Gamma_B(t)^2)^{3/2}h.
 \end{aligned}$$

1462 Then we prove the following Lemma

1485 **Lemma J.6.** *At initialization the following holds with probability $\geq 1 - O(\frac{1}{\text{poly}(p)})$*

1486

$$1487 \quad \frac{\Gamma_B(0)}{\Gamma_1(0)} = O(1)$$

1488

$$1489 \quad \Gamma_i(0) = \sigma_0 \sqrt{\frac{p}{d}} \left(1 \pm O\left(\sqrt{\frac{\log p}{p}}\right) \right), \quad i = 1, 2, 3$$

1490

1491

1492

1493

1494

1495

1496

Proof. We first bound $\Gamma_B(0)$

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

1512

$$\begin{aligned} \Gamma_B(0) &= \sqrt{\sum_{i=4}^{mn} a_i \|\mathbf{W}_0 \mathbf{l}_i\|^2} \\ &= \sqrt{\sum_{i=4}^{mn} a_i \sum_{j=1}^p \|\mathbf{w}_{0,j}^\top \mathbf{l}_i\|^2} \\ &\leq \sqrt{\frac{\phi_{\max}^2}{K-2} \sum_{i=4}^{K+1} \sum_{j=1}^p \|\mathbf{w}_{0,j}^\top \mathbf{l}_i\|^2 + \frac{\sigma_\xi^2}{mn} \sum_{i=K+2}^{mn} \sum_{j=1}^p \|\mathbf{w}_{0,j}^\top \mathbf{l}_i\|^2} \\ &= O\left(\sqrt{\frac{p\sigma_0^2}{d} + \sigma_\xi^2 \frac{p\sigma_0^2}{d}}\right) \quad \textcircled{1} \\ &= O\left(\sigma_0 \sqrt{\frac{p}{d}}\right). \end{aligned}$$

1513 Inequality $\textcircled{1}$ holds with probability $\geq 1 - O(\frac{1}{\text{poly}(mnp)})$. It is obtained by observing that $\mathbf{w}_{0,j}^\top \mathbf{l}_i$'s are independent Gaussian

1514 variables (by the orthogonality of \mathbf{l}_i 's) and applying Lemma J.1 to the sum of $\|\mathbf{w}_{0,j}^\top \mathbf{l}_i\|^2$'s.

1515

1516

By Lemma J.4 and the above, at initialization the following holds with probability $\geq 1 - O(\frac{1}{\text{poly}(p)} + \frac{1}{\text{poly}(mnp)})$

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

Let ψ, ψ_B be constants. Define

1530

1531

1532

1533

1534

1535

1536

1537

1538

1539

$$\begin{aligned} \pi &:= \frac{\Gamma_B(0)}{\Gamma_1(0)} = O(1) \quad \text{by Lemma J.6} \\ \tau &:= (c(1 + 2(1 + \psi)^2) + (\pi + \psi_B)^2)^{3/2} = \Theta(1). \end{aligned}$$

Let γ be a constant satisfying the following

$$\gamma \leq \min \left\{ \sqrt{\frac{(a_1^+ - a_2^+) \psi}{\tau(s + s\psi)}}, \sqrt{\frac{a_1^+ \psi}{\tau(s + s\psi)}}, \sqrt{\frac{a_1^+ \psi_B}{\tau(h + s\psi_B)}}, \sqrt{\frac{a_1^+ - a_2^+}{\tau s}} \right\}.$$

□

1540 Note that $a_1^+ - a_2^+ > 0$ because $\mu^2 + 1 > \phi_1^2$. Additionally, we define the following shorthand

$$\begin{aligned}
 1541 \quad & \epsilon := 4\eta\tau\gamma^2s, \\
 1542 \quad & \epsilon_B := 4\eta\tau\gamma^2h \\
 1543 \quad & \alpha := 1 + 4\eta a_1^+ - \epsilon \\
 1544 \quad & \hat{\alpha} := 1 + 4\eta a_1^+ + \epsilon \\
 1545 \quad & \kappa_2 := \frac{1 + 4\eta a_2^+}{\alpha} < 1 \text{ because } \mu^2 + 1 > \phi_1^2 \\
 1546 \quad & \kappa_3 := \frac{1}{\alpha} \\
 1547 \quad & \kappa_{\perp} := \frac{1}{\alpha} \\
 1548 \quad & \kappa_B := \frac{1}{\alpha}.
 \end{aligned}$$

1556 Now we are ready to prove the following Lemma.

1557 **Lemma J.7.** *If $\forall t \leq T$, $\Gamma_1(t) \leq \gamma$. For any constants ψ, ψ_B , the following holds $\forall t \leq T+1$ with probability $1 - O(\frac{1}{\text{poly}(p)})$,*

- 1558 • $\Gamma_1(t) \geq \alpha^t \Gamma_1(0)$
- 1559 • $\Gamma_1(t) \leq \hat{\alpha}^t \Gamma_1(0)$.
- 1560 • $\Gamma_i(t) \leq (\kappa_i^t + \psi) \Gamma_1(t)$, $i = 2, 3$.
- 1561 • $\Gamma_{\perp}(t) \leq (\kappa_{\perp}^t + \psi) \Gamma_1(t)$.

1562 *Proof.* Let $S(k)$ be the following statement: $\forall t'$ such that $0 \leq t' \leq k$, the following holds

- 1563 • $\Gamma_1(t') \geq \alpha^{t'} \Gamma_1(0)$,
- 1564 • $\Gamma_1(t) \leq \hat{\alpha}^t \Gamma_1(0)$,
- 1565 • $\Gamma_i(t') \leq (\kappa_i^{t'} + \psi) \Gamma_1(t')$, $i = 2, 3$,
- 1566 • $\Gamma_{\perp}(t') \leq (\kappa_{\perp}^{t'} + \psi) \Gamma_1(t')$,
- 1567 • $\Gamma_B(t') \leq (\kappa_B^{t'} \pi + \psi_B) \Gamma_1(t')$.

1568 By Lemma J.6, $S(0)$ holds with high probability. Next we show that, $\forall t \in [0, T+1]$, if $S(t-1)$ holds then $S(t)$ also holds.

1569 By Lemma J.5, the induction hypothesis and $\kappa_2, \kappa_3, \kappa_{\perp}, \kappa_B < 1$, $\Gamma_1(t-1) \leq \gamma$, we have the following

$$1570 \quad \Gamma_1(t) \geq \alpha \Gamma_1(t-1) \tag{37}$$

$$1571 \quad \Gamma_1(t) \leq \hat{\alpha} \Gamma_1(t-1) \tag{38}$$

$$1572 \quad \Gamma_2(t) \leq ((1 + 4\eta a_2^+) (\kappa_2^t + \psi) + \epsilon) \Gamma_1(t)$$

$$1573 \quad \Gamma_3(t) \leq ((\kappa_3^t + \psi) + \epsilon) \Gamma_1(t)$$

$$1574 \quad \Gamma_{\perp}(t) \leq ((\kappa_{\perp}^t + \psi) + \epsilon) \Gamma_1(t)$$

$$1575 \quad \Gamma_B(t) \leq ((\kappa_B^t \pi + \psi_B) + \epsilon_B) \Gamma_1(t).$$

1576 By the construction of our κ 's, α 's, ϵ 's and ψ 's, the last three items in statement $S(t)$ hold. Combining the induction hypothesis with equations 37 and 38 yields the first two items in $S(t)$, which completes the proof. \square

1577 Now we are ready to prove the theorem.

Theorem J.8. If $\sigma_0\sqrt{\frac{p}{d}} = o(1)$ and $\sigma_\xi = o(1)$, with probability at least $1 - O(\frac{m^2n^2}{d} + \frac{1}{\text{poly}(p)}) = 1 - o(1)$, the following holds

- $\|\mathbf{W}_0\mathbf{v}_2\| = o(1)$.
- $\exists t = O(\ln(\frac{1}{\sigma_0}\sqrt{\frac{d}{p}}))$, s.t. $\|\mathbf{W}_t\mathbf{v}_2\| = \Omega(1)$.

Proof. $\|\mathbf{W}_0\mathbf{v}_3\| = o(1)$ follows Lemma J.4 and the assumption that $\sigma_0\sqrt{\frac{p}{d}} = o(1)$. Select a constant ψ such that $\psi < \frac{\mu}{\sqrt{1 + \frac{\sigma_\xi^2}{mn}}}$. Note that $\frac{\mu}{\sqrt{1 + \frac{\sigma_\xi^2}{mn}}} - \psi = \Theta(1)$. Let $T = \lfloor \frac{\ln(\gamma/\Gamma_1(0))}{\ln \alpha} \rfloor = \Theta(\ln \frac{1}{\sigma_0}\sqrt{\frac{d}{p}})$. There are two cases to consider.

- If $\forall t \leq T$, $\Gamma_1(t) \leq \gamma$, by Lemma J.7 we have $\Gamma_1(T+1) \geq \gamma$ and $\Gamma_\perp(T+1) \leq (o(1) + \psi)\Gamma_1(T+1)$. Then

$$\begin{aligned} \|\mathbf{W}_{T+1}\mathbf{v}_2\| &\geq \frac{\mu}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} \|\mathbf{W}_{T+1}\mathbf{l}_1^+\| - \frac{\sqrt{1 + \frac{\sigma_\xi^2}{mn}}}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} \|\mathbf{W}_{T+1}\mathbf{l}_\perp\| \\ &\geq \left(\frac{\mu}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} - \frac{\sqrt{1 + \frac{\sigma_\xi^2}{mn}}}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} \psi - o(1) \right) \gamma \\ &= \Omega(1). \end{aligned}$$

- If $\exists t \leq T$ s.t. $\Gamma_1(t) > \gamma$, we define $T^* = \frac{\ln(\frac{\gamma}{\Gamma_1(0)})}{\ln \alpha}$ and $t^* = \min t$ s.t. $\Gamma_1(t) > \gamma$. It follows that $\forall t \leq t^* - 1$, $\Gamma_1(t) \leq \gamma$. Then we can apply Lemma J.7 to obtain $\Gamma_1(t^*) \leq \hat{\alpha}^{t^*} \Gamma_1(0)$. If $t < T^*$, the above yields $\Gamma_1(t^*) < \gamma$, which contradicts the definition of t^* . Therefore we conclude $t^* \geq T^*$. Lemma J.7 also tells that $\Gamma_\perp(t^*) \leq (\kappa_\perp^{t^*} + \psi)\Gamma_1(t^*)$. Since $t^* \geq T^*$ and $\kappa_\perp < 1$, we have $\kappa_\perp^{t^*} \leq (\frac{\Gamma_1(0)}{\gamma})^{\frac{\ln(1/\kappa_\perp)}{\ln \alpha}} = o(1)$. Therefore $\Gamma_\perp(t^*) \leq (o(1) + \psi)\Gamma_1(t^*)$. By the definition of t^* , $\Gamma_1(t^*) > \gamma$. Then we can lower bound $\|\mathbf{W}_{t^*}\mathbf{v}_2\|$ in the same way as in the previous case

$$\begin{aligned} \|\mathbf{W}_{t^*}\mathbf{v}_2\| &\geq \frac{\mu}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} \|\mathbf{W}_{t^*}\mathbf{l}_1^+\| - \frac{\sqrt{1 + \frac{\sigma_\xi^2}{mn}}}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} \|\mathbf{W}_{t^*}\mathbf{l}_\perp\| \\ &\geq \left(\frac{\mu}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} - \frac{\sqrt{1 + \frac{\sigma_\xi^2}{mn}}}{\sqrt{1 + \mu^2 + \frac{\sigma_\xi^2}{mn}}} \psi - o(1) \right) \gamma \\ &= \Omega(1). \end{aligned}$$

□

Table 6. increasing k improves both subclass and class accuracies on CIFAR-10 RandBit.

| k | Sub Acc | Acc |
|-----|---------|-------|
| 1 | 34.38 | 86.73 |
| 16 | 58.12 | 94.09 |

K. Experimental Setup and Additional Experimental Results

K.1. Datasets

CIFAR-10/100. The two datasets each consist of 60000 32x32 colour images (Krizhevsky et al., 2009). In the case of CIFAR-10, the ‘classes’ refer to the original 10 classes defined in the dataset, while we define ‘subclasses’ as two subclasses: vehicles (airplane, automobile, ship, truck) and animals (bird, cat, deer, dog, frog, horse). On CIFAR-100, we refer to the 10 super-classes (e.g. aquatic mammals, fish, flowers) as our ‘classes’ and the 100 classes as our ‘sub-classes’. These two datasets illustrate a natural setting where class collapse is extremely harmful, as it results in learning representations that do not capture much of the semantically relevant information from the data.

MNIST RandBit. The MNIST RandBit dataset Chen et al. (2021) is created by setting n , the # of bits that specifies how easy the useless feature will be. Larger n makes the feature more discriminative, thus ‘easier’ and more problematic for feature suppression. An extra channel is concatenated to MNIST images where each value in the feature map corresponds to a random integer between 0 and 2^n .

CIFAR-10/100 RandBit. The two datasets are constructed in a similar way as MNIST RandBit, but with images from CIFAR-10/100.

K.2. Training details

For the experiments on CIFAR-10/100 or CIFAR-100 RandBit, we use a ResNet-18 trained with (Momentum) SGD using learning rate = 0.01 and momentum = 0.9. We train with batch size set to 512 for 1000 epochs. For data augmentations, we consider the standard data augmentations from Chen et al. (2020).

For the feature suppression experiments on MNIST RandBit, we directly use the code provided by Chen et al. (2021). We consider a 5-Layer convolutional network. For our data augmentations, we consider the standard set of data augmentations for images and do not alter the useless feature (extra channel concatenated of RandBits).

K.3. Details and additional experiments on varying embedding size

In the experiments presented in Table 2, we vary the width, denoted by w , of the ResNet, which is controlled by the number of convolutional layer filters. For width w , there are $w, 2w, 4w, 8w$ filters in each layer of the four ResNet blocks.

In addition, we explore an alternative way of varying the embedding size, which isolates the effect of the last layer’s embedding size from the size of the lower layers. Specifically, we set the width parameter $w = 4$ and multiply the width of only the last ResNet block by a factor k . It is worth noting that doing this requires a much smaller total number of parameters. Table 6 presents the results on CIFAR-10 RandBit. We observe that increasing k also effectively improves the accuracy. Although the improvement is not as substantial as in the previous case where we increase w , it confirms the same trend predicted by the theory, supporting the conclusion that increasing the embedding size alleviates feature suppression.