

Towards Understanding Ensemble, Knowledge Distillation and Self-Distillation in Deep Learning

Zeyuan Allen-Zhu
zeyuanallen@meta.com
Meta FAIR Labs

Yuanzhi Li
Yuanzhi.Li@mbzuai.ac.ae
Mohamed bin Zayed University of AI

February 15, 2023

(version 3)*

Abstract

We formally study how *ensemble* of deep learning models can improve test accuracy, and how the superior performance of ensemble can be distilled into a single model using *knowledge distillation*. We consider the challenging case where the ensemble is simply an average of the outputs of a few independently trained neural networks with the *same* architecture, trained using the *same* algorithm on the *same* data set, and they only differ by the random seeds used in the initialization.

We show that ensemble/knowledge distillation in *deep learning* works very differently from traditional learning theory (such as boosting or NTKs, neural tangent kernels). To properly understand them, we develop a theory showing that when data has a structure we refer to as “multi-view”, then ensemble of independently trained neural networks can provably improve test accuracy, and such superior test accuracy can also be provably distilled into a single model by training a single model to match the output of the ensemble instead of the true label. Our result sheds light on how ensemble works in deep learning in a way that is completely different from traditional theorems, and how the “dark knowledge” is hidden in the outputs of the ensemble and can be used in distillation. In the end, we prove that self-distillation can also be viewed as implicitly combining ensemble and knowledge distillation to improve test accuracy.

*V1.5 appears on arXiv on this date, and V2/V3 polishes writing. An extended abstract of V3 will appear in ICLR 2023.

1 Introduction

Ensemble [25, 41, 50, 67, 70, 70, 71, 73, 91], also known as model averaging, is one of the oldest and most powerful techniques in practice to improve the performance of deep learning models. By simply averaging the output of merely a few (like 3 or 10) independently trained neural networks of the *same* architecture, using the *same* training method over the *same* training data, it can significantly boost the prediction accuracy over the test set comparing to individual models. The only difference is the randomness used to initialize these neural networks and/or the randomness during training. For example, on the standard CIFAR-100 data set, averaging the output of ten independently trained ResNet-34 can easily offer a 5% improvement in terms of test accuracy. Moreover, it is discovered by Hinton et al. [42] that such superior test-time performance of the ensemble can be transferred into a single model (of the same size as the individual models) using a technique called *knowledge distillation*: that is, simply train a single model to match the output of the ensemble (such as “90% cat + 10% car”, also known as *soft labels*) as opposite to the true data labels, over the same training data.

On the theory side, there are lots of works studying the superior performance of ensemble from principled perspectives [11, 13, 28, 30–33, 36, 43, 46–48, 51, 72, 72, 75, 76]. However, most of these works only apply to: (1). Boosting: where the coefficients associated with the combinations of the single models are actually trained, instead of simply taking average; (2). Bootstrapping/Bagging: the training data are different for each single model; (3). Ensemble of models of different types and architectures; or (4). Ensemble of random features or decision trees.

To the best of our knowledge, *none* of these cited works apply to the particular type of ensemble that is widely used in deep learning: simply take a *uniform* average of the output of the learners, which are neural networks with the *same* architecture and are trained by stochastic gradient descent (SGD) over the *same* training set. In fact, *very critically, for deep learning models*:

- TRAINING AVERAGE DOES NOT WORK: if one directly trains to learn an average of individual neural networks initialized by different seeds, the performance is much worse than ensemble.
- KNOWLEDGE DISTILLATION WORKS: the superior performance of ensemble in deep learning can be distilled into a single model [20, 22, 29, 34, 42, 52, 61].
- SELF-DISTILLATION WORKS: even distilling a single model into another single model of the same size, there is performance boost. [35, 63, 89]

We are unaware of any satisfactory theoretical explanation for the phenomena above. For instance, as we shall argue, some traditional view for why ensemble works, such as ‘ensemble can enlarge the feature space in random feature mappings’, even give contradictory explanations to the above phenomena, thus *cannot* explain knowledge distillation or ensemble in *deep learning*. Motivated by this gap between theory and practice we study the following question for *multi-class* classification:

Our theoretical questions:

How does ensemble improve the test-time performance in deep learning when we simply (unweightedly) average over a few independently trained neural networks? – Especially when all the neural networks have the same architecture, are trained over the same data set using the same standard training algorithm (i.e. gradient descent with the same learning rate and sample regularization) and only differ by the random seeds, and even when all single models already have 100% *training accuracy*? How can such superior test-time performance of ensemble be later “distilled” into a single neural network of the same architecture, simply by training the single model to match the output of the ensemble over the same training data set?

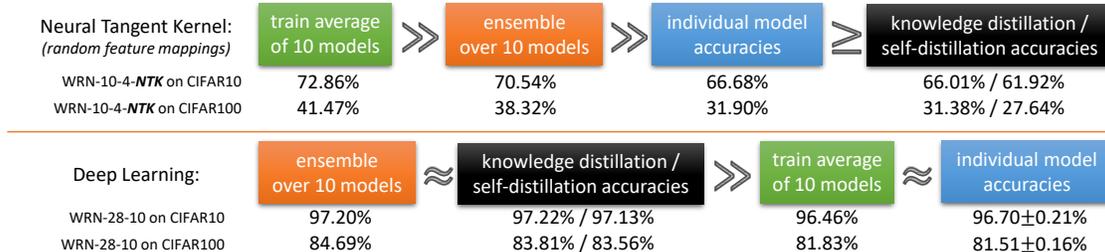


Figure 1: Ensemble in deep learning is very different from ensemble in random feature mappings. Details in Figure 6.

1.1 Our Theoretical Results at a High Level

To the best of our knowledge, this paper makes a first step towards answering these questions in deep learning. On the *theory side*, we prove for certain multi-class classification tasks with a special structure we refer to as **multi-view**, with a training set \mathcal{Z} consisting of N i.i.d. samples from some unknown distribution \mathcal{D} , for certain two-layer convolutional network f with (smoothed-)ReLU activation as learner:

- (Single model has bad test accuracy): there is a value $\mu > 0$ such that when a single model f is trained over \mathcal{Z} using the cross-entropy loss, via gradient descent (GD) starting from random Gaussian initialization, the model can reach zero training error *efficiently*. However, w.h.p. the prediction (classification) error of f over \mathcal{D} is between 0.49μ and 0.51μ .
- (Ensemble provably improves test accuracy): let f_1, f_2, \dots, f_L be $L = \tilde{\Omega}(1)$ independently trained single models as above, then w.h.p. $G = \frac{1}{L} \sum_{\ell} f_{\ell}$ has prediction error $\leq 0.01\mu$ over \mathcal{D} .
- (Ensemble can be distilled into a single model): if we further train (using GD from random initialization) another single model f_0 (same architecture as each f_{ℓ}) to match the output of $G = \frac{1}{L} \sum_{\ell} f_{\ell}$ merely over the same training data set \mathcal{Z} , then f_0 can be trained *efficiently* and w.h.p. f_0 will have prediction error $\leq 0.01\mu$ over \mathcal{D} as well.
- (*Self-distillation* also improves test accuracy): if we further train (using GD from random initialization) another single model f' (same architecture as f_1) to match the output of *the single model* f_1 merely over the same training data set \mathcal{Z} , then f' can be trained *efficiently* and w.h.p. has prediction error at most $\leq 0.26\mu$ over \mathcal{D} . The main idea is that self-distillation is performing “*implicit ensemble + knowledge distillation*”, as we shall argue in Section 4.2.

Thus, on the theory side, we make a step towards understanding ensemble and knowledge distillation in deep learning both computationally (training efficiency) and statistically (generalization error).

1.2 Our Empirical Results at a Glance

We defer discussions of our empirical results to Section 5. However, we highlight some of the empirical findings, as they shall confirm and justify our theoretical approach studying ensemble and knowledge distillation in deep learning. Specifically, we give empirical evidences showing that:

- Knowledge distillation does not work for random feature mappings; and ensemble in deep learning is very different from ensemble in random feature mappings (see Figure 1).
- Special structures in data (such as the “multi-view” structure we shall introduce) is needed for ensemble of neural networks to work.
- The variance due to label noise or the non-convex landscape of training, in the independently-trained models, may not be connected to the superior performance of ensemble in deep learning.

2 Our Methodology and Intuition

2.1 A Failure Attempt Using Random Feature Mappings

The recent advance in deep learning theory shows that under certain circumstances, neural networks can be treated as a linear function over random feature mappings [3, 5, 6, 8, 9, 19, 24, 26, 27, 39, 40, 44, 55, 59, 86, 92]. In particular, the theory shows when $f : \mathbb{R}^{D+d} \rightarrow \mathbb{R}$ is a neural network with inputs $x \in \mathbb{R}^d$ and weights $W \in \mathbb{R}^D$, in some cases, $f(W, x)$ can be approximated by:

$$f(W, x) \approx f(W_0, x) + \langle W - W_0, \nabla_W f(W_0, x) \rangle$$

where W_0 is the random initialization of the neural network, and $\Phi_{W_0}(x) := \nabla_W f(W_0, x)$ is the neural tangent kernel (NTK) feature mapping. This is known as the NTK approach. If this approximation holds, then training a neural network can be approximated by learning a linear function over random features $\Phi_{W_0}(x)$, which is very theory-friendly.

Ensemble works for random features / NTK. Traditional theorems [1, 14, 17, 81] suggest that the ensemble of independently trained random feature models can indeed *significantly improve* test-time performance, as it enlarges the feature space from $\Phi_{W_0}(x)$ to $\{\Phi_{W_0^{(i)}}(x)\}_{i \in [L]}$ for L many independently sampled $W_0^{(i)}$. This can be viewed as a feature selection process [7, 18, 66, 68, 74], and we have confirmed it for NTK in practice, see Figure 1. Motivate by this line of research, we ask:

Can we understand ensemble and knowledge distillation in deep learning as feature selections? (in particular, using the NTK approach?)

Unfortunately, our empirical results provide many counter examples towards those arguments, see discussions below and Figure 1.

Contradiction 1: training average works even better. Although ensemble of linear functions over NTK features with different random seeds: $f_i(x) = \langle W^{(i)}, \Phi_{W_0^{(i)}}(x) \rangle$ does improve test accuracy, however, such improvement is mainly due to the use of a larger set of random features, whose combinations contain functions that generalize better. To see this, we observe that an even superior performance (than the ensemble) can simply be obtained by directly training $F(x) = \frac{1}{L}(f_1 + f_2 + \dots + f_L)$ from random initialization. *In contrast*, recall if $f_i(x)$'s are multi-layer neural networks with different random seeds, then training their average barely gives any better performance comparing to individual networks f_i , as now all the f_i 's are capable of learning the same set of features.

Contradiction 2: knowledge distillation does not work. For NTK feature mappings, we observe that the result obtained by ensemble cannot be distilled at all into individual models, indicating the features selected by ensemble *is not contained* in the feature $\Phi_{W_0^{(i)}}(x)$ of any individual model. In contrast, in actual deep learning, ensemble *does not enlarge feature space*: so an individual neural network is capable of learning the features of the ensemble model.

In sum, ensemble in deep learning may be very different from ensemble in random features. It may be more accurate to study ensemble / knowledge distillation in deep learning as a *feature learning process*, instead of a feature selection process (where the features are prescribed and only their linear combinations are trained). But still, we point out a fundamental difficulty:

Key challenge:

If a single deep learning model is capable of— through knowledge distillation— learning the features of the ensemble model and achieving better test accuracy comparing to training the single model directly (and the same training accuracy, typically at global optimal of 100%), then why the single model *cannot learn* these features directly when we train the model to match the true data labels? What is the **dark knowledge** hidden in the output of ensemble (a.k.a. soft label)¹ comparing to the original hard label?

2.2 Ensemble in Deep Learning: a Feature Learning Process

Before addressing the key challenge, we point out that prior works are very limited with respect to studying neural network training as a feature learning process, due to the extreme non-convexity obstacle in optimization.

Most of the existing works proving that neural networks can learn features only focus on the case *when the input is Gaussian or Gaussian-like* [10, 12, 16, 37, 38, 45, 53, 54, 56, 56–58, 60, 69, 78–80, 84, 85, 87, 90]. However, as we demonstrate in Figure 7 on Page 15,

Ensemble in DL might not improve test accuracy when inputs are Gaussian-like:

Empirically, ensemble *does not* improve test accuracy in deep learning, in certain scenarios when the distribution of the input data is Gaussian or even mixture of Gaussians. This is true over various learner network structures (fully-connected, residual, convolution neural networks) and various labeling functions (when the labels are generated by linear functions, fully-connected, residual, convolutional networks, with/without label noise, with/without classification margin).

Bias variance view of ensemble: Some prior works also try to attribute the benefit of ensemble as reducing the *variance* of individual solutions [15, 62, 64, 82, 83] due to label noise or non-convex landscape of the training objective (so some individual models might simply not be trained very well by over-fitting to the label noise or stuck at a bad local minimal).

However, reducing such variance can reduce a convex test loss (typically cross-entropy), but not necessarily the *test classification error*. Concretely, the synthetic experiments in Figure 7 show that, after applying ensemble over Gaussian-like inputs, the variance of the model outputs is reduced but the test accuracy is *not improved*. We give many more empirical evidences to show that the variance (either from label noise or from the non-convex landscape) is usually not the cause for why ensemble works in deep learning, see Section 5. Moreover, we point out that (see Figure 6) in practice, typically the individual neural networks are trained *equally well*, meaning that they all have perfect training accuracy and almost identical test error, yet ensemble these models still improves the test accuracy significantly.

Hence, to understand the true benefit of ensemble in deep learning in theory, we would like to study a setting that can *approximate* practical deep learning, where:

- The input distribution is more structured than standard Gaussian and there is no label noise. (From above discussions, ensemble cannot work for deep learning distribution-freely, nor even under Gaussian distribution).
- The individual neural networks all are well-trained, in the sense that the training accuracy in the end is 100%, and there is nearly no variance in the test accuracy for individual models. (So training never fails.)

¹For a k -class classification problem, the output of a model $g(x)$ is usually k -dimensional, and represents a soft-max probability distribution over the k target classes. This is known as the *soft label*.



Figure 2: Illustration of images with multiple views (features) in the ImageNet dataset.

We would like to re-elaborate the *key challenge*: ‘ensemble improves test accuracy’ implies that different single models need to learn different sets of features; however, all these models have the same architecture, and trained using the same learning algorithm (SGD with momentum) with identical learning rates, and each the (learned) sets of features in each model lead to the perfect 100% training accuracy and an almost identical test accuracy. Thus, the difference of the features *must not be* due to ‘difference in the data set’, ‘difference in the models’, ‘difference in the training algorithms’, ‘difference in the learning rates’, ‘failure in training occasionally’, ‘failure in generalization in some cases’, etc. Additional principles need to be developed to incorporate the effect of ensemble in deep learning.

In this work, we propose to study a setting of data that we refer to as **multi-view**, where the above two conditions both hold when we train a two-layer neural networks with (smoothed-)ReLU activations. We also argue that the multi-view structure we consider is fairly common in the data sets used in practice, in particular for vision tasks. We give more details below.

2.3 Our Approach: Learning Multi-View Data

Let us first give a thought experiment to illustrate our approach, and we present the precise mathematical definition of the “multi-view” structure in Section 3. Consider a *binary* classification problem and four “features” v_1, v_2, v_3, v_4 . The first two features correspond to the first class label, and the next two features correspond to the second class label. In the data distribution:

- When the label is class 1, then:²

$$\left\{ \begin{array}{ll} \text{both } v_1, v_2 \text{ appears with weight 1, one of } v_3, v_4 \text{ appears with weight 0.1} & \text{w.p. 80\%;} \\ \text{only } v_1 \text{ appears with weight 1, one of } v_3, v_4 \text{ appears with weight 0.1} & \text{w.p. 10\%;} \\ \text{only } v_2 \text{ appears with weight 1, one of } v_3, v_4 \text{ appears with weight 0.1} & \text{w.p. 10\%.} \end{array} \right.$$

- When the label is class 2, then

$$\left\{ \begin{array}{ll} \text{both } v_3, v_4 \text{ appears with weight 1, one of } v_1, v_2 \text{ appears with weight 0.1} & \text{w.p. 80\%;} \\ \text{only } v_3 \text{ appears with weight 1, one of } v_1, v_2 \text{ appears with weight 0.1} & \text{w.p. 10\%;} \\ \text{only } v_4 \text{ appears with weight 1, one of } v_1, v_2 \text{ appears with weight 0.1} & \text{w.p. 10\%.} \end{array} \right.$$

We call the 80% of the data *multi-view data*: these are the data where multiple features exist and can be used to classify them correctly. We call the rest 20% of the data *single-view data*: some features for the correct labels are missing.

Meaningfulness of our multi-view hypothesis. Such “multi-view” structure is very common in many of the datasets where deep learning excels. In vision datasets in particular, as illustrated in Figure 2, a car image can be classified as a car by looking at the headlights, the wheels, or the windows. For a typical placement of a car in images, we can observe all these features, and it suffices to use one of the features to classify it as a car. However, there are some car images taken

²One can for simplicity think of “ v appears with weight α and w appears with weight β ” as $\text{data} = \alpha v + \beta w + \text{noise}$.

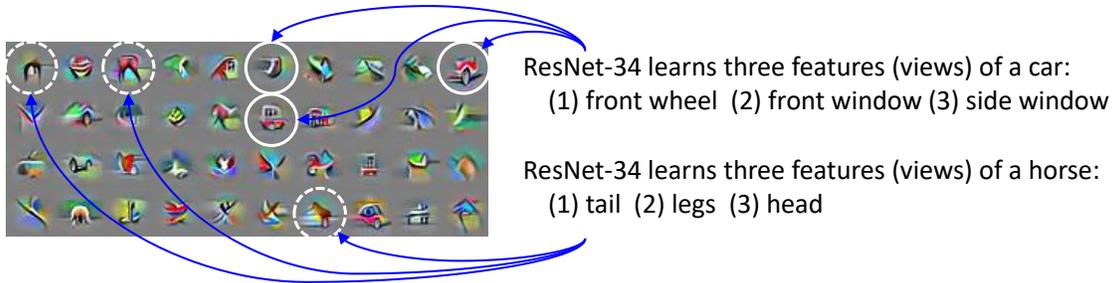


Figure 3: Visualization of the channels in layer-23 of a ResNet-34 trained on CIFAR-10, picture from [4].

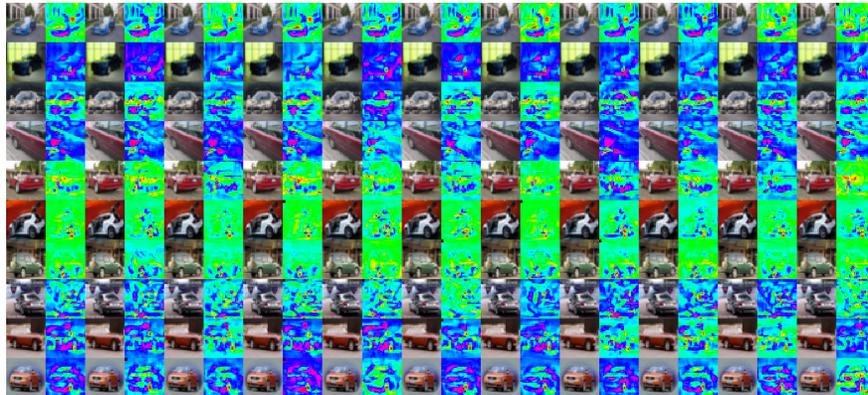


Figure 4: Ten independently trained ResNet-34 models (and their ensemble) detect car images through different reasonings, suggesting that the data has multi views, and independently trained neural networks do utilize this structure. The numerical experiments in Figure 9 also suggest the existence of multi views.

from a particular angle, where one or more of these features are missing. For example, an image of a car facing forward might be missing the wheel feature. Moreover, some car might also have a small fraction of “cat features”: for example, the headlight might appear similar to cat eyes the ear of a cat. This can be used as the “dark knowledge” by the single model to learn from the ensemble.

In Figure 3, we visualize the learned features from an actual neural network to show that they can indeed capture different views. In Figure 4, we plot the “heatmap” for some car images to illustrate that single models (trained from different random seeds) indeed pick up different parts of the input image to classify it as a car. In Figure 9, we manually delete for instance 7/8 of the channels in some intermediate layer of a ResNet, and show that the test accuracy may not be affected by much after ensemble— thus supporting that the multi-view hypothesis can indeed exist even in the intermediate layers of a neural network and ensemble is indeed collecting all these views.

How individual neural networks learn. Under the multi-view data defined above, if we train a neural network using the cross-entropy loss via gradient descent (GD) from random initialization, during the training process of the individual networks, we show that:

- The network will quickly pick up one of the feature $v \in \{v_1, v_2\}$ for the first label, and one of the features $v' \in \{v_3, v_4\}$ for the second label. So, 90% of the training examples, consisting of all the multi-view data and half of the single-view data (those with feature v or v'), are classified correctly. Once classified correctly (with a large margin), these data begin to contribute negligible to gradient by the nature of the cross-entropy loss.
- Next, the network will memorize (using e.g. the noise in the data) the remaining 10% of the

training examples without learning any new features, due to insufficient amount of left-over samples after the first phase, thus achieving training accuracy 100% but test accuracy 90%.

How ensemble improves test accuracy. It is simple why ensemble works. Depending on the randomness of initialization, each individual network will pick up v_1 or v_2 each w.p. 50%. Hence, as long as we ensemble $\tilde{O}(1)$ many independently trained models, w.h.p. their ensemble will pick up both features $\{v_1, v_2\}$ and both features $\{v_3, v_4\}$. Thus, all the data will be classified correctly.

How knowledge distillation works. Perhaps less obvious is how knowledge distillation works. Since ensemble learns all the features v_1, v_2, v_3, v_4 , given a multi-view data with label 1, the ensemble will actually output $\propto (2, 0.1)$, where the 2 comes from features v_1, v_2 and 0.1 comes from one of v_3, v_4 . On the other hand, an individual model learning only one of v_3, v_4 will actually output $\propto (2, 0)$ when the feature v_3 or v_4 in the data does not match the one learned by the model. Hence, by training the individual model to match the output of the ensemble, the individual model is *forced* to learn both features v_3, v_4 , even though it has already perfectly classified the training data.

This is the “dark knowledge” hidden in the output of the ensemble model.

(This theoretical finding is consistent with practice: Figure 8 suggests that models trained from knowledge distillation should have learned most of the features, and further computing their ensemble does not give much performance boost.)

2.4 Significance of Our Technique

Our work belongs to the generic framework where one can prove that certain aspects of the learning algorithm (in this paper, the randomness of the initialization) affects the order where the features are learned, which we believe is also one of the key ingredients to understand the role of the learning algorithm in terms of generalization in deep learning. This is **fundamentally different from convex optimization**, such as kernel method, where (with an ℓ_2 regularization) there is an unique *global minimum* so the choice of optimization algorithm or the random seed of the initialization does not matter (thus, ensemble does not help at all). There are other works that consider other aspects, such as the choice of learning rate [59], that can affect the order where the features are picked in deep learning. In that work [59], the two “features” are asymmetric: a memorizable feature and a generalizable feature, so the learning rate will decide which feature to be picked. In our work, the features are “symmetric”, so the randomness of the initialization will decide which feature to be picked. Our technique is fundamentally different from [59]: they only focus on the NTK setting, where only a linear function over the prescribed sequence of feature mappings is learned. In other words, their features are *not learned* (although their features change over time, following a Gaussian random process which is independent of the learning task); instead, we study a *feature learning* process in this paper. As we have argued and shown empirically, the NTK setting cannot be used to explain ensemble and distillation in deep learning.

We believe that our work extends the reach of traditional optimization and statistical machine learning theory, where typically the statistics (generalization) is separated from optimization (training). As we have pointed out, such “separate” treatment might not be possible to understand (at least ensemble or knowledge distillation in) deep learning.

3 Problem Setup

In this paper, we consider the following data distribution with “multi-view”, that allows us to formally prove our result on ensemble and knowledge distillation for two-layer neural networks. The data distribution is a straight-forward generalization of the intuitive setting in Section 2.3.

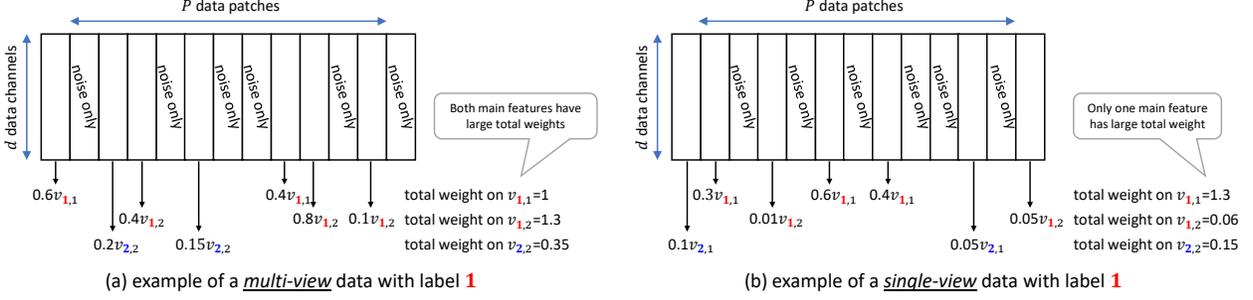


Figure 5: Illustration of a multi-view and a single-view data point; the feature vectors can also be combined with feature noise and random noise, see Def. 3.1.

For simplicity, in the main body, we use example choices of the parameters mainly a function of k (such as $P = k^2$, $\gamma = \frac{1}{k^{1.5}}$, $\mu = \frac{k^{1.2}}{N}$, $\rho = k^{-0.01}$, $\sigma_0 = 1/\sqrt{k}$ as we shall see), and we consider the case when k is sufficiently large. In our formal statements of the theorems in the appendix, we shall give a much larger range of parameters for the theorems to hold.

3.1 Data Distribution and Notations

We consider learning a k -class classification problem over P -patch inputs, where each patch has dimension d . In symbols, each labelled data is represented by (X, y) where $X = (x_1, x_2, \dots, x_P) \in (\mathbb{R}^d)^P$ is the data vector and $y \in [k]$ is the data label. For simplicity, we focus on the case when $P = k^2$, and $d = \text{poly}(k)$ for a large polynomial.

We consider the setting when k is sufficiently large.³ We use “w.h.p.” to denote with probability at least $1 - e^{-\Omega(\log^2 k)}$, and use $\tilde{O}, \tilde{\Theta}, \tilde{\Omega}$ notions to hide polylogarithmic factors in k .

We first assume that each label class $j \in [k]$ has multiple associated features, say *two features for the simplicity of math*, represented by unit **feature vectors** $v_{j,1}, v_{j,2} \in \mathbb{R}^d$. For notation simplicity, we assume that all the features are orthogonal, namely,

$$\forall j, j' \in [k], \forall \ell, \ell' \in [2], \|v_{j,\ell}\|_2 = 1 \quad \text{and} \quad v_{j,\ell} \perp v_{j',\ell'} \quad \text{when} \quad (j, \ell) \neq (j', \ell')$$

although our work also extends to the “incoherent” case trivially. We denote by

$$\mathcal{V} \stackrel{\text{def}}{=} \{v_{j,1}, v_{j,2}\}_{j \in [k]} \quad \text{the set of all features.}$$

We consider the following data and label distribution. Let C_p be a global constant, $s \in [1, k^{0.2}]$ be a sparsity parameter. To be concise, we define the *multi-view distribution* \mathcal{D}_m and *single-view distribution* \mathcal{D}_s together. Due to space limitation, here we hide the specification of the random “noise”, and defer the full definition to Appendix A.⁴

Definition 3.1 (data distributions \mathcal{D}_m and \mathcal{D}_s). *Given $\mathcal{D} \in \{\mathcal{D}_m, \mathcal{D}_s\}$, we define $(X, y) \sim \mathcal{D}$ as follows. First choose the label $y \in [k]$ uniformly at random. Then, the data vector X is generated as follows (also illustrated in Figure 5).*

1. Denote $\mathcal{V}(X) = \{v_{y,1}, v_{y,2}\} \cup \mathcal{V}'$ as the set of feature vectors used in this data vector X , where \mathcal{V}' is a set of features uniformly sampled from $\{v_{j',1}, v_{j',2}\}_{j' \in [k] \setminus \{y\}}$, each with probability $\frac{s}{k}$.

³If we want to work with fixed k , say $k = 2$, our theorem can also be modified to that setting by increasing the number of features per class. In this case, a subset of features per class will be learned by each individual neural network. We keep our current setting with two features to simplify the notations.

⁴At a high level, we shall allow such “noise” to be any feature noise plus Gaussian noise, such as $\text{noise} = \sum_{v' \in \mathcal{V}} \alpha_{p,v'} v' + \xi_p \in \mathbb{R}^d$, where each $\alpha_{p,v'} \in [0, \gamma]$ can be arbitrary, and $\xi_p \sim \mathcal{N}(0, \sigma_p^2 \mathbf{I})$.

2. For each $v \in \mathcal{V}(X)$, pick C_p many disjoint patches in $[P]$ and denote it as $\mathcal{P}_v(X) \subset [P]$ (the distribution of these patches can be arbitrary). We denote $\mathcal{P}(X) = \cup_{v \in \mathcal{V}(X)} \mathcal{P}_v(X)$.
3. If $\mathcal{D} = \mathcal{D}_s$ is the single-view distribution, pick a value $\hat{\ell} = \hat{\ell}(X) \in [2]$ uniformly at random.
4. For each $v \in \mathcal{V}(X)$ and $p \in \mathcal{P}_v(X)$, we set $x_p = z_p v + \text{“noise”} \in \mathbb{R}^d$, where, the random coefficients $z_p \geq 0$ satisfy that:

In the case of multi-view distribution $\mathcal{D} = \mathcal{D}_m$,

- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [1, O(1)]$ when $v \in \{v_{y,1}, v_{y,2}\}$,⁵
- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [\Omega(1), 0.4]$ when $v \in \mathcal{V}(X) \setminus \{v_{y,1}, v_{y,2}\}$,⁶

In the case of single-view distribution $\mathcal{D} = \mathcal{D}_s$,

- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [1, O(1)]$ when $v = v_{y,\hat{\ell}}$,
- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [\rho, O(\rho)]$ when $v = v_{y,3-\hat{\ell}}$,
- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [\Omega(\Gamma), \Gamma]$ when $v \in \mathcal{V}(X) \setminus \{v_{y,1}, v_{y,2}\}$.

5. For each $p \in [P] \setminus \mathcal{P}(X)$, we set x_p to consist only of “noise”.

Remark 3.2. The distribution of how to pick $\mathcal{P}(X)$ and assign $\sum_{p \in \mathcal{P}_v(X)} z_p$ to each patch in $p \in \mathcal{P}_v(X)$ can be arbitrary (and can depend on other randomness in the data as well). In particular, we have allowed **different features** $v_{j,1}, v_{j,2}$ to show up with **different weights** in the data (for example, for multi-view data, some view $v_{y,1}$ can consistently have larger z_p comparing to $v_{y,2}$). Yet, we shall prove that the order to learn these features by the learner network *can still be flipped* depending on the randomness of network initialization.

Interpretation of our data distribution. As we argue more in Appendix A, our setting can be tied to a down-sized version of convolutional networks applied to image classification data. With a small kernel size, good features in an image typically appear only at a few patches, and most other patches are random noise or low-magnitude feature noises. More importantly, our noise parameters shall ensure that, the concept class is *not learnable by linear classifiers or constant degree polynomials*. We believe a (convolutional) neural network with ReLU-like activation is somewhat necessary.

Our final data distribution \mathcal{D} , and the training data set \mathcal{Z} are formally given as follows.

Definition 3.3 (\mathcal{D} and \mathcal{Z}). *The distribution \mathcal{D} consists of data from \mathcal{D}_m w.p. $1-\mu$ and from \mathcal{D}_s w.p. μ . We are given N training samples from \mathcal{D} , and denote the training data set as $\mathcal{Z} = \mathcal{Z}_m \cup \mathcal{Z}_s$ where \mathcal{Z}_m and \mathcal{Z}_s respectively represent multi-view and single-view training data. We write $(X, y) \sim \mathcal{Z}$ as (X, y) sampled uniformly at random from the empirical data set, and denote $N_s = |\mathcal{Z}_s|$. We again for simplicity focus on the setting when $\mu = \frac{1}{\text{poly}(k)}$ and we are given samples $N = k^{1.2}/\mu$ so each label i appears at least $\tilde{\Omega}(1)$ in \mathcal{Z}_s . Our result trivially applies to many other choices of N .*

3.2 Learner Network

We consider a learner network using the following smoothed ReLU activation function $\widetilde{\text{ReLU}}$:

Definition 3.4. *For integer $q \geq 2$ and threshold $\varrho = \frac{1}{\text{polylog}(k)}$, the smoothed function*

$$\widetilde{\text{ReLU}}(z) \stackrel{\text{def}}{=} \begin{cases} 0 & \text{if } z \leq 0; \\ \frac{z^q}{q\varrho^{q-1}} & \text{if } z \in [0, \varrho]; \\ z - (1 - \frac{1}{q})\varrho & \text{if } z \geq \varrho \end{cases}$$

⁵For instance, the marginal distribution of $Z = \sum_{p \in \mathcal{P}_v(X)} z_p$ can be uniform over $[1, 2]$.

⁶For instance, the marginal distribution of $Z = \sum_{p \in \mathcal{P}_v(X)} z_p$ can be uniform over $[0.2, 0.4]$.

Since $\widetilde{\text{ReLU}}$ is smooth we denote its gradient as $\widetilde{\text{ReLU}}'(z)$. We focus on $q = 4$ while our result applies to other constants $q \geq 3$ (see appendix) or most other forms of smoothing. As mentioned in previous section, (smoothed) ReLU has a desired property such that $\widetilde{\text{ReLU}}(z)$ is linear when z is large, but becomes much smaller when z is small. This allows the network to effectively reduce the impact of low-magnitude feature noises from the input patches for better classification.

The *learner network* $F(X) = (F_1(X), \dots, F_k(X)) \in \mathbb{R}^k$ is a two-layer convolutional network parameterized by $w_{i,r} \in \mathbb{R}^d$ for $i \in [k], r \in [m]$, satisfying

$$\forall i \in [k]: \quad F_i(X) = \sum_{r \in [m]} \sum_{p \in [P]} \widetilde{\text{ReLU}}(\langle w_{i,r}, x_p \rangle)$$

Although there exists network with $m = 2$ that can classify the data correctly (e.g. $w_{i,r} = v_{i,r}$ for $r \in [2]$), in this paper, for efficient optimization purpose it is convenient to work on a moderate level of over-parameterization: $m \in [\text{polylog}(k), k]$. Our lower bounds hold for any m in this range and upper bounds hold even for small over-parameterization $m = \text{polylog}(k)$.

Training a single model. We learn the concept class (namely, the labeled data distribution) using gradient descent with learning rate $\eta > 0$, over the cross-entropy loss function L using N training data points $\mathcal{Z} = \{(X_i, y_i)\}_{i \in [N]}$. We denote the empirical loss as:

$$L(F) = \frac{1}{N} \sum_{i \in [N]} L(F; X_i, y_i) = \mathbb{E}_{(X,y) \sim \mathcal{Z}} [L(F; X, y)]$$

where $L(F; X, y) = -\log \frac{e^{Fy(X)}}{\sum_{j \in [k]} e^{F_j(X)}}$. We **randomly initialize** the network F by letting each $w_{i,r}^{(0)} \sim \mathcal{N}(0, \sigma_0^2 I)$ for $\sigma_0^2 = \frac{1}{k}$, which is the most standard initialization people use in practice.

To train a single model, at each iteration t we update using gradient descent (GD):⁷

$$w_{i,r}^{(t+1)} \leftarrow w_{i,r}^{(t)} - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \nabla_{w_{i,r}} L(F^{(t)}; X, y) \quad (3.1)$$

We run the algorithm for $T = \frac{\text{poly}(k)}{\eta}$ iterations. We use $F^{(t)}$ to denote the model F with hidden weights $\{w_{i,r}^{(t)}\}$ at iteration t .

Notations. We denote by $\mathbf{logit}_i(F, X) \stackrel{\text{def}}{=} \frac{e^{F_i(X)}}{\sum_{j \in [k]} e^{F_j(X)}}$. Using this, we can write down

$$\forall i \in [k], r \in [m]: \quad -\nabla_{w_{i,r}} L(F; X, y) = (\mathbf{1}_{i \neq y} - \mathbf{logit}_i(F, X)) \nabla_{w_{i,r}} F_i(X) .$$

4 Main Theorems and Explanations

We now state the main theorems in this paper.⁸ Recall the learner network and its learning process are given in Section 3.2, and the data distribution is in Section 3.1.

Theorem 1 (single model). *For every sufficiently large $k > 0$, every $m \in [\text{polylog}(k), k]$, every $\eta \leq \frac{1}{\text{poly}(k)}$, suppose we train a single model using the gradient descent update (3.1) starting from the random initialization defined in Section 3.2, then after $T = \frac{\text{poly}(k)}{\eta}$ many iterations, with probability $\geq 1 - e^{-\Omega(\log^2 k)}$, the model $F^{(T)}$ satisfies:*

- (training accuracy is perfect): meaning for all $(X, y) \in \mathcal{Z}$, all $i \in [k] \setminus \{y\}$: $F_y^{(T)}(X) > F_i^{(T)}(X)$.

⁷Our result also trivially extends to the case when there is a weight decay (i.e. ℓ_2 regularizer): $w_{i,r}^{(t+1)} \leftarrow (1 - \eta\lambda)w_{i,r}^{(t)} - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \nabla_{w_{i,r}} L(F^{(t)}; X, y)$ as long as λ is not too large. We keep this basic version without weight decay to simplify the analysis.

⁸We shall restate these theorems in the appendix with more details and wider range of parameters.

- (test accuracy is consistently bad): meaning that:

$$\Pr_{(X,y)\sim\mathcal{D}} [\exists i \in [k] \setminus \{y\} : F_y^{(T)}(X) < F_i^{(T)}(X)] \in [0.49\mu, 0.51\mu] .$$

We shall give *technical intuitions* about why Theorem 1 holds in Appendix C. But, at a high-level, we shall construct a “lottery winning” set $\mathcal{M} \subseteq [k] \times [2]$ of cardinality $|\mathcal{M}| \in [k(1 - o(1)), k]$. It only depends on the random initialization of F . Then, with some effort we can prove that, for every $(i, \ell) \in \mathcal{M}$, at the end of the training $F^{(T)}$ will learn feature $v_{i,\ell}$ but not learn feature $v_{i,3-\ell}$. This means for those single-view data (X, y) with $y = i$ and $\widehat{\ell}(X) = 3 - \ell$, the final network $F^{(T)}$ will predict its label wrong. This is why the final test accuracy is around 0.5μ .

Note the property that test accuracy consistently belongs to the range $[0.49\mu, 0.51\mu]$ should be reminiscent of message ⑤ in Figure 6, where multiple single models, although starting from different random initialization, in practice does have a relatively small variance in test accuracies.

Ensemble. Suppose $\{F^{[\ell]}\}_{\ell \in [K]}$ are $K = \widetilde{\Omega}(1)$ independently trained models of F with $m = \text{polylog}(k)$ for $T = O(\frac{\text{poly}(k)}{\eta})$ iterations (i.e., the same setting as Theorem 1 except we only need a small over-parameterization $m = \text{polylog}(k)$). Let us define their ensemble

$$G(X) = \frac{\widetilde{\Theta}(1)}{K} \sum_{\ell} F^{[\ell]}(X) \tag{4.1}$$

Our next theorem states that the ensemble model has much higher test accuracy.

Theorem 2 (ensemble). *In the same setting as Theorem 1 except now we only need a small $m = \text{polylog}(k)$, we have for the ensemble model G in (4.1), with probability at least $1 - e^{-\Omega(\log^2 k)}$:*

- (training accuracy is perfect): meaning for all $(X, y) \in \mathcal{Z}$, for all $i \in [k] \setminus \{y\}$: $G_y(X) > G_i(X)$.
- (test accuracy is almost perfect): meaning that:

$$\Pr_{(X,y)\sim\mathcal{D}} [\exists i \in [k] \setminus \{y\} : G_y(X) < G_i(X)] \leq 0.001\mu .$$

As we discussed in Section 2.3, the reason Theorem 2 holds attributes to the fact that those lottery winning sets \mathcal{M} depend on the random initialization of the networks; and therefore, when multiple models are put together, their “union” of \mathcal{M} shall cover all possible features $\{v_{i,\ell}\}_{(i,\ell) \in [k] \times [2]}$. Moreover, our theorem only requires individual $K = \widetilde{\Omega}(1)$ models for ensemble, which is indeed “averaging the output of a few independently trained models”.

Roadmap. We shall restate and prove the general versions of Theorem 1 and 2 in Appendix E, after establishing core lemmas in Appendix C and D.

4.1 Knowledge Distillation for Ensemble

We consider a knowledge distillation algorithm given the existing ensemble model G (see (4.1)) as follows. For every label $i \in [k]$, let us define the truncated scaled logit as (for $\tau = \frac{1}{\log^2 k}$):

$$\mathbf{logit}_i^\tau(F, X) = \frac{e^{\min\{\tau^2 F_i(X), 1\}/\tau}}{\sum_{j \in [k]} e^{\min\{\tau^2 F_j(X), 1\}/\tau}} \tag{4.2}$$

(This should be reminiscent of the logit function with temperature used by the original knowledge distillation work [42]; we use truncation instead which is easier to analyze.)

Now, we train a new network F from random initialization (where the randomness is independent of all of those used in $F^{[l]}$). At every iteration t , we update each weight $w_{i,r}$ by:

$$w_{i,r}^{(t+1)} = w_{i,r}^{(t)} - \eta \nabla_{w_{i,r}} L(F^{(t)}) - \eta' \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\left(\mathbf{logit}_i^\tau(F^{(t)}, X) - \mathbf{logit}_i^\tau(G, X) \right)^- \nabla_{w_{i,r}} F_i^{(t)}(X) \right] \quad (4.3)$$

Notation. Throughout the paper we denote by $[a]^+ = \max\{0, a\}$ and $[a]^- = \min\{0, a\}$.

This knowledge distillation method (4.3) is almost identical to the one used in the original work [42], except we use a truncation during the training to make it more (theoretically) stable. Moreover, we update the distillation objective using a larger learning rate η' comparing to η of the cross-entropy objective. This is also consistent with the training schedule used in [42].

Let $F^{(t)}$ be the resulting network obtained by (4.3) at iteration t . We have the following theorem:

Theorem 3 (ensemble distillation). *Consider the distillation algorithm (4.3) in which G is the ensemble model defined in (4.1). For every $k > 0$, for $m = \text{polylog}(k)$, for every $\eta \leq \frac{1}{\text{poly}(k)}$, setting $\eta' = \eta \text{poly}(k)$, after $T = \frac{\text{poly}(k)}{\eta}$ many iterations with probability at least $1 - e^{-\Omega(\log^2 k)}$, for at least 90% of the iterations $t \leq T$:*

- (training accuracy is perfect): meaning for all $(X, y) \in \mathcal{Z}$, all $i \in [k] \setminus \{y\}$: $F_y^{(t)}(X) > F_i^{(t)}(X)$.
- (test accuracy is almost perfect): meaning that:

$$\Pr_{(X,y) \sim \mathcal{D}} [\exists i \in [k] \setminus \{y\}: F_y^{(t)}(X) < F_i^{(t)}(X)] \leq 0.001\mu .$$

We shall restate the general version of Theorem 3 in Appendix F, and prove it in Appendix G.

Remark. Theorem 3 necessarily means that the distilled model F has learned all the features $\{v_{i,\ell}\}_{(i,\ell) \in [k] \times [2]}$ from the ensemble model G . This is consistent with our empirical findings in Figure 8: if one trains multiple individual models using knowledge distillation with different random seeds, then their ensemble gives no further performance boost.

4.2 Self Knowledge Distillation as Implicit Ensemble

Self distillation [35, 63] refers to training a single model to match the output of another single model. In this paper we also show that self-distillation can also improve test accuracy under our multi-view setting. Let us consider the following self-distillation algorithm.

Let us now consider $F = F^{(T)}, G = G^{(T)}$ be two single models trained in the same setting as Theorem 1 using independent random initializations (for simplicity, we override to notation a bit, so here G is a single model to be distilled from, instead of the ensemble). We scale them up by a small factor $\tilde{\Theta}(1)$ similar to (4.1). Then, starting from $F^{(T)}$, we apply the following updates for another T' iterations:

$$w_{i,r}^{(t+1)} = w_{i,r}^{(t)} - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left(\left(\mathbf{logit}_i^\tau(F^{(t)}, X) - \mathbf{logit}_i^\tau(G, X) \right)^- \nabla_{w_{i,r}} F_i^{(t)}(X) \right) \quad (4.4)$$

This objective is considered as “self-distillation” since G is an individual model (trained using an identical manner as F , only from a different random initialization). In particular, if $F = G$, then $\mathbf{logit}_i^\tau(F^{(t)}, X) - \mathbf{logit}_i^\tau(G, X) = 0$ so the weights are no longer updated. However, as we will actually prove, this training objective will actually learn an F that *has better generalization comparing to G* .

This time, *for simplicity*, let us make the following additional assumption on the data:

Assumption 4.1 (balanced \mathcal{D}_m). *In Def. 3.1, for multi-view data (X, y) , we additionally assume that the marginal distributions of $\sum_{p \in \mathcal{P}_v(X)} z_p^q \in [1, 1 + o(1)]$ for $v \in \{v_{y,1}, v_{y,2}\}$.*

The rationale for this assumption is quite simple. Suppose \mathcal{M} is the aforementioned “lottery winning” set of training a single model without knowledge distillation. Assumption 4.1 will ensure that each $(i, 1)$ and $(i, 2)$ will belong to \mathcal{M}_F with relatively equal probability. If we train two models F and G , their combined lottery winning set $\mathcal{M}_F \cup \mathcal{M}_G$ shall be of cardinality around $\frac{3}{2}k(1 - o(1))$. Therefore, if we can distill the knowledge of G to F , the test accuracy can be improved from $1 - \frac{1}{2}\mu$ to $1 - \frac{1}{4}\mu$. See the following theorem:⁹

Theorem 4 (self-distillation). *Under this additional Assumption 4.1, consider the distillation algorithm (4.4) where G is an independently trained single model (in the same setting as Theorem 1). For every $k > 0$, every $m \in [\text{polylog}(k), k]$, every $\eta \leq \frac{1}{\text{poly}(k)}$, after $T' = \frac{\text{poly}(k)}{\eta}$ many iterations of algorithm (4.4), with probability at least $1 - e^{-\Omega(\log^2 k)}$:*

- (training accuracy is perfect): meaning for all $(X, y) \in \mathcal{Z}$, all $i \neq y$: $F_y^{(T+T')}(X) > F_i^{(T+T')}(X)$.
- (test accuracy is better): meaning that:

$$\Pr_{(X,y) \sim \mathcal{D}} [\exists i \in [k] \setminus \{y\} : F_y^{(T+T')}(X) < F_i^{(T+T')}(X)] \leq 0.26\mu$$

Recall from Theorem 1, all individual models should have test error at least $0.49\mu \gg 0.26\mu$. Hence the model F generalizes better (comparing to both the original model F before self-distillation and the individual model G) after self-distillation to the individual model G . We shall restate the general version of Theorem 4 in Appendix F, and prove it in Appendix G.

Why does self-distillation improve test accuracy? Self-distillation is performing implicit ensemble + knowledge distillation. As the main idea of behind the proof, we actually show that self-distillation is performing implicit ensemble together with knowledge distillation. In particular, let $\mathcal{M}_F, \mathcal{M}_G \subseteq V$ be the features learned by individual models F, G starting from (independent) random initializations $W_F^{(0)}$ and $W_G^{(0)}$ respectively when trained on the original data set, now, if we further train the individual model F to match the output of individual model G , F is actually going to learn a larger set of features $\mathcal{M}_G \cup \mathcal{M}_F$, where features in \mathcal{M}_F come from gradient of the original objective, and features in \mathcal{M}_G come from the gradient of the knowledge distillation objective w.r.t. G . This is equivalent to first ensemble F and G , then train an additional model H from random initialization to match the ensemble— Self-distillation implicitly merge “ensemble individual models F, G and distill the ensemble to another individual model H ” into “ensemble individual models F, G and distill the ensemble to the individual model F ” since F and H have the same structure. Then eventually it is merged directly into “training an individual model F via self-distillation to match the output of an individual model G ”.

5 Our Empirical Results at a High Level

On the *empirical side*, to further justify our approach studying ensemble and knowledge distillation indeep learning, we show:

⁹One can trivially relax Assumption 4.1 so that the two views have different distributions but with a constant ratio between their expectations; in this way the improved accuracy is no longer $1 - \frac{1}{4}\mu$ but shall depend on this constant ratio. For simplicity of this paper, we do not state the result of that more general case.

finite-width neural kernel models	CIFAR10 test accuracy					CIFAR100 test accuracy				
	single model (best of 10)	ensemble (over 10)	train $\sum_{\ell} f_{\ell}$ (over 10)	knowledge distillation	self-distill	single model (best of 10)	ensemble (over 10)	train $\sum_{\ell} f_{\ell}$ (over 10)	knowledge distillation	self-distill
SimpleCNN-10-3-NTK	64.36%	67.38%	69.37%	64.63%	65.24%	out of memory <small>◦ due to memory restriction, trained $\sum_{\ell} f_{\ell}$ over fewer than 10 models.</small>				
ResNet10-2-NTK	69.15%	73.29%	74.71%	68.82%	66.09%					
ResNet16-2-NTK	68.32%	73.79%	74.62% (over 7) [◦]	66.12%	70.61%					
ResNet16-5-NTK	74.21%	78.46%	out of memory	70.23%	75.66%					
ResNet10-10-NTK	76.66%	80.39%	out of memory	77.25%	74.46%					
SimpleCNN10-6-NTK'	59.92%	63.43%	65.69%	59.12%	57.81%	18.99%	26.54%	28.28%	18.27%	18.40%
ResNet10-4-NTK'	66.68%	70.54%	72.86%	66.01%	62.91%	31.90%	38.32%	41.47%	31.38%	27.64%
SimpleCNN-10-6-GP	30.48%	35.33%	40.08%	29.43%	29.10%	9.82%	11.82%	12.22%	8.95%	9.33%
ResNet-10-4-GP	42.17%	48.60%	53.17%	39.45%	41.63%	18.89%	22.92%	25.88%	16.91%	16.59%



Message ①: for neural kernel methods, ensemble helps on improving test accuracies, but ensemble is *not better than* training the sum of the individuals directly. In other words, the benefit of using ensemble here merely comes from the richer set of prescribed features.

Message ②: for neural kernel methods, the superior test performance of ensemble *cannot be distilled* into a single model.

Message ③: for neural kernel methods, self-distillation is generally *no better than* a single model's test performance.

neural networks	single model (over 10)	ensemble (over 10)	train $\sum_{\ell} f_{\ell}$ (over 10)	knowledge distillation	self-distill	single model (over 10)	ensemble (over 10)	train $\sum_{\ell} f_{\ell}$ (over 10)	knowledge distillation	self-distill
	ResNet-28-2	95.22±0.14%	96.33%	95.02%	96.16%	95.78%	76.38±0.23%	81.13%	73.18%	79.03%
ResNet-34	93.65±0.19%	94.97%	93.12%	94.59%	94.21%	71.66±0.43%	76.85%	68.88%	73.74%	73.14%
ResNet-34-2	95.45±0.14%	96.55%	95.00%	96.08%	95.86%	77.01±0.35%	81.48%	72.99%	79.23%	79.07%
ResNet-16-10	96.08±0.16%	96.80%	95.88% (over 6) [◦]	96.81%	96.62%	80.03±0.17%	83.18%	80.53% (over 6) [◦]	82.67%	82.25%
ResNet-22-10	96.44±0.09%	97.12%	96.41% (over 5) [◦]	97.09%	97.05%	81.17±0.23%	84.33%	81.59% (over 5) [◦]	83.71%	83.26%
ResNet-28-10	96.70±0.21%	97.20%	96.46% (over 4) [◦]	97.22%	97.13%	81.51±0.16%	84.69%	81.83% (over 4) [◦]	83.81%	83.56%



Message ④: for neural nets, ensemble helps on improving test accuracies, **and** this accuracy gain *cannot* be matched by training the sum of the individuals directly. In other words, the benefit of using ensemble comes from somewhere other than enlarging the model.

Message ⑤: for neural nets, the superior test performance of ensemble *can be distilled* into single model by a large extent.

Message ⑥: for neural nets, self-distillation *clearly improves* the test performance of single models.

Message ⑦: for neural nets, the superior performance of ensemble *does not* come from the variance of test accuracies in single models.

Figure 6: Comparing the performances of (1) training 10 independent single models f_1, \dots, f_{10} , (2) their ensemble, (3) training $f_1 + \dots + f_{10}$ directly, (4) knowledge distillation of the ensemble into a single model, and (5) training a single model using self-distillation.

(NTK' = the original finite-width neural network first-order approximation [6], NTK = the more popular variant where for each output label one learns a different linear function over the NTK features (e.g. [8]), and GP = training only the last layer of a finite-width random neural network [23]. All the neural networks in these experiments are trained to $\sim 100\%$ training accuracy, and the single model performances match the state-of-the-art for these models on CIFAR-10/100. For experiment details, see Appendix B.1.)

- **Ensemble (i.e. model averaging) in deep learning works very differently from ensemble in random feature mappings** — in particular, different from the neural tangent kernel (NTK) approach [3, 5, 6, 8, 9, 19, 24, 26, 27, 39, 40, 44, 55, 59, 86, 92].

Let us do a thought experiment. If ensemble works, can we obtain the same test performance of ensemble by training the sum of over L neural networks $G = \sum_{\ell \in [L]} f_{\ell}$ directly? (As opposite to training each f_{ℓ} independently and then average them.) Actually, this “direct training” approach in deep learning is unable to improve the test accuracy even comparing to single models, not to say the ensemble model. See Figure 6.

In contrast, when each f_{ℓ} is a linear function over random feature mappings (e.g., the NTK feature mappings given by the random initialization of the network), although the ensemble of these random feature models *does improve* test accuracy, training directly over $F = \sum_{\ell \in [L]} f_{\ell}$ gives even superior test accuracy comparing to the ensemble. See Figure 6.

- **Knowledge distillation works for “ensemble of neural networks” but does not work**

		no label noise				with 10% label noise			
		uniform sampling		rejection sampling		uniform sampling		rejection sampling	
		gaussian input	mixture of gaussian	gaussian input	mixture of gaussian	gaussian input	mixture of gaussian	gaussian input	mixture of gaussian
data without margin	linear	80.3% (79.6%)	80.7% (80.1%)	78.9% (78.6%)	80.7% (80.7%)	74.3% (74.1%)	73.6% (74.0%)	72.9% (72.2%)	74.2% (73.7%)
	fc2	67.7% (65.1%)	67.7% (64.9%)	66.3% (64.5%)	67.6% (66.9%)	64.3% (63.2%)	70.1% (66.7%)	64.6% (63.5%)	66.2% (63.3%)
	fc3	68.9% (69.0%)	64.0% (64.4%)	76.8% (76.6%)	73.2% (73.1%)	66.5% (66.4%)	63.0% (62.4%)	72.5% (72.0%)	78.1% (78.6%)
	res3	69.1% (68.0%)	70.7% (71.2%)	69.3% (69.0%)	69.9% (69.4%)	68.7% (65.9%)	66.9% (63.8%)	68.1% (68.1%)	68.8% (69.5%)
	conv2	65.4% (65.7%)	67.0% (66.8%)	68.3% (68.2%)	68.3% (68.2%)	67.1% (66.2%)	65.1% (65.5%)	65.8% (66.0%)	67.5% (67.9%)
	conv3	68.7% (68.5%)	70.7% (71.2%)	77.8% (77.1%)	80.3% (80.3%)	67.5% (68.2%)	67.7% (67.5%)	73.6% (73.4%)	71.8% (72.1%)
	resconv3	78.3% (78.3%)	79.6% (79.3%)	83.8% (82.6%)	82.1% (81.8%)	74.1% (73.9%)	73.4% (73.1%)	78.7% (78.5%)	79.2% (78.5%)
	linear	79.0% (78.0%)	79.0% (77.2%)	78.4% (77.3%)	80.0% (80.0%)	82.1% (81.7%)	80.7% (80.0%)	81.6% (80.2%)	84.1% (82.4%)
	fc2	80.6% (79.0%)	80.4% (78.4%)	78.4% (76.6%)	78.4% (77.0%)	77.4% (75.3%)	73.7% (73.9%)	74.7% (72.2%)	75.7% (71.4%)
data with margin	fc3	76.0% (75.8%)	80.4% (80.1%)	76.9% (77.0%)	73.3% (72.9%)	70.7% (70.8%)	73.9% (74.6%)	70.4% (70.5%)	67.5% (66.4%)
	res3	80.7% (80.8%)	84.7% (83.9%)	84.6% (84.0%)	84.4% (83.7%)	75.5% (74.0%)	76.9% (76.4%)	76.8% (74.8%)	76.4% (74.5%)
	conv2	70.6% (70.3%)	74.5% (73.6%)	67.8% (67.8%)	69.6% (69.4%)	68.8% (67.5%)	73.3% (71.7%)	67.0% (66.6%)	67.6% (67.2%)
	conv3	76.2% (76.2%)	75.3% (76.1%)	79.6% (79.1%)	84.3% (83.6%)	72.2% (72.1%)	81.2% (81.3%)	73.0% (72.3%)	74.4% (75.1%)
	resconv3	92.1% (91.7%)	92.1% (92.3%)	93.9% (93.8%)	95.5% (95.1%)	85.2% (85.1%)	83.5% (84.4%)	87.3% (86.8%)	85.6% (85.9%)

Figure 7: When data is Gaussian-like, and when the target label is generated by some fully-connected(fc) / residual(res) / convolutional(conv) network, *ensemble does not* improve test accuracy. “xx % (yy %)” means xx% accuracy for single model and yy% for ensemble. More experiments in Appendix B.4 (Figure 10 and 11).

for “ensemble of random feature mappings” on standard data sets.

When f_ℓ is a linear function over random feature mappings (e.g., the NTK feature mappings), the superior test performance of ensemble *cannot be distilled* into a single model. In contrast, in deep learning, such superior performance can be distilled into a single model using [42]. The situation is similar for *self-distillation*, where it hardly works on improving test performance for neural kernel methods, but works quite well for real neural networks. See Figure 6. Together with the first point, experiments suggest that to understand the benefit of ensemble and knowledge distillation in deep learning, it is perhaps **inevitable to study deep learning as a feature learning process**, instead of feature selection process (e.g. NTK or other neural kernel methods) where only the linear combinations of prescribed features are trained.

- Some prior works attribute the benefit of ensemble to **reducing the variance** of individual solutions [15, 62, 64, 82, 83] due to label noise or non-convex landscape of the training objective. We observe that this *may not be the cause for “ensemble in deep learning” to work*.
 - For standard deep learning datasets (e.g. CIFAR-10/100), individual neural networks (e.g. ResNets) trained by SGD typically have already converged to global optimas with 100% training accuracy and nearly zero training loss (no failure in training).
 - For standard deep learning datasets (e.g. CIFAR-10/100), ensemble helps even when there is essentially no label noise. In contrast, in our synthetic experiment Figure 7, ensemble does not help on Gaussian-like data even when there is label noise.
 - For standard neural networks (e.g. ResNets) trained on standard data set (e.g. CIFAR-10/100), when all the individual models are well-trained with the same learning rate/weight decay and only differ by their random seeds, there is *almost no variance* in test accuracy for individual models (e.g. 0.1 ~ 0.4% std on CIFAR-100, see Figure 6). Hence with high probability, all individual models are learned *almost equally well* (no failure in generalization), yet ensemble still offers a huge benefit in test performance.
 - For neural networks trained on our Gaussian-like data, there is relatively *higher variance* (e.g. 0.5 ~ 1.0% std in test accuracies, see Figure 12 in the appendix), yet ensemble offers no benefit at all.
 - For individual neural networks trained using knowledge distillation with different random seeds, ensemble does not improve their test accuracy by much (see Figure 8) — despite

	CIFAR10 test accuracy				CIFAR100 test accuracy			
	single model (over 10)	ensemble (over 10)	10 runs of knowledge distill	ensemble over knowledge distill	single model (over 10)	ensemble (over 10)	10 runs of knowledge distill	ensemble over knowledge distill
ResNet-28-2	95.22±0.14%	96.33%	95.89±0.07%	96.21%	76.38±0.23%	81.13%	78.94±0.21%	80.35%
ResNet-34	93.65±0.19%	94.97%	94.37±0.13%	94.88%	71.66±0.43%	76.85%	73.57±0.34%	75.60%
ResNet-34-2	95.45±0.14%	96.55%	96.00±0.12%	96.42%	77.01±0.35%	81.48%	79.43±0.23%	81.56%
ResNet-16-10	96.08±0.16%	96.80%	96.73±0.07%	96.76%	80.03±0.17%	83.18%	82.51±0.14%	83.36%
ResNet-22-10	96.44±0.09%	97.12%	97.01±0.09%	97.09%	81.17±0.23%	84.33%	83.54±0.19%	84.27%
ResNet-28-10	96.70±0.21%	97.20%	97.06±0.08%	97.24%	81.51±0.16%	84.69%	83.75±0.16%	84.87%

Message ①: an ensemble over *single models* (independently trained) can be distilled into a single model with moderate accuracy loss.
Message ②: an ensemble over models *after knowledge distillation* *does not improve accuracy by much* – in fact, not exceeding the ensemble accuracy of the original single models ③ – despite the training objective is still non-convex and different random seeds are used. This means, knowledge distillation models (i.e. simply matching the soft labels) *have learned most of the features* from the ensemble, and have less variety comparing to the original single models. This also means that “(huge) non-convexity” in neural networks and SGD with “different random seeds” *even together do not guarantee ensemble advantage unconditionally*; the structure of the data (and hard labels) is extremely important for ensemble to work as we mainly focus on in this paper.

Figure 8: Single models (+ their ensemble) vs. Knowledge distillations (+ their ensemble). Details in Appendix B.2.

that the knowledge distillation objective is “as non-convex as” the original training objective and only the training labels are changed from hard to soft labels.

- **Special structure in data** (such as the “multi-view” structure we shall introduce) is **arguably necessary for ensemble to work**. Over certain data sets with no multi-view structure, ensemble *does not improve the test-time performance in deep learning* — despite having a non-convex training objective and different random seeds are used. See Figure 7. In contrast, real-life data sets such as CIFAR-10/100 do have the multi-view structure, moreover, standard neural networks such as ResNet do utilize this structure during the training process in the same way as we show in our theory. See Figure 4.
- For neural networks, **knowledge distillation has learned most of the features from the ensemble, and the use of *hard labels* to train individual models is a key for why ensemble works in deep learning**.

Specifically, as in Figure 8, if one evaluates an ensemble over models that are independently at random trained from knowledge distillation (i.e., using soft labels), its performance does not exceed the ensemble over the original single models. This means, models trained via knowledge distillation have learned most of the features from the ensemble, and has less variety comparing to the original models. We shall see this is consistent with our theoretical result.

6 Conclusion and Discussion

In this work, we have shown, to the best of our knowledge, the first theoretical result towards understanding how ensemble work in deep learning. As our main contribution, we provide empirical evidence that ensemble might work very differently in deep learning comparing to ensemble in random feature models. Moreover, ensemble does not always improve test accuracy in deep learning, especially when the input data comes from Gaussian-like distribution.

Motivated by these empirical observations, we propose a generic structure of the data we refer to as multi-view, and prove that ensemble improves test accuracy for two-layer neural networks in this setting. Moreover, we also prove that ensemble model can be distilled into a single model. Meaning that, through training a single model to “simulate” the output of the ensemble over the

CIFAR100	# input channels		original	split to 2	split to 4	split to 8	avg over 2	avg over 4	avg over 8
ResNet-28 (a)	16		70.44±0.29%	68.77±0.25%	66.70±0.66%	-	69.00±0.43%	66.45±0.15%	-
ResNet-28 (b)	32		70.49±0.29%	67.62±0.89%	63.28±0.50%	-	67.99±0.15%	63.89±0.31%	-
ResNet-28-2 (a)	32	single	76.09±0.23%	74.50±0.68%	72.47±1.78%	70.84±1.32%	75.31±0.23%	73.69±0.34%	71.60±0.34%
ResNet-28-2 (b)	64	model	76.12±0.23%	74.88±0.22%	72.81±0.29%	69.21±0.49%	74.58±0.33%	72.71±0.29%	68.85±0.28%
ResNet-28-4 (a)	64	test	79.10±0.18%	78.57±0.29%	77.94±0.43%	76.88±0.35%	78.42±0.35%	78.14±0.16%	77.52±0.20%
ResNet-28-4 (b)	128	accuracy	78.53±0.16%	77.72±0.20%	76.62±0.29%	74.93±0.40%	77.95±0.22%	76.88±0.33%	75.17±0.25%
ResNet-28-10 (a)	160		81.23±0.23%	81.03±0.17%	80.53±0.09%	80.12±0.26%	80.58±0.28%	81.06±0.22%	80.63±0.22%
ResNet-28-10 (b)	320		80.76±0.27%	80.41±0.24%	80.09±0.16%	79.02±0.22%	80.54±0.23%	80.03±0.16%	79.38±0.27%
ResNet-28 (a)	16		75.52%	74.07%	73.63%	-	74.05%	70.98%	-
ResNet-28 (b)	32		74.47%	73.58%	72.17%	-	71.97%	68.03%	-
ResNet-28-2 (a)	32	ensemble	80.33%	79.73%	79.58%	78.75%	79.24%	78.19%	76.31%
ResNet-28-2 (b)	64	model	79.63%	80.18%	79.17%	78.20%	78.42%	76.81%	72.90%
ResNet-28-4 (a)	64	test	82.64%	82.81%	82.56%	82.24%	82.26%	82.12%	81.71%
ResNet-28-4 (b)	128	accuracy	81.84%	82.06%	81.89%	81.74%	81.28%	80.63%	79.14%
ResNet-28-10 (a)	160		84.05%	84.08%	83.65%	83.51%	83.79%	84.12%	83.69%
ResNet-28-10 (b)	320		83.10%	83.40%	83.81%	83.53%	83.21%	83.00%	82.19%

Figure 9: Justify the multi-view hypothesis in practice. We regard some intermediate layer of a pre-trained ResNet as “input” with multiple channels (this pre-trained network stays fixed and shared for all individual models). Then, we train a new model either starting from this input (i.e. the “original” column), or from a fraction of the input (i.e., “split into 4” means using only 1/4 of the input channels), or from an average of the input (i.e., “average over 4” means averaging every four channels). Details in Appendix B.3.

Observation 1. Even when we significantly collapse the input channels (through averaging or throwing away most of them), most of the single model test accuracies do not drop by much. Moreover, it’s known [65] that in ResNet, most channels are indeed learning different features (views) of the input, also see Figure 3 for an illustration. This indicates that many data can be classified correctly using different views.

Observation 2. Even when single model accuracy drops noticeably, ensemble accuracy does not change by much. We believe this is a strong evidence that there are multiple views in the data (even at intermediate layers), and **ensemble can collect all of them even when some models have missing views.**

same training data set, single model is able to match the test accuracy of the ensemble, and thus being superior to any single model that is clean, directly trained on the original data’s labels.

We believe that our framework can be applied to other settings as well, for example, data augmentation using random cropping could be potentially regarded as another way to enforce the network to learn “multi-views”. We hope that our new theoretical insights on how neural networks pick up features during training can also help in practice design new, principled approach to improve test accuracy of a neural network, potentially matching that of the ensemble.

APPENDIX I: MISSING DETAILS

In Section A, we give a formal definition of the data distribution: this expands the earlier Section 3.1 by giving more discussions and the full specifications of the noise parameters.

In Section B, we give the experiment setups and some additional experiments.

Appendix II gives the full proofs, but it will start with Section C for technical intuitions and the proof plan.

A Data Distribution and Notations (Full Version)

We consider learning a k -class classification problem over P -patch inputs, where each patch has dimension d . In symbols, each labelled data is represented by (X, y) where $X = (x_1, x_2, \dots, x_P) \in (\mathbb{R}^d)^P$ is the data vector and $y \in [k]$ is the data label. For simplicity, we focus on the case when $P = k^2$, and $d = \text{poly}(k)$ for a large polynomial.

We consider the setting when k is sufficiently large.¹⁰ We use “w.h.p.” to denote with probability at least $1 - e^{-\Omega(\log^2 k)}$, and use $\tilde{O}, \tilde{\Theta}, \tilde{\Omega}$ notions to hide polylogarithmic factors in k .

We first assume that each label class $j \in [k]$ has multiple associated features, say *two features for the simplicity of math*, represented by unit *feature vectors* $v_{j,1}, v_{j,2} \in \mathbb{R}^d$. For notation simplicity, we assume that all the features are orthogonal, namely,

$$\forall j, j' \in [k], \forall \ell, \ell' \in [2], \|v_{j,\ell}\|_2 = 1 \quad \text{and} \quad v_{j,\ell} \perp v_{j',\ell'} \quad \text{when} \quad (j, \ell) \neq (j', \ell')$$

although our work also extends to the “incoherent” case trivially. We denote by

$$\mathcal{V} \stackrel{\text{def}}{=} \{v_{j,1}, v_{j,2}\}_{j \in [k]} \quad \text{the set of all features.}$$

We now consider the following data and label distribution. Let C_p be a global constant, $s \in [1, k^{0.2}]$ be a global parameter to control feature sparsity, $\sigma_p = \frac{1}{\sqrt{d \text{polylog}(k)}}$ be a parameter to control magnitude of the *random noise*, $\gamma = \frac{1}{k^{1.5}}$ be a parameter to control the *feature noise*. (Our proof in the appendix holds for a wider range of γ .)

To be concise, we define the *multi-view distribution* \mathcal{D}_m and *single-view distribution* \mathcal{D}_s together.

Definition 3.1 (data distributions \mathcal{D}_m and \mathcal{D}_s). *Given $\mathcal{D} \in \{\mathcal{D}_m, \mathcal{D}_s\}$, we define $(X, y) \sim \mathcal{D}$ as follows. First choose the label $y \in [k]$ uniformly at random. Then, the data vector X is generated as follows (also illustrated in Figure 5).*

1. Denote $\mathcal{V}(X) = \{v_{y,1}, v_{y,2}\} \cup \mathcal{V}'$ as the set of feature vectors used in this data vector X , where \mathcal{V}' is a set of features uniformly sampled from $\{v_{j',1}, v_{j',2}\}_{j' \in [k] \setminus \{y\}}$, each with probability $\frac{s}{k}$.

◇ comment: (X, y) shall be primarily supported on two main features $v_{y,1}, v_{y,2}$ and $\sim O(s)$ minor features

2. For each $v \in \mathcal{V}(X)$, pick C_p many disjoint patches in $[P]$ and denote it as $\mathcal{P}_v(X) \subset [P]$ (the distribution of these patches can be arbitrary). We denote $\mathcal{P}(X) = \cup_{v \in \mathcal{V}(X)} \mathcal{P}_v(X)$.

◇ comment: the weights of X on each feature v shall be written on patches in $\mathcal{P}_v(X)$

3. If $\mathcal{D} = \mathcal{D}_s$ is the single-view distribution, pick a value $\hat{\ell} = \hat{\ell}(X) \in [2]$ uniformly at random.

¹⁰If we want to work with fixed k , say $k = 2$, our theorem can also be modified to that setting by increasing the number of features per class. In this case, a subset of features per class will be learned by each individual neural network. We keep our current setting with two features to simplify the notations.

4. For each $v \in \mathcal{V}(X)$ and $p \in \mathcal{P}_v(X)$, we set

$$x_p = z_p v + \sum_{v' \in \mathcal{V}} \alpha_{p,v'} v' + \xi_p \in \mathbb{R}^d$$

Above, each $\alpha_{p,v'} \in [0, \gamma]$ is the feature noise, and $\xi_p \sim \mathcal{N}(0, \sigma_p^2 \mathbf{I})$ is an (independent) random Gaussian noise. The coefficients $z_p \geq 0$ satisfy that:

In the case of multi-view distribution $\mathcal{D} = \mathcal{D}_m$,

- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [1, O(1)]$ when $v \in \{v_{y,1}, v_{y,2}\}$,
and the marginal distribution of $\sum_{p \in \mathcal{P}_v(X)} z_p$ is left-close,¹¹
- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [\Omega(1), 0.4]$ when $v \in \mathcal{V}(X) \setminus \{v_{y,1}, v_{y,2}\}$,
and the marginal distribution of $\sum_{p \in \mathcal{P}_v(X)} z_p$ is right-close.¹²

◇ comment: total weights on features $v_{y,1}, v_{y,2}$ are larger than those on minor features $\mathcal{V}(X) \setminus \{v_{y,1}, v_{y,2}\}$

In the case of single-view distribution $\mathcal{D} = \mathcal{D}_s$,

- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [1, O(1)]$ when $v = v_{y,\hat{\ell}}$,
- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [\rho, O(\rho)]$ when $v = v_{y,3-\hat{\ell}}$, ◇ comment: we consider $\rho = k^{-0.01}$ for simplicity
- $\sum_{p \in \mathcal{P}_v(X)} z_p \in [\Omega(\Gamma), \Gamma]$ when $v \in \mathcal{V}(X) \setminus \{v_{y,1}, v_{y,2}\}$. we consider $\Gamma = \frac{1}{\text{polylog}(k)}$ for simplicity

◇ comment: total weight on feature $v_{y,\hat{\ell}}$ is much larger than those on $v_{y,3-\hat{\ell}}$ or minor features

5. For each $p \in [P] \setminus \mathcal{P}(X)$, we set:

$$x_p = \sum_{v' \in \mathcal{V}} \alpha_{p,v'} v' + \xi_p$$

where $\alpha_{p,v'} \in [0, \gamma]$ is the feature noise and $\xi_p \sim \mathcal{N}(0, \frac{\gamma^2 k^2}{d} \mathbf{I})$ is (independent) Gaussian noise.

Remark A.1. The distribution of how we pick $\mathcal{P}(X)$ and how to assign $\sum_{p \in \mathcal{P}_v(X)} z_p$ to each patch in $p \in \mathcal{P}_v(X)$ can be arbitrary (and can depend on other randomness in the data as well). Except the marginal distributions of the sum of some z_p 's are left-close or right-close, but we also do not have other restrictions on how z_p 's are distributed within the summation. In particular, we have allowed **different features** $v_{j,1}, v_{j,2}$ to show up with **different weights** in the data (for example, for multi-view data, some view $v_{y,1}$ can consistently have larger z_p comparing to $v_{y,2}$.) Yet, we shall prove that the order to learn these features by the learner network *can still be flipped* depending on the randomness of network initialization. We also do not have any restriction on the distribution of the feature noise $\alpha_{p,v'}$ (they can depend on other randomness of the data distribution as well).

Generality and significance of our data distribution. Our setting is *tied to a down-sized version of convolutional networks* applied to image classification data. With a small kernel size, good features of an image typically appear only at a few patches,¹³ and most other patches are simply random noise or low-magnitude feature noises that are less relevant to the label.

¹¹We say a distribution p over a real interval $[a, b]$ for constants a, b is *left-close*, if there is a $\varepsilon \leq \frac{1}{\text{polylog}(k)}$ such that $\Pr_{z \sim p}[z \leq a + \varepsilon] \geq \frac{1}{\text{polylog}(k)}$, and is *right-close* if $\Pr_{z \sim p}[z \geq b - \varepsilon] \geq \frac{1}{\text{polylog}(k)}$. For instance, here $Z = \sum_{p \in \mathcal{P}_v(X)} z_p$ can be a uniform distribution over $[1, 2]$. This assumption is simply to avoid the case when the distribution is too skewed.

¹²For instance, $Z = \sum_{p \in \mathcal{P}_v(X)} z_p$ can be a uniform distribution over $[0.2, 0.4]$.

¹³For example, in image classification when the image is of size 64×64 , at the first layer, each patch can be a sub-image of size $d = 48 = 4 \times 4 \times 3$ (3 RGB channels), and there are 256 patches. At the second layer, we can have higher dimension d per patch, such as $d = 4 \times 4 \times 64$ when more channels are introduced. In convolutional networks, there are typically over-laps between patches, we point out that our setting is more general: In fact, for example for a data $X = (a, b, c, d)$ with patches (a, b, c) and (b, c, d) , we can simply define $x_1 = (a, b, c)$ and $x_2 = (b, c, d)$. Moreover, our X can also be viewed as intermediate output of the previous convolution layer in a convolution network.

More importantly, the above concept class (namely, labeled data distribution in Def. 3.1) is *not learnable by linear classifiers or constant degree polynomials*. Indeed, if we only use a linear classifier, then the total accumulated (low-magnitude) feature noise from all patches can be as large as $\gamma P \gg 1$ by our choice of $\gamma = \frac{1}{k^{1.5}}$ and $P = k^2$. This is much larger than the magnitude of the signal. On the other hand, by Markov brother’s inequality, low-degree polynomials also lack the power to be approximately linear (to fit the signal) when the input is large, while being sub-linear when there are low-magnitude feature noises. We also conjecture that one can prove this concept class is not efficiently learnable by kernel methods in general, using the recent development of kernel lower bounds [2, 4]. Thus, we believe a (convolutional) neural network with ReLU-like activation is in some sense necessary to learn this concept class.

Our final data distribution \mathcal{D} , and the training data set \mathcal{Z} are formally given as follows.

Definition 3.3 (\mathcal{D} and \mathcal{Z}). *We assume that the final distribution \mathcal{D} consists of data from \mathcal{D}_m w.p. $1 - \mu$ and from \mathcal{D}_s w.p. μ . We are given N training samples from \mathcal{D} , and denote the training data set as $\mathcal{Z} = \mathcal{Z}_m \cup \mathcal{Z}_s$ where \mathcal{Z}_m and \mathcal{Z}_s respectively represent multi-view and single-view training data. We write $(X, y) \sim \mathcal{Z}$ as (X, y) sampled uniformly at random from the empirical data set, and denote $N_s = |\mathcal{Z}_s|$. We again for simplicity focus on the setting when $\mu = \frac{1}{\text{poly}(k)}$ and we are given samples $N = k^{1.2}/\mu$ so each label i appears at least $\tilde{\Omega}(1)$ in \mathcal{Z}_s . Although our result trivially applies to other choices of N .*

B Experiment Details

Our real-life experiments use the CIFAR-10/100 datasets [49]. The SimpleCNN architecture we have used comes from [8], and the (pre-activation) ResNet architecture we have used comes from the wide resnet work [88]. For instance, SimpleCNN-10-3 stands for the 10-layer architecture in [8] but widened by a factor of 3; and ResNet-34-2 stands for the 34-layer wide resnet architecture in [88] and the widening factor is 2.

For training regular neural networks, it is well-known that SGD with momentum and 0.1 learning rate is a state-of-the-art training method. We use batch size 125, train for 140 epochs, and decay the learning rate thrice at epochs 80, 100 and 120 each by a factor 0.2.¹⁴ We use standard random crop, random flip, normalization, and cutout augmentation [77] for the training data.

For training neural-kernel models (NTKs), we find Adam a better training algorithm with an initial learning rate 0.001. We use batch size 50, train for 200 epochs, and decay the learning rate twice at epochs 140 and 170 each by a factor 0.2. We use ZCA data preprocessing which has been reported very helpful for improving neural kernel methods’ performance together with cutout augmentation [77].¹⁵

B.1 Real-Life Data: Single Model vs Ensemble vs Distillation

For our experiment in **Figure 6**, we compare the performance of neural kernel methods vs. real neural networks on the standard CIFAR-10/100 datasets.

When presenting the single-model accuracies in Figure 6, we simply run the training algorithms 10 times from independently randomly initialized seeds. The NTK models we present the best accuracies among the 10 runs, and for ResNet models we present the mean and standard deviations.

¹⁴Some standard training batch for such parameter settings can be found on <https://github.com/bearpaw/pytorch-classification>.

¹⁵ZCA data preprocessing does not help regular neural net training.

When presenting the ensemble accuracies in Figure 6, we simply take an average of the 10 independently training models’ outputs and use that to predict test labels.

When presenting the “directly train $\sum_{\ell} f_{\ell}$ ” result in Figure 6, we directly train a larger network consists of averaging 10 single models (separately, independently initialized). We use the same training algorithm as that for training single models. For some of the NTK models, our 16GB GPU memory sometimes only allows us to train an average of *fewer than 10 single models*; and when we do so, we have put a \diamond remark in Figure 6.¹⁶

When presenting the “knowledge distillation” result in Figure 6, we adopt the original knowledge distillation objective of [42]. It is very similar to (4.3) that we used in our theoretical proof (see (4.3)). It has a weight parameter for the ratio between standard cross-entropy vs the distillation objective (known as $\frac{\eta'}{\eta}$ in (4.3)), and they have a temperature parameter that controls the distillation objective (that is very similar to our τ parameter in (4.3)). We have tuned both parameters in a reasonable range to get the best distillation accuracy.

When presenting the “self-distillation” result in Figure 6, we divide the training process of a single model into two stages: in the first stage it uses the original cross-entropy loss with hard training labels and records the best model in the checkpoints, and in the second stage it trains another single model from random initialization using the distillation objective of [42] to match the output of the previously recorded best model.

Remark B.1. We confirm two more experimental findings that we did not include in Figure 6. First, one can repeat self-distillation multiple times but the test accuracy gain becomes very incremental. Second, one can alternatively use a three-stage process for self-distillation like we did in our theoretical result (see Section 4.2): namely, train two independent single models F and G , and then continue to train G by distilling it to match the output of F . The resulting test accuracy is extremely close to that of the two-stage process.

B.2 Real-life Data: Ensemble over Distillations of Ensemble

For our experiment in **Figure 8**, we have studied the process of (1) training 10 independent single models, (2) evaluating their ensemble, (3) training 10 independent single models using knowledge distillations to match the outputs of (2), and (4) evaluating their ensemble.

The process of (1) and (2) are identical to that in Section B.1.

To present (3), we first apply parameter tuning for the knowledge distillation objective (see Section B.1). Then, we fix the best-selected parameters and perform knowledge distillation 10 times. In other words, these 10 runs differ *only* in the random seeds used in their initialization and SGD, but are *identical* in learning rate, weight decay, knowledge distillation parameters, and all other parameters.

Finally, (4) is a simple (unweighted) average over the 10 models produced by (3).

B.3 Real-life Data: Justifying the Multi-View Assumption

We also perform an experiment in **Figure 9** to justify that in real-life training, there is strong evidence that there are multiple views of the data— even at some intermediate layers— to justify the image labels.

Recall that ResNet has three blocks of layers. In the (a) version of the experiment, we take a pre-trained model, and view its output at the end of the *first* block as “input”, to train a new model where the trainable parameters are the second and third blocks. In the (b) version of the

¹⁶In the case of ResNet16-5-NTK and ResNet10-10-NTK, it only allows us to train an average of 2 models; since this is somewhat meaningless, we simply report “out of memory” in Figure 6.

experiment, we view the pre-trained model’s output at the end of the *second* block as “input”, to train a new model where the trainable parameters are in the third block only.

Specifically, we consider ResNet-28- M version (a) and (b) for $M \in \{1, 2, 4, 10\}$. For instance, the new “input” has $N = 32$ channels for the case of “ResNet-28-2 version (a)”, and has $N = 320$ channels for the case of “ResNet-28-10 version (b).”

For each of the settings above, we

- split the input into 8 chunks (with $N/8$ channels) and train 1 model each, totaling 8 models;
- split the input into 4 chunks (with $N/4$ channels) and train 2 models each, totaling 8 models;
- split the input into 2 chunks (with $N/2$ channels) and train 4 models each, totaling 8 models;
- average the input into $N/8$ channels (by averaging over every 8 channels) and train 8 models;
- average the input into $N/4$ channels (by averaging over every 4 channels) and train 8 models;
- average the input into $N/2$ channels (by averaging over every 2 channels) and train 8 models.

We call those 8 models “single models” and present their accuracies in the first half of the rows of Figure 9.

Next, we also present the ensemble accuracy of these 8 single models in the second half of the rows of Figure 9. (Note for the 8 single models, also use 8 different seeds for the upper-layer pre-trained models. This allows us to compare the ensemble accuracies in a more fair manner.)

B.4 Synthetic Data: Whether Ensemble Improves Accuracy over Gaussian-Like Data

Recall in **Figure 7** we have shown that ensemble does not seem to improve test accuracy on Gaussian-like data. We explain how we perform this experiment.

Synthetic data generation. We generate synthetic data with $k = 10$ labels.

- We consider inputs that are generated as either Gaussian or mixture of Gaussian with different means.
- We consider inputs that are either uniformly generated, or generated through rejection sampling (so as to make different labels to have roughly the same number of data).
- We consider data that are either without label noise, or with 10% of the label randomly flipped.
- We consider data that are generated from a relatively small (but unknown to the learner) ground-truth network, that are either linear, or fully-connected (e.g. fc2 for 2 layers), or convolutional (e.g. conv3 for 3 layers), or residual (e.g. res3 for 3-layered residual and resconv3 for 3-layered residual convolutional).
- We consider data that are either generated as above, or generated with margin across labels.
- Finally, for each of the settings above, we select a dimension d so that the single-model testing accuracy is around 60% \sim 80%.

Learner networks. We also consider fully-connected, convolutional, as well as residual networks with $m = 200$ neurons to learn the given data distribution. For each data/learner pair, we use SGD with momentum 0.9, and tune the learning rate together with weight decay parameters so as to maximize test accuracy. We run for 10 single models and compare their (best) accuracy to their ensemble accuracy.

Result: single vs ensemble. Our detailed comparison tables are in Figure 10 (for non-convolutional inputs) and Figure 11 (for convolutional inputs). To make the result more easily interpretable, we

have included in Figure 7 an abbreviated table which, for each data distribution, picks the best single and best ensemble model across all learner networks. It is clear from these reported results that, for a plethora of settings of Gaussian-like datasets, the accuracy given by ensemble barely exceeds that of single models.

Result: accuracy consistency on single models. For our synthetic datasets, we have also computed the mean and standard deviation for the 10 trained single models from different random initializations. We observe that their standard deviation is also negligible comparing the already not-so-great accuracy: for instance, a standard deviation of 1.0% is quite small comparing to a 70% test accuracy model on the test data. See Figure 12.

		fc2 learner	fc3 learner	fc4 learner	res3 learner	res4 learner	input dimension						
without label noise	uniform sampling	gauss, linear	79.4% (79.4%)	0.0%	79.9% (79.2%)	-0.7%	78.0% (77.8%)	-0.3%	80.3% (79.2%)	-1.1%	80.1% (79.6%)	-0.5%	100
		gauss, fc2	65.4% (64.9%)	-0.4%	65.0% (64.2%)	-0.8%	64.2% (64.4%)	0.2%	67.7% (65.1%)	-2.6%	64.7% (64.2%)	-0.5%	100
		gauss, fc3	66.6% (66.7%)	0.1%	68.9% (68.2%)	-0.7%	68.6% (69.0%)	0.5%	67.7% (67.5%)	-0.2%	68.4% (68.1%)	-0.3%	52
		gauss, res3	67.1% (67.0%)	-0.1%	68.8% (68.0%)	-0.8%	66.7% (66.3%)	-0.5%	69.1% (67.1%)	-2.0%	67.5% (67.8%)	0.3%	100
		gauss, linear, margin	79.0% (77.9%)	-1.2%	78.8% (77.8%)	-1.0%	77.8% (77.3%)	-0.5%	77.8% (77.3%)	-0.5%	78.7% (78.0%)	-0.8%	292
		gauss, fc2, margin	80.6% (78.7%)	-1.9%	78.9% (79.0%)	0.2%	78.6% (78.5%)	-0.1%	78.8% (78.7%)	-0.1%	79.0% (78.8%)	-0.3%	100
		gauss, fc3, margin	75.2% (75.8%)	0.5%	76.0% (75.2%)	-0.8%	75.9% (75.6%)	-0.4%	76.0% (75.4%)	-0.7%	75.8% (75.4%)	-0.3%	100
		gauss, res3, margin	80.5% (80.8%)	0.3%	80.4% (80.0%)	-0.3%	80.4% (79.2%)	-1.2%	80.7% (80.8%)	0.1%	80.1% (79.5%)	-0.6%	100
		mixture, linear	80.7% (80.0%)	-0.7%	80.5% (79.7%)	-0.8%	78.8% (78.4%)	-0.4%	80.7% (80.1%)	-0.6%	79.8% (79.1%)	-0.8%	100
		mixture, fc2	65.1% (64.5%)	-0.5%	65.2% (64.9%)	-0.4%	64.2% (64.5%)	0.3%	67.7% (64.7%)	-3.0%	64.1% (63.5%)	-0.6%	100
		mixture, fc3	62.8% (62.8%)	0.0%	64.0% (63.7%)	-0.3%	63.4% (64.2%)	0.8%	63.3% (62.7%)	-0.6%	63.7% (64.4%)	0.7%	80
		mixture, res3	69.9% (70.4%)	0.5%	70.7% (71.2%)	0.5%	70.2% (70.4%)	0.1%	70.2% (71.2%)	1.0%	70.6% (70.7%)	0.1%	100
	mixture, linear, margin	77.5% (77.2%)	-0.3%	77.8% (77.0%)	-0.8%	78.0% (76.4%)	-1.6%	77.8% (76.9%)	-0.9%	79.0% (76.6%)	-2.4%	292	
	mixture, fc2, margin	79.9% (78.0%)	-1.9%	78.8% (78.4%)	-0.5%	76.5% (76.3%)	-0.2%	80.4% (77.4%)	-3.0%	77.4% (77.6%)	0.2%	100	
	mixture, fc3, margin	79.1% (78.6%)	-0.5%	79.9% (80.1%)	0.2%	79.5% (79.4%)	-0.1%	79.6% (79.6%)	0.0%	80.4% (79.4%)	-1.0%	100	
	mixture, res3, margin	82.8% (82.8%)	0.0%	84.7% (83.9%)	-0.8%	83.7% (83.7%)	0.0%	84.1% (83.5%)	-0.6%	83.1% (82.9%)	-0.3%	100	
	rejection sampling	gauss, linear	78.9% (78.6%)	-0.4%	78.7% (78.5%)	-0.2%	76.7% (76.7%)	0.0%	78.9% (78.3%)	-0.7%	78.6% (78.4%)	-0.2%	100
		gauss, fc2	66.3% (64.1%)	-2.2%	66.0% (64.5%)	-1.5%	64.2% (63.8%)	-0.4%	66.0% (64.3%)	-1.8%	63.9% (63.6%)	-0.3%	100
		gauss, fc3	69.0% (68.6%)	-0.4%	75.5% (75.2%)	-0.3%	76.8% (76.6%)	-0.2%	72.6% (72.7%)	0.1%	74.8% (74.7%)	-0.1%	24
		gauss, res3	69.3% (68.3%)	-1.0%	69.3% (69.0%)	-0.3%	67.7% (67.9%)	0.2%	68.8% (68.2%)	-0.6%	68.3% (68.8%)	0.5%	100
		gauss, linear, margin	77.7% (77.3%)	-0.4%	77.6% (77.0%)	-0.6%	76.9% (76.4%)	-0.5%	78.4% (77.2%)	-1.2%	77.5% (77.3%)	-0.2%	292
		gauss, fc2, margin	78.4% (76.6%)	-1.8%	77.0% (76.5%)	-0.5%	76.5% (76.5%)	0.0%	76.8% (76.4%)	-0.4%	76.5% (76.2%)	-0.3%	100
		gauss, fc3, margin	72.8% (72.2%)	-0.6%	76.6% (76.2%)	-0.4%	76.9% (77.0%)	0.2%	75.1% (75.7%)	0.7%	75.1% (75.9%)	0.8%	100
		gauss, res3, margin	83.2% (82.7%)	-0.5%	84.6% (84.0%)	-0.6%	83.2% (82.4%)	-0.8%	82.9% (83.0%)	0.1%	83.8% (83.2%)	-0.6%	100
mixture, linear		80.5% (80.0%)	-0.5%	80.4% (80.0%)	-0.4%	80.2% (79.9%)	-0.3%	80.7% (80.7%)	0.0%	79.4% (79.4%)	0.1%	100	
mixture, fc2		67.6% (66.6%)	-1.0%	67.0% (66.8%)	-0.2%	66.7% (66.4%)	-0.3%	66.3% (66.0%)	-0.3%	66.3% (66.9%)	0.6%	100	
mixture, fc3		65.6% (65.4%)	-0.2%	73.1% (73.1%)	0.0%	73.2% (72.7%)	-0.5%	68.6% (68.8%)	0.2%	70.6% (70.5%)	-0.1%	32	
mixture, res3		68.3% (67.8%)	-0.5%	69.4% (69.4%)	0.0%	68.4% (68.6%)	0.2%	69.9% (68.4%)	-1.5%	69.0% (68.5%)	-0.5%	100	
mixture, linear, margin	79.9% (79.3%)	-0.6%	79.7% (79.5%)	-0.2%	79.0% (79.2%)	0.2%	79.9% (80.0%)	0.1%	80.0% (79.7%)	-0.3%	292		
mixture, fc2, margin	75.9% (75.7%)	-0.2%	76.9% (76.9%)	0.0%	76.4% (76.3%)	-0.1%	78.4% (77.0%)	-1.5%	76.2% (76.2%)	0.0%	100		
mixture, fc3, margin	71.9% (71.7%)	-0.3%	73.3% (72.8%)	-0.5%	72.8% (72.9%)	0.1%	71.9% (72.2%)	0.3%	72.1% (71.8%)	-0.3%	100		
mixture, res3, margin	82.6% (81.7%)	-0.9%	84.4% (83.7%)	-0.7%	82.4% (82.4%)	0.1%	82.4% (82.1%)	-0.3%	83.4% (82.9%)	-0.5%	100		
With 10% label noise	uniform sampling	gauss, linear	74.3% (74.1%)	-0.2%	72.2% (71.1%)	-1.1%	68.0% (68.7%)	0.7%	70.8% (70.4%)	-0.4%	68.9% (69.3%)	0.4%	100
		gauss, fc2	63.9% (63.2%)	-0.8%	64.3% (62.0%)	-2.3%	60.3% (61.1%)	0.8%	63.1% (62.5%)	-0.6%	61.3% (61.6%)	0.3%	100
		gauss, fc3	66.1% (65.6%)	-0.5%	66.5% (66.4%)	-0.1%	64.7% (66.0%)	1.3%	66.4% (66.3%)	-0.2%	66.0% (65.3%)	-0.7%	40
		gauss, res3	66.7% (65.9%)	-0.8%	66.4% (65.9%)	-0.4%	64.1% (65.2%)	1.1%	68.7% (65.4%)	-3.3%	63.9% (63.6%)	-0.3%	80
		gauss, linear, margin	82.1% (81.7%)	-0.4%	80.9% (80.5%)	-0.4%	78.4% (78.9%)	0.5%	80.0% (80.3%)	0.3%	79.8% (80.1%)	0.3%	144
		gauss, fc2, margin	75.9% (75.2%)	-0.8%	76.0% (75.1%)	-0.9%	73.4% (74.6%)	1.1%	77.4% (75.3%)	-2.1%	73.7% (74.8%)	1.1%	100
		gauss, fc3, margin	70.7% (70.8%)	0.1%	70.7% (70.3%)	-0.4%	68.4% (69.0%)	0.6%	69.9% (70.1%)	0.2%	69.3% (69.5%)	0.2%	100
		gauss, res3, margin	74.4% (74.0%)	-0.4%	73.1% (72.7%)	-0.3%	71.0% (71.5%)	0.5%	75.5% (72.2%)	-3.3%	70.3% (71.1%)	0.8%	100
		mixture, linear	73.3% (73.1%)	-0.2%	73.6% (74.0%)	0.4%	69.1% (70.2%)	1.1%	72.9% (73.1%)	0.3%	71.8% (72.6%)	0.8%	100
		mixture, fc2	69.1% (66.7%)	-2.4%	68.6% (66.5%)	-2.1%	65.6% (65.8%)	0.2%	70.1% (66.6%)	-3.5%	64.1% (64.6%)	0.5%	80
		mixture, fc3	63.0% (61.9%)	-1.2%	61.3% (62.4%)	1.1%	60.2% (60.0%)	-0.2%	61.5% (60.9%)	-0.6%	58.8% (58.7%)	0.1%	80
		mixture, res3	66.8% (63.6%)	-3.2%	66.9% (63.5%)	-3.4%	61.0% (62.2%)	1.2%	66.0% (63.8%)	-2.3%	61.6% (62.9%)	1.3%	100
	mixture, linear, margin	80.7% (79.7%)	-1.0%	79.4% (80.0%)	0.7%	78.1% (77.0%)	-1.1%	78.8% (79.0%)	0.2%	78.5% (78.6%)	0.1%	144	
	mixture, fc2, margin	73.7% (73.9%)	0.1%	73.2% (73.2%)	0.0%	72.2% (72.9%)	0.7%	72.6% (73.1%)	0.5%	72.7% (73.1%)	0.3%	100	
	mixture, fc3, margin	73.5% (73.5%)	0.0%	73.9% (74.6%)	0.8%	71.9% (73.1%)	1.3%	71.9% (72.7%)	0.8%	71.8% (72.1%)	0.3%	100	
	mixture, res3, margin	75.9% (76.2%)	0.4%	76.7% (76.4%)	-0.3%	76.9% (76.4%)	-0.5%	75.1% (74.9%)	-0.2%	74.5% (74.9%)	0.3%	100	
	rejection sampling	gauss, linear	72.9% (72.2%)	-0.8%	72.7% (71.5%)	-1.2%	66.6% (67.1%)	0.5%	72.2% (70.5%)	-1.8%	72.0% (67.7%)	-4.3%	100
		gauss, fc2	62.7% (62.2%)	-0.5%	62.8% (62.2%)	-0.6%	60.6% (60.4%)	-0.2%	63.7% (62.2%)	-1.5%	64.6% (63.5%)	-1.1%	80
		gauss, fc3	65.2% (64.5%)	-0.7%	69.9% (69.1%)	-0.8%	72.5% (72.0%)	-0.5%	67.2% (67.3%)	0.1%	69.6% (69.5%)	-0.1%	24
		gauss, res3	67.6% (66.5%)	-1.1%	68.1% (68.1%)	0.0%	66.9% (67.8%)	1.0%	64.3% (64.9%)	0.6%	64.5% (64.5%)	0.0%	80
		gauss, linear, margin	81.6% (80.2%)	-1.4%	81.3% (78.8%)	-2.5%	76.7% (76.7%)	0.0%	78.1% (78.0%)	-0.1%	77.4% (77.4%)	0.1%	144
		gauss, fc2, margin	72.1% (72.2%)	0.1%	72.1% (71.5%)	-0.6%	71.5% (71.6%)	0.1%	74.2% (71.6%)	-2.6%	74.7% (71.3%)	-3.5%	100
		gauss, fc3, margin	66.6% (66.9%)	0.3%	70.4% (70.5%)	0.1%	64.4% (64.5%)	0.1%	65.4% (66.1%)	0.6%	63.8% (64.5%)	0.7%	100
		gauss, res3, margin	76.8% (73.6%)	-3.2%	74.1% (74.3%)	0.3%	74.3% (74.8%)	0.5%	70.6% (71.1%)	0.5%	70.1% (71.1%)	1.0%	100
mixture, linear		74.2% (73.7%)	-0.5%	73.8% (73.5%)	-0.3%	72.5% (72.6%)	0.1%	73.5% (72.4%)	-1.1%	69.6% (70.7%)	1.1%	100	
mixture, fc2		62.3% (62.3%)	0.0%	66.0% (63.3%)	-2.8%	66.2% (62.0%)	-4.2%	62.3% (63.1%)	0.8%	61.6% (61.5%)	-0.1%	80	
mixture, fc3		71.9% (71.3%)	-0.5%	75.3% (74.8%)	-0.5%	78.1% (78.6%)	0.5%	73.4% (73.5%)	0.1%	74.4% (73.9%)	-0.5%	20	
mixture, res3		66.0% (65.9%)	-0.1%	68.5% (69.5%)	1.1%	68.8% (69.4%)	0.6%	64.6% (65.9%)	1.3%	63.5% (64.9%)	1.4%	80	
mixture, linear, margin	84.1% (81.5%)	-2.6%	82.0% (82.3%)	0.3%	82.7% (82.4%)	-0.3%	80.8% (80.7%)	-0.1%	80.5% (80.6%)	0.1%	144		
mixture, fc2, margin	75.7% (71.4%)	-4.3%	71.6% (70.9%)	-0.7%	71.0% (71.4%)	0.4%	69.9% (69.7%)	-0.2%	69.0% (70.1%)	1.1%	100		
mixture, fc3, margin	65.2% (65.5%)	0.3%	66.8% (66.4%)	-0.5%	66.1% (66.3%)	0.2%	67.5% (64.0%)	-3.5%	62.6% (62.8%)	0.2%	100		
mixture, res3, margin	75.6% (72.4%)	-3.2%	74.1% (74.5%)	0.4%	72.4% (73.2%)	0.8%	76.4% (72.6%)	-3.8%	75.0% (71.9%)	-3.1%	100		

Figure 10: For synthetic Gaussian-like data, ensemble barely helps on improving test accuracy. In this table, we give a closer look at how ensemble performs with respect to each individual learner network. (For a more concise view, see Figure 7; for convolutional data, see Figure 11.)

	fc2 learner	fc3 learner	fc4 learner	res3 learner	res4 learner	conv2 learner	conv3 learner	resconv3 learner	input dimension
uniform sampling	gauss, conv2	62.4% (61.6%)	61.6% (61.4%)	60.3% (59.6%)	61.4% (60.7%)	65.3% (65.0%)	65.3% (65.0%)	65.0% (65.7%)	65.4% (65.1%)
	gauss, conv3	62.8% (62.4%)	61.5% (61.3%)	61.5% (61.3%)	61.7% (61.4%)	65.3% (65.0%)	65.3% (65.0%)	65.7% (66.5%)	68.5% (68.5%)
	gauss, resconv3	63.2% (63.2%)	63.1% (62.2%)	62.2% (62.7%)	63.3% (63.2%)	66.2% (66.2%)	66.2% (66.2%)	66.2% (66.2%)	76.6% (76.6%)
	gauss, conv2, margin	66.9% (67.1%)	66.7% (66.0%)	66.0% (66.3%)	66.5% (66.8%)	66.6% (66.8%)	66.6% (66.8%)	66.6% (66.8%)	70.6% (70.3%)
	gauss, conv3, margin	67.3% (67.2%)	67.9% (67.0%)	66.4% (66.3%)	68.0% (67.8%)	68.0% (67.8%)	68.0% (67.8%)	68.0% (67.8%)	75.1% (74.9%)
	gauss, resconv3, margin	76.4% (76.5%)	76.6% (76.1%)	77.1% (76.3%)	76.6% (75.7%)	76.4% (75.4%)	76.4% (75.4%)	76.4% (75.4%)	87.8% (87.8%)
	mixture, conv2	63.0% (62.5%)	62.7% (62.4%)	62.9% (62.1%)	63.8% (63.2%)	63.8% (63.2%)	63.8% (63.2%)	63.8% (63.2%)	67.0% (66.8%)
	mixture, conv3	65.5% (65.1%)	65.5% (65.4%)	64.5% (64.9%)	64.6% (64.3%)	64.2% (63.8%)	64.2% (63.8%)	64.2% (63.8%)	70.7% (70.0%)
	mixture, resconv3	70.4% (70.1%)	70.8% (70.0%)	69.5% (69.1%)	70.0% (70.4%)	70.0% (70.4%)	70.0% (70.4%)	70.0% (70.4%)	79.3% (79.3%)
	mixture, conv2, margin	64.9% (64.5%)	63.9% (64.1%)	63.8% (63.6%)	64.4% (64.0%)	63.8% (63.6%)	64.4% (64.0%)	63.8% (63.6%)	73.3% (73.3%)
without label noise	gauss, conv2	65.3% (65.0%)	65.0% (65.5%)	65.3% (65.1%)	65.1% (65.4%)	65.1% (65.0%)	65.1% (65.0%)	65.1% (65.0%)	68.3% (68.2%)
	gauss, conv3	66.1% (66.0%)	66.4% (66.6%)	66.0% (66.4%)	66.6% (66.2%)	66.6% (66.8%)	66.6% (66.8%)	66.6% (66.8%)	74.9% (74.9%)
	gauss, resconv3	68.2% (67.7%)	67.5% (68.4%)	68.2% (68.5%)	68.6% (68.1%)	68.6% (68.1%)	68.6% (68.1%)	68.6% (68.1%)	82.6% (82.6%)
	gauss, conv2, margin	62.7% (62.5%)	63.0% (62.9%)	62.6% (63.3%)	62.8% (62.1%)	61.7% (61.9%)	61.7% (61.9%)	61.7% (61.9%)	67.8% (67.8%)
	gauss, conv3, margin	63.6% (63.1%)	63.5% (63.3%)	63.3% (63.6%)	63.3% (62.9%)	62.6% (62.6%)	62.6% (62.6%)	62.6% (62.6%)	76.8% (76.4%)
	gauss, resconv3, margin	72.1% (72.1%)	72.0% (71.4%)	71.1% (71.9%)	72.2% (71.7%)	71.0% (72.2%)	71.0% (72.2%)	71.0% (72.2%)	92.6% (92.3%)
	mixture, conv2	64.5% (64.5%)	65.1% (64.7%)	64.1% (65.2%)	64.7% (64.3%)	64.6% (65.0%)	64.6% (65.0%)	64.6% (65.0%)	68.2% (68.2%)
	mixture, conv3	67.5% (67.7%)	68.1% (67.8%)	68.3% (68.2%)	68.0% (68.0%)	68.4% (68.8%)	68.0% (68.0%)	68.0% (68.0%)	75.5% (75.5%)
	mixture, resconv3	66.6% (66.4%)	67.6% (66.8%)	66.3% (66.0%)	66.5% (66.1%)	66.2% (67.0%)	66.2% (67.0%)	66.2% (67.0%)	81.5% (81.3%)
	mixture, conv2, margin	64.2% (64.3%)	63.6% (64.0%)	63.2% (62.7%)	64.2% (64.1%)	63.7% (64.0%)	63.7% (64.0%)	63.7% (64.0%)	69.1% (68.4%)
uniform sampling	gauss, conv2	73.0% (72.5%)	73.5% (72.7%)	73.0% (73.1%)	73.2% (73.0%)	73.4% (73.1%)	73.2% (73.0%)	73.4% (73.1%)	84.3% (83.6%)
	gauss, conv3	62.2% (61.6%)	61.8% (61.5%)	60.8% (60.7%)	62.1% (61.7%)	61.5% (61.2%)	61.5% (61.2%)	61.5% (61.2%)	66.2% (66.1%)
	gauss, resconv3	65.1% (65.2%)	64.0% (63.7%)	63.0% (63.5%)	64.6% (64.4%)	63.0% (63.5%)	64.6% (64.4%)	63.0% (63.5%)	74.1% (72.5%)
	gauss, conv2, margin	65.7% (65.4%)	64.2% (63.9%)	62.2% (63.0%)	64.4% (64.6%)	62.8% (64.0%)	64.4% (64.6%)	62.8% (64.0%)	70.8% (71.3%)
	gauss, conv3, margin	73.7% (73.3%)	71.8% (71.7%)	70.2% (70.0%)	71.5% (71.2%)	70.6% (70.8%)	71.5% (71.2%)	70.6% (70.8%)	84.0% (82.6%)
	mixture, conv2	63.5% (63.2%)	63.3% (62.9%)	61.2% (61.7%)	63.1% (62.9%)	61.2% (62.9%)	63.1% (62.9%)	61.2% (62.9%)	65.1% (65.5%)
	mixture, conv3	63.1% (63.4%)	62.4% (62.4%)	59.7% (60.5%)	62.1% (62.5%)	61.2% (62.5%)	62.1% (62.5%)	61.2% (62.5%)	67.7% (67.5%)
	mixture, resconv3	63.0% (62.9%)	61.7% (61.9%)	58.8% (60.9%)	61.6% (62.2%)	60.5% (62.4%)	61.6% (62.2%)	60.5% (62.4%)	72.8% (73.1%)
	mixture, conv2, margin	68.2% (67.9%)	68.0% (68.1%)	66.8% (66.0%)	68.1% (68.4%)	67.5% (68.0%)	68.1% (68.4%)	67.5% (68.0%)	80.9% (80.9%)
	mixture, conv3, margin	75.4% (75.4%)	74.7% (75.4%)	72.4% (73.6%)	75.3% (75.2%)	73.6% (73.6%)	75.3% (75.2%)	73.6% (73.6%)	81.2% (81.3%)
without label noise	gauss, conv2	62.9% (62.4%)	62.4% (62.5%)	61.6% (60.9%)	62.7% (62.4%)	62.6% (62.7%)	62.9% (62.5%)	62.6% (62.7%)	65.0% (64.8%)
	gauss, conv3	64.3% (63.7%)	64.4% (64.3%)	63.8% (64.1%)	64.5% (63.6%)	64.4% (64.9%)	64.5% (63.6%)	64.4% (64.9%)	69.9% (69.9%)
	gauss, resconv3	61.3% (61.1%)	62.7% (62.1%)	60.3% (60.9%)	62.1% (62.2%)	61.3% (62.1%)	62.1% (62.2%)	61.3% (62.1%)	76.6% (75.9%)
	gauss, conv2, margin	65.7% (64.6%)	65.5% (65.5%)	63.3% (64.1%)	65.0% (65.0%)	65.6% (65.8%)	65.0% (65.8%)	65.6% (65.8%)	66.7% (66.6%)
	gauss, conv3, margin	67.6% (67.0%)	66.9% (67.3%)	65.7% (65.8%)	66.6% (66.2%)	66.3% (66.2%)	66.6% (66.2%)	66.3% (66.2%)	72.7% (72.3%)
	gauss, resconv3, margin	64.9% (64.6%)	64.1% (64.1%)	62.8% (62.8%)	64.2% (64.7%)	63.9% (64.2%)	64.2% (64.7%)	63.9% (64.2%)	82.9% (85.3%)
	mixture, conv2	64.5% (64.2%)	64.5% (64.8%)	63.4% (64.4%)	64.9% (65.0%)	64.8% (64.8%)	64.9% (65.0%)	64.8% (64.8%)	66.3% (66.3%)
	mixture, conv3	66.6% (67.6%)	65.8% (66.4%)	63.8% (63.8%)	66.0% (65.7%)	65.4% (65.1%)	66.0% (65.7%)	65.4% (65.1%)	77.5% (77.9%)
	mixture, resconv3	64.7% (64.4%)	63.4% (64.0%)	62.1% (62.1%)	63.6% (63.6%)	62.7% (63.4%)	63.6% (63.6%)	62.7% (63.4%)	67.6% (67.2%)
	mixture, conv2, margin	65.5% (65.7%)	64.9% (64.8%)	63.1% (65.0%)	65.3% (65.4%)	65.0% (65.2%)	65.3% (65.4%)	65.0% (65.2%)	73.3% (73.3%)
mixture, conv3, margin	67.3% (67.3%)	66.3% (67.1%)	66.8% (67.5%)	66.1% (67.1%)	65.0% (66.3%)	66.1% (67.1%)	65.0% (66.3%)	84.2% (84.7%)	

Figure 11: For synthetic Gaussian-like data, ensemble barely helps on improving test accuracy. In this table, we give a closer look at how ensemble performs with respect to each individual learner network, when the data is generated from a *target convolutional network*. This gives evidence that, having convolutional data may alone be necessary for ensemble to work either.

		learner networks								
		synthetic data	fc2	fc3	fc4	res3	res4	conv2	conv3	resconv3
without label noise	uniform sampling	gauss, linear	79.0±0.2%	79.4±0.3%	77.5±0.4%	79.5±0.3%	79.4±0.3%	74.4±0.6%	73.7±0.7%	74.3±0.5%
		gauss, fc2	65.0±0.2%	64.5±0.3%	63.8±0.3%	64.7±1.2%	64.3±0.3%	62.1±0.4%	62.1±0.5%	62.5±0.5%
		gauss, fc3	66.4±0.1%	68.4±0.3%	68.0±0.3%	67.3±0.3%	67.3±0.5%	59.4±0.3%	62.6±0.8%	61.1±0.6%
		gauss, res3	66.4±0.3%	67.3±0.6%	66.1±0.6%	67.3±0.7%	66.8±0.5%	60.9±0.4%	60.3±0.7%	60.9±0.5%
		gauss, linear, margin	78.2±0.3%	78.0±0.3%	77.2±0.4%	77.0±0.3%	77.8±0.4%	70.3±1.0%	69.6±0.8%	69.7±0.9%
		gauss, fc2, margin	78.5±0.8%	78.4±0.3%	78.2±0.4%	78.3±0.4%	78.3±0.3%	75.0±0.3%	74.6±0.4%	74.0±0.3%
		gauss, fc3, margin	75.0±0.1%	75.2±0.4%	75.2±0.4%	75.2±0.4%	75.3±0.4%	67.9±0.4%	68.2±0.7%	66.6±0.8%
		gauss, res3, margin	80.1±0.3%	80.0±0.3%	77.7±3.0%	80.2±0.3%	79.6±0.3%	71.9±0.4%	71.3±0.6%	71.8±0.5%
		mixture, linear	80.0±0.3%	79.8±0.3%	78.1±0.5%	80.2±0.3%	79.3±0.3%	75.3±0.5%	72.8±0.5%	74.8±0.5%
		mixture, fc2	64.2±0.4%	64.2±0.5%	63.0±0.5%	64.4±1.2%	63.5±0.4%	60.9±0.3%	60.6±0.5%	59.9±0.4%
		mixture, fc3	62.3±0.3%	63.4±0.4%	63.1±0.3%	62.5±0.5%	63.0±0.4%	58.3±0.3%	60.2±0.6%	60.2±0.6%
		mixture, res3	69.5±0.4%	70.3±0.4%	69.9±0.3%	69.7±0.3%	69.8±0.5%	63.2±0.4%	61.5±0.6%	63.1±0.5%
		mixture, linear, margin	76.8±0.3%	76.8±0.5%	76.4±0.7%	77.0±0.5%	76.9±0.8%	69.7±0.7%	69.8±0.6%	69.7±1.0%
		mixture, fc2, margin	77.2±1.1%	77.3±0.6%	75.9±0.4%	77.1±1.2%	76.4±0.4%	74.0±0.4%	73.2±0.4%	73.9±0.5%
	mixture, fc3, margin	78.4±0.3%	79.6±0.2%	79.1±0.4%	79.1±0.3%	79.1±0.6%	71.8±0.6%	72.0±0.4%	72.1±0.5%	
	mixture, res3, margin	82.5±0.2%	83.7±0.5%	82.5±1.2%	83.2±0.5%	82.7±0.3%	74.7±0.4%	73.9±0.6%	73.4±0.5%	
	rejection sampling	gauss, linear	78.2±0.5%	78.0±0.4%	76.2±0.3%	78.4±0.3%	78.0±0.3%	73.9±0.5%	72.3±0.9%	73.4±0.7%
		gauss, fc2	63.7±1.0%	63.9±0.8%	63.1±0.7%	63.8±0.9%	63.0±0.4%	59.6±0.2%	59.3±0.6%	59.4±0.7%
		gauss, fc3	68.4±0.3%	75.1±0.2%	76.1±0.3%	72.2±0.3%	73.6±0.6%	56.8±0.3%	65.1±0.4%	63.6±0.5%
		gauss, res3	67.9±0.5%	68.9±0.3%	67.0±0.5%	67.9±0.3%	67.5±0.4%	58.3±0.3%	59.7±0.6%	59.6±0.7%
		gauss, linear, margin	77.2±0.3%	76.9±0.3%	76.3±0.4%	77.3±0.4%	77.0±0.3%	70.4±0.9%	69.5±0.5%	68.5±0.8%
		gauss, fc2, margin	76.4±0.8%	76.4±0.4%	76.2±0.3%	76.4±0.3%	76.1±0.3%	72.2±0.4%	72.0±0.7%	72.0±0.5%
		gauss, fc3, margin	72.5±0.2%	76.2±0.3%	75.8±0.5%	74.6±0.3%	74.6±0.4%	58.9±0.8%	63.7±1.0%	62.9±0.5%
		gauss, res3, margin	82.5±0.3%	83.7±0.4%	82.4±0.4%	82.7±0.1%	82.6±0.5%	72.5±0.5%	72.2±0.6%	73.7±0.7%
		mixture, linear	80.2±0.3%	80.0±0.2%	78.4±1.6%	80.1±0.3%	78.8±0.3%	74.5±0.6%	73.9±1.0%	74.7±0.5%
		mixture, fc2	66.2±0.5%	66.3±0.4%	66.0±0.4%	66.0±0.3%	65.6±0.5%	62.6±0.5%	62.0±0.6%	62.6±0.6%
		mixture, fc3	65.1±0.3%	72.2±0.7%	72.8±0.3%	67.7±0.4%	69.8±0.6%	52.1±0.3%	61.8±0.4%	60.0±0.4%
		mixture, res3	67.7±0.3%	68.9±0.2%	67.7±0.4%	67.6±0.9%	67.9±0.5%	59.8±0.4%	59.5±0.6%	59.8±0.5%
mixture, linear, margin		79.3±0.3%	79.2±0.3%	78.4±0.3%	79.5±0.2%	79.5±0.3%	71.2±0.7%	71.0±0.7%	70.7±1.2%	
mixture, fc2, margin		75.6±0.2%	76.2±0.4%	76.0±0.3%	76.0±0.9%	75.7±0.3%	72.8±0.4%	71.9±0.5%	72.4±0.7%	
mixture, fc3, margin	71.4±0.3%	72.8±0.5%	71.5±0.6%	71.5±0.4%	71.4±0.4%	62.0±0.3%	63.4±0.8%	61.9±0.4%		
mixture, res3, margin	81.4±0.5%	83.1±0.5%	81.7±0.3%	81.6±0.4%	82.3±0.5%	71.3±0.6%	71.0±0.7%	71.3±0.4%		
With 10% label noise	uniform sampling	gauss, linear	73.2±0.5%	71.1±0.5%	67.3±0.5%	69.8±0.4%	68.2±0.5%	67.3±0.4%	66.0±0.7%	67.5±0.4%
		gauss, fc2	62.8±0.5%	61.7±1.0%	59.6±0.5%	61.7±0.6%	60.0±0.7%	61.3±0.6%	60.7±0.3%	61.1±0.4%
		gauss, fc3	65.7±0.2%	66.2±0.2%	63.9±0.6%	65.6±0.5%	65.0±0.6%	60.5±0.4%	61.4±0.5%	60.9±0.3%
		gauss, res3	66.4±0.2%	65.6±0.5%	63.4±0.5%	64.8±1.4%	62.6±0.6%	62.3±0.9%	61.7±0.8%	62.5±0.6%
		gauss, linear, margin	81.5±0.3%	80.1±0.4%	77.9±0.3%	79.4±0.4%	78.8±0.5%	75.0±0.4%	73.8±0.6%	74.9±0.6%
		gauss, fc2, margin	75.3±0.3%	74.6±0.5%	73.1±0.3%	74.5±1.1%	73.2±0.4%	73.6±0.2%	72.9±0.6%	73.1±0.6%
		gauss, fc3, margin	70.4±0.2%	70.3±0.3%	67.9±0.5%	69.2±0.3%	68.1±0.7%	66.4±0.3%	66.0±0.5%	66.2±0.5%
		gauss, res3, margin	72.8±0.7%	72.3±0.4%	70.1±0.4%	71.0±1.6%	69.4±0.6%	68.4±0.4%	67.3±0.7%	68.6±0.5%
		mixture, linear	72.3±0.5%	72.6±0.5%	68.2±0.6%	71.2±1.0%	70.3±0.9%	69.4±0.5%	67.4±0.9%	68.0±0.6%
		mixture, fc2	66.2±1.1%	65.2±1.2%	62.2±1.4%	64.4±2.0%	63.5±0.4%	64.2±0.3%	63.9±0.4%	63.7±0.6%
		mixture, fc3	61.1±0.8%	59.9±0.7%	59.0±0.5%	59.3±0.9%	57.9±0.6%	57.9±0.5%	58.1±0.7%	57.5±0.3%
		mixture, res3	63.2±1.3%	62.2±1.9%	60.3±0.4%	61.3±1.8%	60.3±0.7%	59.2±0.3%	60.0±0.8%	59.8±0.9%
		mixture, linear, margin	79.2±0.5%	78.8±0.4%	76.5±1.6%	77.9±0.6%	77.5±0.6%	72.4±0.6%	71.8±0.9%	72.0±0.8%
		mixture, fc2, margin	72.8±0.5%	72.5±0.3%	71.2±0.6%	71.9±0.5%	71.8±0.4%	70.6±0.5%	69.8±0.5%	70.2±0.4%
	mixture, fc3, margin	72.4±0.5%	72.9±0.6%	70.8±0.7%	71.4±0.4%	70.6±0.5%	69.5±0.5%	69.2±0.6%	69.1±0.6%	
	mixture, res3, margin	75.5±0.3%	75.8±0.6%	73.9±1.7%	74.2±0.6%	73.6±0.6%	72.0±0.4%	71.7±0.7%	71.3±0.6%	
	rejection sampling	gauss, linear	71.8±0.5%	70.2±1.0%	66.0±0.4%	68.6±1.4%	65.8±2.3%	65.1±0.7%	64.0±0.9%	65.3±0.7%
		gauss, fc2	62.2±0.3%	62.0±0.5%	59.2±0.6%	60.7±1.1%	60.1±1.7%	59.9±0.5%	59.0±0.4%	59.8±0.4%
		gauss, fc3	64.6±0.4%	69.3±0.4%	71.3±0.7%	66.9±0.2%	68.3±0.9%	54.7±0.5%	62.7±1.1%	60.0±0.7%
		gauss, res3	66.5±0.5%	67.5±0.3%	66.1±0.5%	63.8±0.4%	63.3±0.7%	59.2±0.4%	58.6±0.9%	59.7±0.6%
		gauss, linear, margin	80.3±0.6%	78.3±1.2%	76.2±0.4%	77.5±0.4%	76.9±0.2%	73.4±0.8%	72.2±0.6%	73.4±1.0%
		gauss, fc2, margin	71.8±0.2%	71.6±0.3%	70.1±1.2%	71.0±1.2%	69.9±1.7%	69.4±0.6%	68.9±0.5%	69.5±0.5%
		gauss, fc3, margin	66.2±0.2%	69.5±0.5%	63.5±0.5%	64.6±0.4%	63.1±0.5%	58.3±0.5%	60.5±1.0%	59.3±0.5%
		gauss, res3, margin	73.4±1.2%	73.6±0.3%	73.5±0.5%	70.1±0.3%	69.1±0.6%	67.6±0.3%	68.1±0.6%	68.4±0.8%
		mixture, linear	72.3±0.8%	70.8±1.2%	69.2±3.1%	70.0±1.3%	68.9±0.6%	67.9±0.5%	66.3±0.5%	67.4±0.7%
		mixture, fc2	61.9±0.3%	62.2±1.4%	59.8±2.3%	61.6±0.7%	59.6±0.8%	59.3±0.7%	59.8±0.4%	59.9±0.7%
		mixture, fc3	71.4±0.2%	74.9±0.3%	77.2±0.6%	72.6±0.4%	73.6±0.5%	61.7±0.2%	69.0±0.2%	67.3±0.5%
		mixture, res3	64.9±0.6%	67.6±0.7%	67.0±1.0%	63.9±0.4%	62.2±0.7%	60.4±0.6%	60.6±0.8%	59.8±0.7%
mixture, linear, margin		81.1±1.1%	81.5±0.4%	80.7±2.0%	80.3±0.3%	79.9±0.4%	75.0±1.2%	73.4±1.2%	74.1±1.1%	
mixture, fc2, margin		70.8±1.8%	70.5±0.6%	69.9±0.6%	69.0±0.6%	68.2±0.6%	68.6±0.7%	68.2±0.6%	67.5±0.7%	
mixture, fc3, margin	64.7±0.3%	65.4±0.7%	65.0±0.6%	63.3±1.5%	61.5±0.6%	60.7±0.7%	61.4±0.6%	60.0±0.7%		
mixture, res3, margin	71.6±1.4%	73.3±0.6%	71.2±0.7%	70.6±2.1%	69.3±2.2%	68.1±0.4%	67.8±0.7%	67.6±0.6%		

Figure 12: For synthetic Gaussian-like data, we compute the standard deviations of single models test accuracies over 10 runs. This standard deviation is generally bigger than what we see on the CIFAR10/100 datasets (see Figure 6), yet ensemble still offers nearly no benefit.

APPENDIX II: COMPLETE PROOFS

C Single Model: Proof Plan and Induction Hypothesis

Our main proof relies on an induction hypothesis for every iteration $t = 0, 1, 2, \dots, T$. Before we state it, let us introduce several notations. Let us denote

$$\Lambda_i^{(t)} \stackrel{\text{def}}{=} \max_{r \in [m], \ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \quad \text{and} \quad \Lambda_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \max_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \quad (\text{C.1})$$

Suppose $m \leq \text{poly}(k)$. For every $i \in [k]$, let us denote

$$\mathcal{M}_i^{(0)} \stackrel{\text{def}}{=} \left\{ r \in [m] \mid \exists \ell \in [2]: \langle w_{i,r}^{(0)}, v_{i,\ell} \rangle \geq \Lambda_{i,\ell}^{(0)} \left(1 - O\left(\frac{1}{\log k}\right) \right) \right\}$$

Intuition. If a neuron $r \in [m]$ is not in $\mathcal{M}_i^{(0)}$, it means that for both $\ell \in \{1, 2\}$, the correlation $\langle w_{i,r}^{(0)}, v_{i,\ell} \rangle$ at the random initialization is, by a non-trivial factor, smaller than $\Lambda_{i,\ell}^{(0)}$ — the largest correlation between $\langle w_{i,r'}, v_{i,\ell} \rangle$ among all neurons. In words, this means the magnitude of $v_{i,1}$ and $v_{i,2}$ inside the random initialization $w_{i,r}^{(0)}$ is non-trivially lagging behind, comparing to other neurons. We shall prove that, through the course of the training, those neurons r will lose the lottery and not learn anything useful for the output label $i \in [k]$. (This corresponds to Induction Hypothesis C.3i later.)

Fact C.1. *With probability at least $1 - e^{-\Omega(\log^5 k)}$, we have $|\mathcal{M}_i^{(0)}| \leq m_0 \stackrel{\text{def}}{=} O(\log^5 k)$.*

(The proof of Fact C.1 follows from standard analysis on Gaussian variables, see Proposition H.1.)

Suppose we denote by $S_{i,\ell} \stackrel{\text{def}}{=} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y=i} \sum_{p \in P_{v_{i,\ell}}(X)} z_p^q \right]$. Then, define

$$\mathcal{M} \stackrel{\text{def}}{=} \left\{ (i, \ell^*) \in [k] \times [2] \mid \Lambda_{i,\ell^*}^{(0)} \geq \Lambda_{i,3-\ell^*}^{(0)} \left(\frac{S_{i,3-\ell^*}}{S_{i,\ell^*}} \right)^{\frac{1}{q-2}} \left(1 + \frac{1}{\log^2(m)} \right) \right\} \quad (\text{C.2})$$

Intuition. If $(i, \ell) \in \mathcal{M}$, we shall prove that the feature $v_{i,\ell}$ has a higher chance than $v_{i,3-\ell}$ to be learned by the model. (This is because, after an appropriate scaling factor defined by the training data, $v_{i,\ell}$ correlates more with the network’s random initialization comparing to $v_{i,3-\ell}$.)

Our next proposition states that, for every $i \in [k]$, with decent probability at least one of $(i, 1)$ or $(i, 2)$ shall be in \mathcal{M} . But more importantly, our later Induction Hypothesis C.3e ensures that, during the entire training process, if $(i, 3 - \ell) \in \mathcal{M}$, then $v_{i,\ell}$ *must be missing from the learner network*. They together imply that test accuracy on single-view data cannot exceed 49.99%, as one of the views shall be missing.

On the other hand, our next proposition also ensures that the “weaker” feature among the two, still has some non-negligible chance to be picked up by the random initialization. This is behind the reason that why ensemble works in our later proofs.

Proposition C.2. *Suppose $m \leq \text{poly}(k)$. We have the following properties about \mathcal{M} .*

- For every $i \in [k]$, at most one of $(i, 1)$ or $(i, 2)$ is in \mathcal{M} (obvious).
- For every $i \in [k]$, suppose $S_{i,\ell} \geq S_{i,3-\ell}$, then $\Pr [(i, 3 - \ell) \in \mathcal{M}] \geq m^{-O(1)}$.
- For every $i \in [k]$, $\Pr [(i, 1) \in \mathcal{M} \text{ or } (i, 3) \in \mathcal{M}] \geq 1 - o(1)$.

(Proposition C.2 is a result of the anti-concentration of the maximum of Gaussian, see Appendix H.)

We are now ready to state our induction hypothesis.

Induction Hypothesis C.3. For every $\ell \in [2]$, for every $r \in [m]$, for every $(X, y) \in \mathcal{Z}_m$ and $i \in [k]$, or for every $(X, y) \in \mathcal{Z}_s$ and $i \in [k] \setminus \{y\}$:

(a) For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, we have: $\langle w_{i,r}^{(t)}, x_p \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \pm \tilde{O}(\sigma_0)$.

(b) For every $p \in \mathcal{P}(X) \setminus (\mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X))$, we have: $|\langle w_{i,r}^{(t)}, x_p \rangle| \leq \tilde{O}(\sigma_0)$.

(c) For every $p \in [P] \setminus \mathcal{P}(X)$, we have: $|\langle w_{i,r}^{(t)}, x_p \rangle| \leq \tilde{O}(\sigma_0 \gamma k)$.

In addition, for every $(X, y) \in \mathcal{Z}_s$, every $i \in [k]$, every $r \in [m]$, every $\ell \in [2]$,

(d) For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, we have: $\langle w_{i,r}^{(t)}, x_p \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p + \langle w_{i,r}^{(t)}, \xi_p \rangle \pm \tilde{O}(\sigma_0 \gamma k)$

(e) For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, if $(i, 3 - \ell) \in \mathcal{M}$ we have: $|\langle w_{i,r}^{(t)}, x_p \rangle| \leq \tilde{O}(\sigma_0)$.

(f) For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, if $r \in [m] \setminus \mathcal{M}_i^{(0)}$ we have: $|\langle w_{i,r}^{(t)}, x_p \rangle| \leq \tilde{O}(\sigma_0)$.

Moreover, we have for every $i \in [k]$,

(g) $\Lambda_i^{(t)} \geq \Omega(\sigma_0)$ and $\Lambda_i^{(t)} \leq \tilde{O}(1)$.

(h) for every $r \in [m]$, every $\ell \in [2]$, it holds that $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq -\tilde{O}(\sigma_0)$.

(i) for every $r \in [m] \setminus \mathcal{M}_i^{(0)}$, every $\ell \in [2]$, it holds that $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \leq \tilde{O}(\sigma_0)$.

Intuition. The first three items in Induction Hypothesis C.3 essentially say that, when studying the correlation between $w_{i,r}$ with a multi-view data, or between $w_{i,r}$ with a single-view data (but $y \neq i$), the correlation is about $\langle w_{i,r}, v_{i,1} \rangle$ and $\langle w_{i,r}, v_{i,2} \rangle$ and the remaining terms are sufficiently small. (Of course, this requires a careful proof.) We shall later prove that at least one of $\Lambda_{i,1}^{(t)}$ or $\Lambda_{i,2}^{(t)}$ is large after training. Therefore, using the first three items, we can argue that all multi-view data are classified correctly.

The middle three items in Induction Hypothesis C.3 essentially say that, when studying the correlation between $w_{i,r}$ with a single-view data (X, y) with $y = i$, then the correlation also has a significant noise term $\langle w_{i,r}^{(t)}, \xi_p \rangle$. This term shall become useful for us to argue that single-view data can be all memorized (through for instance memorizing the noise).

The remaining items in Induction Hypothesis C.3 are just some regularization statements.

D Single Model: Technical Proofs

We devote this section to prove that Induction Hypothesis C.3 holds for every iteration $t \leq T$, and in the next Section E, we state how the induction hypothesis easily implies our main theorems for single model and ensemble model.

Parameter D.1. We state the range of parameters for our proofs in this section to hold.

- $\varrho = \frac{1}{\text{polylog}(k)}$ (recall ϱ is the threshold for the smoothed ReLU activation)
- $\Gamma = \frac{1}{\text{polylog}(k)}$ (recall Γ controls off-target feature magnitude in Def. 3.1)
- $q \geq 3$ and $\sigma_0^{q-2} = \frac{1}{k}$ (recall $w_{i,r}^{(0)} \sim \mathcal{N}(0, \sigma_0^2 I)$ gives the initialization magnitude)
- $N_s \leq \tilde{O}(k/\rho)$ and $N_s \leq \frac{k^2}{s} \rho^{q-1}$. (recall N_s is the size of single-view training data)
- $\gamma \leq \tilde{O}(\frac{\sigma_0}{k})$ and $\gamma^q \leq \tilde{\Theta}(\frac{1}{k^{q-1} m P})$ (recall γ controls feature noise in Def. 3.1)
- $\text{polylog}(k) \leq s \leq k^{0.2}$ (recall s controls feature sparsity in Def. 3.1)
- $\rho^{q-1} \geq \frac{1}{k}$ (recall ρ controls on-target feature magnitude of single-view data in Def. 3.1)

- $N \geq N_s \cdot \text{poly}(k)$, $\eta T \geq N \cdot \text{poly}(k)$, and $\sqrt{d} \geq \eta T \cdot \text{poly}(k)$.
- $\text{polylog}(k) \leq m \leq \tilde{O}(\frac{1}{s\sigma_0^q})$.

Example. A reasonable set of parameters is, up to polylogarithmic factors:

$$q = 4, \quad \sigma_0 = \frac{1}{\sqrt{k}}, \quad \rho = \frac{1}{k^{0.2}}, \quad m \leq k, \quad s \leq k^{0.2}, \quad N_s \leq k^{1.2}, \quad P \leq k^2, \quad \gamma \leq \frac{1}{k^{1.5}}.$$

Theorem D.2. Under Parameter D.1, for any $m \in [\tilde{\Omega}(1), \tilde{O}(\frac{1}{s\sigma_0^q})]$ and sufficiently small $\eta \leq \frac{1}{\text{poly}(k)}$, our Induction Hypothesis C.3 holds for all iterations $t = 0, 1, \dots, T$.

D.1 Gradient Calculations and Function Approximation

Gradient calculations. Recall $\mathbf{logit}_i(F, X) \stackrel{\text{def}}{=} \frac{e^{F_i(X)}}{\sum_{j \in [k]} e^{F_j(X)}}$. Recall also

Fact D.3. Given data point $(X, y) \in \mathcal{D}$, for every $i \in [k]$, $r \in [m]$,

$$-\nabla_{w_{i,r}} L(F; X, y) = (1 - \mathbf{logit}_i(F, X)) \sum_{p \in [P]} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) x_p \quad \text{when } i = y \quad (\text{D.1})$$

$$-\nabla_{w_{i,r}} L(F; X, y) = -\mathbf{logit}_i(F, X) \sum_{p \in [P]} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) x_p \quad \text{when } i \neq y \quad (\text{D.2})$$

Now, we also have the following observations:

Claim D.4. If Induction Hypothesis C.3 holds at iteration t , and if $s \leq \tilde{O}(\frac{1}{\sigma_0^q m})$ and $\gamma \leq \tilde{O}(\frac{1}{(\sigma_0 k (mP)^{1/q})})$, then

- for every $(X, y) \in \mathcal{Z}$ and every $i \in [k]$: $\mathbf{logit}_i(F^{(t)}, X) = O\left(\frac{e^{O(\Lambda_i^{(t)})m_0}}{e^{O(\Lambda_i^{(t)})m_0 + k}}\right)$
- for every $(X, y) \in \mathcal{Z}_s$ and $i \in [k] \setminus \{y\}$: $\mathbf{logit}_i(F^{(t)}, X) = O\left(\frac{1}{k}\right) (1 - \mathbf{logit}_y(F^{(t)}, X))$

Proof of Claim D.4. Recall $F_i^{(t)}(X) = \sum_{r \in [m]} \sum_{p \in [P]} \widetilde{\text{ReLU}}(\langle w_{i,r}^{(t)}, x_p \rangle)$. For every $r \notin \mathcal{M}_i^{(0)}$, using Induction Hypothesis C.3i we have

$$\sum_{p \in [P]} \widetilde{\text{ReLU}}(\langle w_{i,r}^{(t)}, x_p \rangle) \leq \tilde{O}(\sigma_0^q) \cdot s + \tilde{O}((\sigma_0 \gamma k)^q) \cdot P \leq \frac{1}{m \text{polylog}(k)}$$

so they sum up to at most $\frac{1}{\text{polylog}(k)}$. For any $r \in \mathcal{M}_i^{(0)}$, we have

$$\begin{aligned} \sum_{p \in [P]} \widetilde{\text{ReLU}}(\langle w_{i,r}^{(t)}, x_p \rangle) &\leq \sum_{\ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \cdot \left(\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \right) + \tilde{o}(\sigma_0) + \tilde{O}(\sigma_0^q) \cdot s + \tilde{O}((\sigma_0 \gamma k)^q) \cdot P \\ &\leq \sum_{\ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \cdot \left(\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \right) + O\left(\frac{1}{m_0}\right) \end{aligned}$$

Recall from Def. 3.1 we have $\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \leq O(1)$; and furthermore when $(X, y) \in \mathcal{Z}_s$ and $i \in [k] \setminus \{y\}$ we have $\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \leq \Gamma = \frac{1}{\text{polylog}(k)}$. In the former case, we have

$$0 \leq F_i^{(t)}(X) \leq m_0 \cdot \Lambda_i^{(t)} \cdot O(1) + O(1)$$

and this proves the first bound; in the latter case we have

$$0 \leq F_i^{(t)}(X) \leq m_0 \cdot \Lambda_i^{(t)} \cdot \Gamma + O(1) \leq O(1) \quad (\text{D.3})$$

and this proves the second bound. \square

Definition D.5. For each data point X , we consider a value $V_{i,r,\ell}(X)$ given as:

$$V_{i,r,\ell}(X) \stackrel{\text{def}}{=} \mathbb{1}_{v_{i,\ell} \in \mathcal{V}(X)} \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) z_p$$

Definition D.6. We define four error terms that shall be used frequently in our proofs.

$$\begin{aligned} \mathcal{E}_1 &:= \widetilde{O}(\sigma_0^{q-1}) \gamma s & \mathcal{E}_{2,i,r}(X) &:= O(\gamma(V_{i,r,1}(X) + V_{i,r,2}(X))) \\ \mathcal{E}_3 &:= \widetilde{O}(\sigma_0 \gamma k)^{q-1} \gamma P & \mathcal{E}_{4,j,\ell}(X) &:= \widetilde{O}(\sigma_0)^{q-1} \mathbb{1}_{v_{j,\ell} \in \mathcal{V}(X)} \end{aligned}$$

we first bound the positive gradient (namely for $i = y$):

Claim D.7 (positive gradient). *Suppose Induction Hypothesis C.3 holds at iteration t . For every $(X, y) \in \mathcal{Z}$, every $r \in [m]$, every $\ell \in [2]$, and $i = y$, we have*

- (a) $\langle -\nabla_{w_{i,r}} L(F^{(t)}; X, y), v_{i,\ell} \rangle \geq (V_{i,r,\ell}(X) - \widetilde{O}(\sigma_p P)) (1 - \mathbf{logit}_i(F^{(t)}, X))$
- (b) $\langle -\nabla_{w_{i,r}} L(F^{(t)}; X, y), v_{i,\ell} \rangle \leq (V_{i,r,\ell}(X) + \mathcal{E}_1 + \mathcal{E}_3) (1 - \mathbf{logit}_i(F^{(t)}, X))$
- (c) For every $j \in [k] \setminus \{i\}$,

$$|\langle \nabla_{w_{i,r}} L(F^{(t)}, X, y), v_{j,\ell} \rangle| \leq (1 - \mathbf{logit}_i(F^{(t)}, X)) (\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X))$$

Proof of Claim D.7. We drop the superscript (t) for notational simplicity. Using the gradient formula from (D.1) (in the case of $i = y$), and the orthogonality among feature vectors, we have

$$\begin{aligned} \langle -\nabla_{w_{i,r}} L(F; X, y), v_{j,\ell} \rangle &= (1 - \mathbf{logit}_i(F, X)) \times \\ &\left(\mathbb{1}_{v_{j,\ell} \in \mathcal{V}(X)} \sum_{p \in \mathcal{P}_{v_{j,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) z_p + \sum_{p \in [P]} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) \alpha_{p,v_{j,\ell}} \pm \sum_{p \in [P]} |\langle v_{j,\ell}, \xi_p \rangle| \right) \end{aligned}$$

Using the randomness of ξ_p , we have (recalling $\sigma_p = \frac{1}{\sqrt{d \text{polylog}(k)}}$ and $\gamma \leq \frac{1}{k}$)

$$\sum_{p \in [P]} |\langle v_{j,\ell}, \xi_p \rangle| \leq \widetilde{O}(\sigma_p \cdot s + \frac{\gamma k}{\sqrt{d}} \cdot P) \ll \widetilde{O}(\sigma_p \cdot P)$$

When $j = i$ we have $v_{i,\ell} \in \mathcal{V}(X)$ so this proves Claim D.7a. Using Induction Hypothesis C.3, we have

- For every $p \in \mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X)$, we have: $\widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) \in [0, 1]$.
- For every $p \in \mathcal{P}(X) \setminus (\mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X))$, we have: $\widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) \in [0, \widetilde{O}(\sigma_0^{q-1})]$.
- For every $p \in [P] \setminus \mathcal{P}(X)$, we have: $\widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) \in [0, \widetilde{O}((\sigma_0 \gamma k)^{q-1})]$.

Using the sparsity from Def. 3.1, we have $|\mathcal{P}(X) \setminus (\mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X))| \leq \widetilde{O}(s)$. Combining this with $\alpha_{p,v} \leq \gamma$, and setting $j = i$, this proves Claim D.7b.

Finally, when $j \neq i$, using Induction Hypothesis C.3 we additionally have

- When $v_{j,\ell} \in \mathcal{V}(X)$ and $p \in \mathcal{P}_{v_{j,\ell}}(X)$, we have $\widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) \leq \widetilde{O}(\sigma_0^{q-1})$

- For $p \in \mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X)$, we have a more precise bound using Induction Hypothesis C.3a:

$$\begin{aligned} & \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) \leq \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, v_{i,\ell} \rangle + \tilde{o}(\sigma_0)) \\ & \leq \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, v_{i,\ell} \rangle + \tilde{o}(\sigma_0)) z_p \leq \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) z_p + \tilde{o}(\sigma_0) = V_{i,r,\ell}(X) \end{aligned}$$

Putting them together proves Claim D.7c. \square

We also have the following claim about the negative gradient (namely for $i \neq y$), whose proof is completely symmetric to that of Claim D.7 so we ignore here.

Claim D.8 (negative gradient). *Suppose Induction Hypothesis C.3 holds at iteration t . For every $(X, y) \in \mathcal{Z}$, every $r \in [m]$, every $\ell \in [2]$, and $i \in [k] \setminus \{y\}$, we have*

- (a) $\langle -\nabla_{w_{i,r}} L(F^{(t)}, X, y), v_{i,\ell} \rangle \geq -\mathbf{logit}_i(F^{(t)}, X) \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)} V_{i,r,\ell}(X) \right)$
- (b) For every $j \in [k]$: $\langle -\nabla_{w_{i,r}} L(F^{(t)}, X, y), v_{j,\ell} \rangle \leq \mathbf{logit}_i(F^{(t)}, X) \tilde{O}(\sigma_p P)$
- (c) For every $j \in [k] \setminus \{i\}$: $\langle -\nabla_{w_{i,r}} L(F^{(t)}, X, y), v_{j,\ell} \rangle \geq -\mathbf{logit}_i(F^{(t)}, X) (\mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X))$

Function approximations. Let us denote

$$\Phi_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \sum_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \quad \text{and} \quad \Phi_i^{(t)} \stackrel{\text{def}}{=} \sum_{\ell \in [2]} \Phi_{i,\ell}^{(t)} \quad (\text{D.4})$$

One can easily derive that

Claim D.9 (function approximation). *Suppose Induction Hypothesis C.3 holds at iteration t and supposes $\leq \tilde{O}(\frac{1}{\sigma_0^q m})$ and $\gamma \leq \tilde{O}(\frac{1}{\sigma_0 k (mP)^{1/q}})$. Let $Z_{i,\ell}^{(t)}(X) \stackrel{\text{def}}{=} \mathbf{1}_{v_{i,\ell} \in \mathcal{V}(X)} \left(\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \right)$, we have:*

- for every t , every $(X, y) \in \mathcal{Z}_m$ and $i \in [k]$, or for every $(X, y) \in \mathcal{Z}_s$ and $i \in [k] \setminus \{y\}$,

$$\begin{aligned} F_i^{(t)}(X) &= \sum_{\ell \in [2]} \left(\Phi_{i,\ell}^{(t)} \times Z_{i,\ell}^{(t)}(X) \right) \pm \tilde{O}(\sigma_0 + \sigma_0^q s m + (\sigma_0 \gamma k)^q \cdot P m) \\ &= \sum_{\ell \in [2]} \left(\Phi_{i,\ell}^{(t)} \times Z_{i,\ell}^{(t)}(X) \right) \pm O\left(\frac{1}{\text{polylog}(k)}\right) \end{aligned}$$

- for every $(X, y) \sim \mathcal{D}$, with probability at least $1 - e^{-\Omega(\log^2 k)}$ it satisfies for every $i \in [k]$,

$$F_i^{(t)}(X) = \sum_{\ell \in [2]} \left(\Phi_{i,\ell}^{(t)} \times Z_{i,\ell}^{(t)}(X) \right) \pm O\left(\frac{1}{\text{polylog}(k)}\right)$$

D.2 Useful Claims as Consequences of the Induction Hypothesis

In this sub-section we state some consequences of our Induction Hypothesis C.3. They shall be useful in our later proof of the induction hypothesis.

D.2.1 Correlation Growth

Claim D.10 (growth). *Suppose Induction Hypothesis C.3 holds at iteration t , then for every $i \in [k]$, suppose $\Lambda_i^{(t)} \leq O(1/m_0)$, then it satisfies*

$$\Lambda_i^{(t+1)} = \Lambda_i^{(t)} + \Theta\left(\frac{\eta}{k}\right) \widetilde{\text{ReLU}}'(\Lambda_i^{(t)})$$

Proof of Claim D.10. Recall $\Lambda_i^{(t)} \stackrel{\text{def}}{=} \max_{r \in [m], \ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$.

Now, let us take any $r \in [m]$ and $\ell \in [2]$ so that $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \tilde{\Omega}(\sigma_0)$. We first show a lower bound on the increment. By Claim D.7 and Claim D.8,

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbb{1}_{y=i} \left(V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P) \right) \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right. \\ &\quad \left. - \mathbb{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbb{1}_{v_{i,\ell} \in \mathcal{P}(X)} V_{i,r,\ell}(X) \right) \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \end{aligned} \quad (\text{D.5})$$

Recall $V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) z_p$. Using Induction Hypothesis C.3, we know that as long as $(X, y) \in \mathcal{Z}_m$, or when $(X, y) \in \mathcal{Z}_s$ but $i \neq y$, it satisfies

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \pm \tilde{o}(\sigma_0) \right) z_p$$

- When $i = y$ is the correct label, at least when $(X, y) \in \mathcal{Z}_m$, we have $\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \geq 1$, and together with $|\mathcal{P}_{v_{i,\ell}}| \leq C_p = O(1)$, this tells us $V_{i,r,\ell}(X) \geq \Omega(1) \cdot \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \right)$.
- When $i \neq y$ is the wrong label and when $v_{i,\ell} \in \mathcal{P}(X)$, we can use $z_p \leq O(1)$ to derive that $V_{i,r,\ell}(X) \leq O(1) \cdot \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \right)$.

Together with $\mathbf{logit}_i(F^{(t)}, X) \leq O(\frac{1}{k})$ from Claim D.4, we can derive that

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbb{1}_{y=i} \cdot \Omega(1) - O(1) \cdot \mathbb{1}_{y \neq i} \mathbb{1}_{v_{i,\ell} \in \mathcal{P}(X)} \frac{1}{k} \right] \cdot \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \right) \\ &\quad - \eta \tilde{O} \left(\frac{\sigma_p P + \mathcal{E}_1 + \mathcal{E}_3}{k} \right) \end{aligned}$$

Finally, recall the property of our distribution $\Pr_{(X,y) \sim \mathcal{D}} [v_{i,\ell} \in \mathcal{P}(X) \mid i \neq y] = \frac{s}{k} \ll o(1)$, we derive that

$$\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle \geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \Omega \left(\frac{\eta}{k} \right) \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \right)$$

As for the lower bound, using Claim D.7 and Claim D.8 again, we have

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\leq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbb{1}_{y=i} (V_{i,r,\ell}(X) + \mathcal{E}_1 + \mathcal{E}_3) \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right. \\ &\quad \left. - \mathbb{1}_{y \neq i} \left(\tilde{O}(\sigma_p P) \right) \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \end{aligned}$$

so a completely symmetric argument also shows

$$\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle \leq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + O \left(\frac{\eta}{k} \right) \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \right) \quad \square$$

Claim D.10 immediately gives the following corollary (using $\Omega(\sigma_0) \leq \Lambda_i^{(0)} \leq \tilde{O}(\sigma_0)$):

Claim D.11. *Suppose Induction Hypothesis C.3 holds for every iteration. Define thresholds (noticing $\Lambda_{\emptyset}^- \leq \Lambda_{\emptyset}^+$):*

$$\Lambda_{\emptyset}^- \stackrel{\text{def}}{=} \Theta \left(\frac{\varrho}{\log k} \right) = \tilde{\Theta}(1) \quad \text{and} \quad \Lambda_{\emptyset}^+ \stackrel{\text{def}}{=} \Theta \left(\frac{1}{m_0} \right) = \tilde{\Theta}(1)$$

Let $T_{0,i}$ be the first iteration so that $\Lambda_i^{(t)} \geq 2\Lambda_{\emptyset}^-$, and $T_0 \stackrel{\text{def}}{=} \Theta \left(\frac{k}{\eta \sigma_0^{q-2}} \right)$ (noticing $T_0 \geq T_{0,i}$) Then,

- for every $i \in [k]$ and $t \geq T_0$, it satisfies $\Lambda_i^{(t)} \geq \Lambda_{\emptyset}^+$

- for every $i \in [k]$ and $t \geq T_{0,i}$, it satisfies $\Lambda_i^{(t)} \geq \Lambda_{\emptyset}^-$

D.2.2 Single-View Error Till the End

In this subsection we present a claim to bound the “convergence” (namely, the $(1 - \mathbf{logit}_y(F^{(t)}, X))$ part) for every single-view data from T_0 till the end.

Claim D.12 (single view till the end). *Suppose Induction Hypothesis C.3 holds for all iterations $t < T$ and $\gamma \leq \tilde{O}(\sigma_0 k)$. We have that*

- (a) for every $(X, y) \in \mathcal{Z}_s$, for every $r \in [m]$, every $\ell \in [2]$, every $p \in \mathcal{P}_{v_{y,\ell}}(X)$

$$\sum_{t=T_0}^T (1 - \mathbf{logit}_y(F^{(t)}, X)) \widetilde{\text{ReLU}}'(\langle w_{y,r}, x_p \rangle) \leq \tilde{O}\left(\frac{N}{\eta}\right)$$

- (b) for every $(X, y) \in \mathcal{Z}_s$,

$$\sum_{t=T_0}^T (1 - \mathbf{logit}_y(F^{(t)}, X)) \leq \tilde{O}\left(\frac{N}{\eta \rho^{q-1}}\right)$$

Before proving Claim D.12, we first establish a simple claim to bound how the (correlation with the) noise term grows on single view data. This is used to show that the learner learns most single-view data through *memorization*.

Claim D.13 (noise lower bound). *Suppose Induction Hypothesis C.3 holds at iteration t .*

- (a) For every $(X, y) \in \mathcal{Z}_s$, every $\ell \in [2]$, for every $p \in \mathcal{P}_{v_{y,\ell}}(X)$,

$$\langle w_{i,r}^{(t+1)}, \xi_p \rangle \geq \langle w_{i,r}^{(t)}, \xi_p \rangle - \frac{\eta}{\sqrt{d}} + \tilde{\Omega}\left(\frac{\eta}{N}\right) \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) (1 - \mathbf{logit}_i(F^{(t)}, X^*)) \geq \dots \geq -\frac{\eta T}{\sqrt{d}}$$

- (b) For every $(X, y) \in \mathcal{Z}_s$, every $\ell \in [2]$,

$$\begin{aligned} \sum_{p \in \mathcal{P}_{v_{y,\ell}}(X)} \langle w_{y,r}^{(t+1)}, \xi_p \rangle &\geq \sum_{p \in \mathcal{P}_{v_{y,\ell}}(X)} \langle w_{y,r}^{(t)}, \xi_p \rangle - \frac{O(\eta)}{\sqrt{d}} \\ &\quad + \tilde{\Omega}\left(\frac{\eta}{N}\right) \widetilde{\text{ReLU}}'(\rho \cdot \langle w_{y,r}^{(t)}, v_{y,\ell} \rangle - \tilde{O}(\eta T / \sqrt{d} + \sigma_0 \gamma k)) (1 - \mathbf{logit}_y(F^{(t)}, X)) \end{aligned}$$

Proof of Claim D.13. For every $(X^*, y^*) \in \mathcal{Z}_s$, every $i \in [k]$, every $\ell \in [2]$, and every $p^* \in \mathcal{P}_{v_{i,\ell}}(X^*)$, one can calculate that

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, \xi_{p^*} \rangle &= \langle w_{i,r}^{(t)}, \xi_{p^*} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbf{1}_{y=i} \left(\sum_{p \in [P]} \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) \langle x_p, \xi_{p^*} \rangle \right) (1 - \mathbf{logit}_i(F^{(t)}, X)) \right. \\ &\quad \left. - \mathbf{1}_{y \neq i} \left(\sum_{p \in [P]} \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) \langle x_p, \xi_{p^*} \rangle \right) \mathbf{logit}_i(F^{(t)}, X) \right] \end{aligned}$$

Note when $X \neq X^*$, we have $|\langle x_p, \xi_{p^*} \rangle| \leq \tilde{O}(\sigma_p) \leq o(\frac{1}{\sqrt{d}})$; and when $X = X^*$ but $p \neq p^*$, we also have $|\langle x_p, \xi_{p^*} \rangle| \leq \tilde{O}(\sigma_p) \leq o(\frac{1}{\sqrt{d}})$. Therefore, when $i = y^*$,

$$\langle w_{i,r}^{(t+1)}, \xi_{p^*} \rangle = \langle w_{i,r}^{(t)}, \xi_{p^*} \rangle + \tilde{\Theta}\left(\frac{\eta}{N}\right) \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_{p^*} \rangle) (1 - \mathbf{logit}_i(F^{(t)}, X^*)) \pm \frac{\eta}{\sqrt{d}} \quad (\text{D.6})$$

Using the non-negativity of $\widetilde{\text{ReLU}}'$ we arrive at the first conclusion. Next, using Induction Hypothesis C.3d, we have $\langle w_{i,r}^{(t)}, x_{p^*} \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_{p^*} + \langle w_{i,r}^{(t)}, \xi_{p^*} \rangle \pm \tilde{O}(\sigma_0 \gamma k)$. Also, recall from Def. 3.1 that $\sum_{p^* \in \mathcal{P}_{v_{i,\ell}}(X^*)} z_{p^*} \geq$

$\Omega(\rho)$. Therefore, when summing over constantly many $p^* \in P_{v_i, \ell}(X^*)$ we have

$$\begin{aligned} \sum_{p^* \in P_{v_i, \ell}(X^*)} \langle w_{i,r}^{(t+1)}, \xi_{p^*} \rangle &\geq \sum_{p^* \in P_{v_i, \ell}(X^*)} \langle w_{i,r}^{(t)}, \xi_{p^*} \rangle - \frac{O(\eta)}{\sqrt{d}} \\ &\quad + \widetilde{\Omega}\left(\frac{\eta}{N}\right) \widetilde{\text{ReLU}}' \left(\rho \cdot \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle - \widetilde{O}(\eta T / \sqrt{d} + \sigma_0 \gamma k) \right) \left(1 - \mathbf{logit}_i(F^{(t)}, X^*) \right) \end{aligned}$$

This arrives at our second conclusion. \square

Proof of Claim D.12. We now prove Claim D.12 using Claim D.13. Let us denote $i = y$.

Claim D.12a is in fact a direct corollary of Claim D.13a, because once the summation has reached $\widetilde{\Omega}\left(\frac{N}{\eta}\right)$ at some iteration $t = t_0$, then according to Claim D.13a, we must have already satisfied

$$\forall t \geq t_0: \quad \langle w_{i,r}^{(t)}, \xi_p \rangle \geq \text{polylog}(k)$$

but according to $\langle w_{i,r}^{(t)}, x_p \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p + \langle w_{i,r}^{(t)}, \xi_p \rangle \pm \widetilde{O}(\sigma_0 \gamma k)$ from Induction Hypothesis C.3c, and using $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq -1$ from Induction Hypothesis C.3h, we immediately have

$$F_i^{(t)}(X) \geq \langle w_{i,r}^{(t)}, x_{p^*} \rangle - O(1) \geq \text{polylog}(k)$$

while at the same time, one can easily derive (recall (D.3)) that $F_j^{(t)}(X) \leq m_0 \cdot \Lambda_i^{(t)} \cdot \Gamma \leq O(1)$ for every $j \neq i$. Therefore, we have $1 - \mathbf{logit}_y(F^{(t)}, X) \leq e^{-\log^5 k}$ for every $t \geq t_1$. This proves the Claim D.12a.

Next, we move to Claim D.12b. We prove by way of contradiction and suppose

$$\sum_{t \geq T_0} \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \geq \widetilde{\Omega}\left(\frac{N}{\eta \rho^{q-1}}\right)$$

Using $\Lambda_i^{(t)} = \max_{r \in [m], \ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \geq \widetilde{\Omega}(1)$ from Claim D.11 and the definition of $\mathcal{M}_i^{(0)}$, we have

$$\sum_{(r,\ell) \in \mathcal{M}_i^{(0)} \times [2]} \mathbb{1}_{\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \widetilde{\Omega}(1)} \sum_{t \geq T_0} \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \geq \widetilde{\Omega}\left(\frac{N}{\eta \rho^{q-1}}\right)$$

Note that when $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \widetilde{\Omega}(1)$ and $\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle \geq \text{polylog}(k)$ simultaneously hold, there must exist some $p^* \in \mathcal{P}_{v_{i,\ell}}(X)$ so that $\langle w_{i,r}^{(t)}, \xi_{p^*} \rangle \geq \text{polylog}(k)$, but according to Induction Hypothesis C.3d, we have (noticing $\widetilde{\text{ReLU}}$ is in the linear regime now because $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq 0$)

$$F_i^{(t)}(X) \geq \langle w_{i,r}^{(t)}, x_{p^*} \rangle - O(1) = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_{p^*} + \langle w_{i,r}^{(t)}, \xi_{p^*} \rangle \pm \widetilde{O}(\sigma_0 \gamma k) - O(1) \geq \text{polylog}(k)$$

In contrast, one can derive (recall (D.3)) that $F_j^{(t)}(X) \leq m_0 \cdot \Lambda_i^{(t)} \cdot \Gamma \leq O(1)$ for every $j \neq i$. This means $1 - \mathbf{logit}_i(F^{(t)}, X) \ll e^{-\log^5 k}$. In other words,

$$\sum_{t \geq T_0} \sum_{(r,\ell) \in \mathcal{M}_i^{(0)} \times [2]} \mathbb{1}_{\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \widetilde{\Omega}(1)} \mathbb{1}_{\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle \leq \text{polylog}(k)} \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \geq \widetilde{\Omega}\left(\frac{N}{\eta \rho^{q-1}}\right)$$

Now we partition the iterations between T_0 and T into $4m_0$ stages of consecutive iterations, denoted

by $\mathcal{T}_1, \dots, \mathcal{T}_{4m_0}$, so that each of them have a similar partial sum in the above summation. In symbols:

$\forall g \in [4m_0]$:

$$\sum_{t \in \mathcal{T}_g} \sum_{(r, \ell) \in \mathcal{M}_i^{(0)} \times [2]} \mathbb{1}_{\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \tilde{\Omega}(1)} \mathbb{1}_{\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle \leq \text{polylog}(k)} \left(1 - \mathbf{logit}_i \left(F^{(t)}, X\right)\right) \geq \tilde{\Omega} \left(\frac{N}{\eta \rho^{q-1}}\right) \quad (\text{D.7})$$

Let us first look at stage \mathcal{T}_1 . By averaging, there exists some $(r, \ell) = (r_1^*, \ell_1^*) \in \mathcal{M}_i^{(0)} \times [2]$ so that

$$\sum_{t \in \mathcal{T}_1} \mathbb{1}_{\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \tilde{\Omega}(1)} \mathbb{1}_{\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle \leq \text{polylog}(k)} \left(1 - \mathbf{logit}_i \left(F^{(t)}, X\right)\right) \geq \tilde{\Omega} \left(\frac{N}{\eta \rho^{q-1}}\right)$$

Applying Claim D.13b, we know that after stage \mathcal{T}_1 (namely, for any $t \in [\max \mathcal{T}_1, T]$), it satisfies

$$\text{for } (r, \ell) = (r_1^*, \ell_1^*) \quad \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle \geq \tilde{\Omega} \left(\frac{N}{\eta \rho^{q-1}}\right) \cdot \tilde{\Omega} \left(\frac{\eta}{N} \rho^{q-1}\right) > \text{polylog}(k)$$

Continuing to stage \mathcal{T}_2 , by averaging again, we can find some other $(r, \ell) = (r_2^*, \ell_2^*) \in \mathcal{M}_i^{(0)} \times [2]$ so that

$$\sum_{t \in \mathcal{T}_2} \mathbb{1}_{\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \tilde{\Omega}(1)} \mathbb{1}_{\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle \leq \text{polylog}(k)} \left(1 - \mathbf{logit}_i \left(F^{(t)}, X\right)\right) \geq \tilde{\Omega} \left(\frac{N}{\eta \rho^{q-1}}\right)$$

From the conclusion of the previous stage, it must satisfy that $(r_2^*, \ell_2^*) \neq (r_1^*, \ell_1^*)$. A similar derivation also tells us that after stage \mathcal{T}_2 (namely, for any $t \in [\max \mathcal{T}_2, T]$), it satisfies

$$\text{for } (r, \ell) = (r_2^*, \ell_2^*) \quad \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle > \text{polylog}(k)$$

We continue this argument until we finish stage \mathcal{T}_{2m_0} . At this point, we know for every $t \in [\max \mathcal{T}_{2m_0}, T]$

$$\text{for all } (r, \ell) \in \mathcal{M}_i^{(0)} \times [2] \quad \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \langle w_{i,r}^{(t)}, \xi_p \rangle > \text{polylog}(k)$$

This contradicts (D.7) for any $g > 2m_0$. □

D.2.3 Multi-View Error Till the End

In this subsection we present a claim to bound the ‘‘convergence’’ (namely, the $(1 - \mathbf{logit}_y(F^{(t)}, X))$ part) for the average multi-view data from T_0 till the end.

Claim D.14 (multi-view till the end). *Suppose Induction Hypothesis C.3 holds for every iteration $t < T$, and suppose $N_s \leq \frac{k^2 \rho^{q-1}}{s}$, then*

$$\sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[1 - \mathbf{logit}_y \left(F^{(t)}, X\right)\right] \leq \tilde{O} \left(\frac{k}{\eta}\right) + \tilde{O} \left(\frac{s N_s}{\eta k \rho^{q-1}}\right) \leq \tilde{O} \left(\frac{k}{\eta}\right)$$

In fact, Claim D.14 is a direct corollary of the following claim, combined with $\Lambda_i^{(t)} = \tilde{O}(1)$ from Induction Hypothesis C.3g, and with the convergence Claim D.12b for single-view data.

Claim D.15. *Suppose Induction Hypothesis C.3 holds at iteration t and $t \geq T_0$, then*

$$\begin{aligned} \sum_{i \in [k]} \Lambda_i^{(t+1)} &\geq \sum_{i \in [k]} \Lambda_i^{(t)} + \Omega(\eta) \times \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right] \\ &\quad - \eta O \left(\frac{s N_s}{k N} \right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right] \end{aligned}$$

Proof of Claim D.15. Recall $\Lambda_i^{(t)} \stackrel{\text{def}}{=} \max_{r \in [m], \ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$. Let us take r, ℓ to be this argmax so that Claim D.11 tells us $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \Lambda_{\varnothing}^+ = \tilde{\Theta}(1)$. By Claim D.7 and Claim D.8,

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbf{1}_{y=i} \left(V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P) \right) \left(1 - \mathbf{logit}_i \left(F^{(t)}, X \right) \right) \right. \\ &\quad \left. - \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)} V_{i,r,\ell}(X) \right) \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \end{aligned}$$

Recall $V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) z_p$. Using Induction Hypothesis C.3a, we know that as long as $(X, y) \in \mathcal{Z}_m$, or when $(X, y) \in \mathcal{Z}_s$ but $i \neq y$, it satisfies

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \pm \tilde{o}(\sigma_0)) z_p$$

Since $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \Lambda_{\varnothing}^+ = \Theta(\frac{1}{m_0}) \gg \varrho$ (see Claim D.11) and since $|\mathcal{P}_{v_{i,\ell}}(X)| \leq O(1)$, for most of $p \in \mathcal{P}_{v_{i,\ell}}$ we must be already in the linear regime of $\widetilde{\text{ReLU}}$ so

$$0.9 \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \leq V_{i,r,\ell}(X) \leq \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p$$

According to our choice of the distribution (see Def. 3.1):

- When $(X, y) \in \mathcal{Z}_m$ and $y = i$, we have $V_{i,r,\ell}(X) \geq 0.9$.
- When $(X, y) \in \mathcal{Z}_s$ and $y = i$, we have $V_{i,r,\ell}(X) \geq 0$.
- When $(X, y) \in \mathcal{Z}_m$, $y \neq i$ and $v_{i,\ell} \in \mathcal{P}(X)$, we have $V_{i,r,\ell}(X) \leq 0.4$.
- When $(X, y) \in \mathcal{Z}_s$, $y \neq i$ and $v_{i,\ell} \in \mathcal{P}(X)$, we have $V_{i,r,\ell}(X) \leq \Gamma \leq 1$.

Together, we derive that

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[0.89 \cdot \mathbf{1}_{y=i} \left(1 - \mathbf{logit}_i \left(F^{(t)}, X \right) \right) \right] \\ &\quad - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + 0.41 \mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)} \right) \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \\ &\quad - O \left(\frac{\eta N_s}{N} \right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbf{1}_{y=i} \cdot \tilde{O}(\sigma_p P) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \\ &\quad - O \left(\frac{\eta N_s}{k N} \right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)} \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \end{aligned} \tag{D.8}$$

Above, we have applied Claim D.4 which says for $(X, y) \in \mathcal{Z}_s$, it holds that $\mathbf{logit}_i(F^{(t)}, X) \leq O(\frac{1}{k}) (1 - \mathbf{logit}_y(F^{(t)}, X))$.

Finally, substituting $\mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)}$ with the naive upper bound $\mathcal{E}_1 + \mathcal{E}_3 + 0.41 \mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)} \leq 0.41$, we have

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[0.89 \cdot \mathbf{1}_{y=i} \left(1 - \mathbf{logit}_i \left(F^{(t)}, X \right) \right) - 0.41 \cdot \mathbf{1}_{y \neq i} \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \\ &\quad - O \left(\frac{\eta N_s}{k N} \right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\left(k \mathbf{1}_{y=i} \cdot \tilde{O}(\sigma_p P) + \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)} \right) \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \end{aligned}$$

Summing up over all $i \in [k]$, and using $v_{i,\ell} \in \mathcal{P}(X)$ with probability $\frac{s}{k}$ when $i \neq y$, we finish the proof. \square

D.2.4 Multi-View Individual Error

Our next claim states that up to a polynomial factor, the error on any individual multi-view data is bounded by the training error.

Claim D.16 (multi-view individual error). *For every $t \geq 0$, every $(X, y) \in \mathcal{Z}_m$,*

$$1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \leq \tilde{O} \left(\frac{k^4}{s^2} \right) \cdot \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right]$$

(The same also holds w.p. $\geq 1 - e^{-\Omega(\log^2 k)}$ for every $(X, y) \sim \mathcal{D}_m$ on the left hand side.)

Furthermore, if $\mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [1 - \mathbf{logit}_y (F^{(t)}, X)] \leq \frac{1}{k^4}$ is sufficiently small, we have $0.4\Phi_i^{(t)} - \Phi_j^{(t)} \leq -\Omega(\log k)$ for every pair $i, j \in [k]$.

Proof. For a data point $(X, y) \in \mathcal{Z}_m$, let us denote by $\mathcal{H}(X)$ be the set of all $i \in [k] \setminus \{y\}$ such that,

$$\sum_{\ell \in [2]} \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \geq 0.8 - \frac{1}{100 \log(k)}, \quad \sum_{\ell \in [2]} \sum_{p \in \mathcal{P}_{v_{y,\ell}}(X)} z_p \leq 2 + \frac{1}{100 \log(k)}$$

Now, suppose $1 - \mathbf{logit}_y (F^{(t)}, X) = \xi(X)$, then using $\min\{1, \beta\} \leq 2(1 - \frac{1}{1+\beta})$, we have

$$\min \left\{ 1, \sum_{i \in [k] \setminus \{y\}} e^{F_i^{(t)}(X) - F_y^{(t)}(X)} \right\} \leq 2\xi(X)$$

By Claim D.9 and our definition of $\mathcal{H}(X)$, this implies that

$$\min \left\{ 1, \sum_{i \in \mathcal{H}(X)} e^{0.4\Phi_i^{(t)} - \Phi_y^{(t)}} \right\} \leq 4\xi(X)$$

If we denote by $\psi = \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [1 - \mathbf{logit}_y (F^{(t)}, X)]$, then

$$\begin{aligned} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\min \left\{ 1, \sum_{i \in \mathcal{H}(X)} e^{0.4\Phi_i^{(t)} - \Phi_y^{(t)}} \right\} \right] &\leq 4\psi \\ \implies \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\sum_{i \in \mathcal{H}(X)} \min \left\{ \frac{1}{k}, e^{0.4\Phi_i^{(t)} - \Phi_y^{(t)}} \right\} \right] &\leq 4\psi \end{aligned}$$

Notice that we can rewrite the LHS so that

$$\begin{aligned} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\sum_{j \in [k]} \mathbb{1}_{j=y} \sum_{i \in [k]} \mathbb{1}_{i \in \mathcal{H}(X)} \min \left\{ \frac{1}{k}, e^{0.4\Phi_i^{(t)} - \Phi_j^{(t)}} \right\} \right] &\leq 4\psi \\ \implies \sum_{j \in [k]} \sum_{i \in [k]} \mathbb{1}_{i \neq j} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{j=y} \mathbb{1}_{i \in \mathcal{H}(X)} \right] \min \left\{ \frac{1}{k}, e^{0.4\Phi_i^{(t)} - \Phi_j^{(t)}} \right\} &\leq 4\psi \end{aligned}$$

Note however, that for every $i \neq j \in [k]$, the probability of generating a multi-view sample $(X, y) \in \mathcal{Z}_m$ with $y = j$ and $i \in \mathcal{H}(X)$ is at least $\tilde{\Omega}(\frac{1}{k} \cdot \frac{s^2}{k^2})$. This implies

$$\sum_{j \in [k]} \sum_{i \in [k] \setminus i} \min \left\{ \frac{1}{k}, e^{0.4\Phi_i^{(t)} - \Phi_j^{(t)}} \right\} \leq \tilde{O} \left(\frac{k^3}{s^2} \psi \right) \quad (\text{D.9})$$

Finally, using $1 - \frac{1}{1+\beta} \leq \min\{1, \beta\}$, it is easy to see for every $(X, y) \in \mathcal{Z}_m$

$$1 - \mathbf{logit}_y \left(F^{(t)}, X \right) = \min \left\{ 1, \sum_{i \in [k] \setminus \{y\}} 2e^{0.4\Phi_i^{(t)} - \Phi_y^{(t)}} \right\} \leq k \cdot \sum_{i \in [k] \setminus \{y\}} \min \left\{ \frac{1}{k}, e^{0.4\Phi_i^{(t)} - \Phi_y^{(t)}} \right\} \leq \tilde{O} \left(\frac{k^4}{s^2} \psi \right)$$

We complete the proof.

Note that if one replaces $(X, y) \in \mathcal{Z}_m$ with $(X, y) \in \mathcal{D}_m$, we also have

$$1 - \mathbf{logit}_y \left(F^{(t)}, X \right) = \min \left\{ 1, \sum_{i \in [k] \setminus \{y\}} 2e^{0.4\Phi_i^{(t)} - \Phi_y^{(t)}} \right\}$$

with high probability, so the same result also holds.

Note also (D.9) implies if $\mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [1 - \mathbf{logit}_y (F^{(t)}, X)] \leq \frac{1}{k^4}$ is sufficiently small, we have $0.4\Phi_i^{(t)} - \Phi_j^{(t)} \leq -\Omega(\log k)$ for every pair $i \neq j$. Using the non-negativity of $\Phi_i^{(t)}$, we know the relationship holds also when $i = j$. □

D.2.5 Multi-View Error in Stage 2

As we shall see later, our final proof is divided into three stages for each index $i \in [k]$: the first stage is for $t \leq T_{0,i}$, the second stage is for all $t \in [T_{0,i}, T_0]$, and the third iteration is for $t > T_0$. We have the following claim to bound the maximum error of multi-view data during the second stage.

Claim D.17 (multi-view stage 2). *Suppose Induction Hypothesis C.3 holds for every iteration $t \leq T_0$, and $\Upsilon = \tilde{\Theta}(\frac{1}{k^{0.2}})$ is a parameter. Then, for every $i \in [k]$*

(a)

$$\sum_{t=T_{0,i}}^{T_0} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [\mathbf{1}_{y=i} (1 - \mathbf{logit}_y (F^{(t)}, X))] \leq O\left(\frac{s}{k} T_0 \Upsilon\right) + \tilde{O}\left(\frac{1}{\eta}\right)$$

(b) for every $t \in [T_{0,i}, T_0]$, every $j \in [k] \setminus \{i\}$, every $\ell \in [2]$,

$$\begin{aligned} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [\mathbf{1}_{y \neq i} \mathbf{logit}_i (F^{(t)}, X)] &\leq O\left(\frac{1}{k}\right) \\ \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [\mathbf{1}_{y \neq i} \mathbf{1}_{v_{j,\ell} \in \mathcal{P}(X)} \mathbf{logit}_i (F^{(t)}, X)] &\leq O\left(\frac{s}{k^2}\right) \end{aligned}$$

In order to prove Claim D.17 we first establish the following claim.

Claim D.18. *Let Υ be any $\Upsilon \in [\frac{1}{k}, \frac{1}{s}]$, and recall $\Phi_i^{(t)} \stackrel{\text{def}}{=} \sum_{r \in [m], \ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$ from (D.4). Then, letting $T_1 \stackrel{\text{def}}{=} \tilde{\Theta}(\frac{k^{2.5}\Upsilon^{2.5}}{\eta})$ and suppose Induction Hypothesis C.3 holds for all iterations $t \leq T_1$. Then,*

$$\forall t \leq T_1, \quad \forall i \in [k]: \quad e^{0.4\Phi_i^{(t)}} \leq k\Upsilon.$$

(Note when $\Upsilon \leq \tilde{O}(k^{-0.2})$ we have $T_0 \leq T_1$.)

Proof of Claim D.18. Recall from Induction Hypothesis C.3i that for those $r \in [m] \setminus \mathcal{M}_i^{(0)}$ and $\ell \in [2]$, it satisfies $[\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \leq \tilde{O}(\sigma_0)$ and their summation does not exceed $\frac{1}{\text{poly}(\log(k))}$ due to our choice of m . Thus, to prove this claim, it suffices to slightly abuse the notation and think of

$$\widehat{\Phi}^{(t)} = \max_{i \in [k]} \sum_{(r,\ell) \in \mathcal{M}_i^{(0)} \times [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$$

Let i be this argmax in $\widehat{\Phi}^{(t)}$. For every $(r, \ell) \in \mathcal{M}_i^{(0)} \times [2]$, by Claim D.7 and Claim D.8,

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\leq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbf{1}_{y=i} (V_{i,r,\ell}(X) + \mathcal{E}_1 + \mathcal{E}_3) (1 - \mathbf{logit}_i (F^{(t)}, X)) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} (\tilde{O}(\sigma_p P)) \mathbf{logit}_i (F^{(t)}, X) \right] \end{aligned} \tag{D.10}$$

where recall

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) z_p \leq \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \leq O(1)$$

Therefore,

$$\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle \leq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + O(\eta) \left(\mathbb{E}_{(X,y) \sim \mathcal{Z}} [\mathbb{1}_{y=i}(1 - \mathbf{logit}_y(F^{(t)}, X))] + \tilde{O}(\sigma_p P) \right)$$

Note that single-view data contribute to at most $O(\frac{\eta N_s}{kN})$ on the RHS of (D.10), so we only focus on those $(X, y) \in \mathcal{Z}_m$. By Claim D.9 we know $F_y^{(t)}(X) \geq \Phi_y^{(t)} - \frac{1}{\text{polylog}(k)}$ and for $j \in [k] \setminus \{y\}$,

- W.p. $1 - (1 - \frac{s}{k})^2$, both $v_{j,1}, v_{j,2} \notin \mathcal{P}(X)$, and in this case $F_j^{(t)}(X) \leq \frac{1}{\text{polylog}(k)}$;
- W.p. $(1 - \frac{s}{k})^2$, at least one of $v_{j,1}, v_{j,2} \in \mathcal{P}(X)$, and in this case $F_j^{(t)}(X) \leq 0.4\Phi_j^{(t)} + \frac{1}{\text{polylog}(k)}$;

Together, and using the inequality $1 - \mathbf{logit}_y(F^{(t)}, X) \leq \frac{\sum_{j \neq y} e^{F_j^{(t)}(X)}}{e^{F_y^{(t)}(X)}}$, and conclude that

$$\mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [\mathbb{1}_{i=y}(1 - \mathbf{logit}_y(F^{(t)}, X))] \leq \frac{1}{k} \cdot O\left(\frac{se^{0.4\hat{\Phi}^{(t)}} + k}{e^{\hat{\Phi}^{(t)}}}\right)$$

Summing up over all $(r, \ell) \in \mathcal{M}_i^{(0)} \times [2]$ with $|\mathcal{M}_i^{(0)}| \leq m_0 \leq \tilde{O}(1)$, we have

$$\hat{\Phi}^{(t+1)} \leq \hat{\Phi}^{(t)} + \eta \frac{1}{k} \tilde{O}\left(\frac{se^{0.4\hat{\Phi}^{(t)}} + k}{e^{\hat{\Phi}^{(t)}}} + \frac{N_s}{N}\right)$$

This implies whenever $e^{\hat{\Phi}^{(t)}} = \Omega(k^{2.5}\Upsilon^{2.5})$, we have

$$\hat{\Phi}^{(t+1)} \leq \hat{\Phi}^{(t)} + \eta \tilde{O}\left(\frac{1}{k^{2.5}\Upsilon^{2.5}}\right)$$

This finishes the proof of Claim D.18. \square

Proof of Claim D.17. We first prove Claim D.17a and the proof of Claim D.17b is only simpler.

Using (D.8) in the proof of Claim D.15, we know as long as $t \geq T_{0,i}$ (so $\Lambda_i^{(t)} \geq \Lambda_\emptyset$),¹⁷

$$\begin{aligned} \Lambda_i^{(t+1)} &\geq \Lambda_i^{(t)} + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\tilde{\Omega}(1) \cdot \mathbb{1}_{y=i} \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right] \\ &\quad - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y \neq i} (\mathcal{E}_1 + \mathcal{E}_3 + 0.4\mathbb{1}_{v_{i,1} \text{ or } v_{i,2} \in \mathcal{P}(X)}) \mathbf{logit}_i(F^{(t)}, X) \right] \\ &\quad - O\left(\frac{\eta N_s}{N}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbb{1}_{y=i} \cdot \tilde{O}(\sigma_p P) \left(1 - \mathbf{logit}_y(F^{(t)}, X) \right) \right] \\ &\quad - O\left(\frac{\eta N_s}{kN}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbb{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbb{1}_{v_{i,1} \text{ or } v_{i,2} \in \mathcal{P}(X)} \right) \left(1 - \mathbf{logit}_y(F^{(t)}, X) \right) \right] \end{aligned}$$

Applying Claim D.9, we have for $(X, y) \in \mathcal{Z}_m$ and $i \neq y$, if at least one of $\{v_{i,1}, v_{i,2}\}$ is in $\mathcal{P}(X)$, then $F_i^{(t)}(X) \leq 0.4\Phi_i^{(t)} + \frac{1}{\text{polylog}(k)}$. Therefore,

$$\begin{aligned} &\mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y \neq i} \mathbb{1}_{v_{i,1} \text{ or } v_{i,2} \in \mathcal{P}(X)} \mathbf{logit}_i(F^{(t)}, X) \right] \\ &\leq \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y \neq i} \mathbb{1}_{v_{i,1} \text{ or } v_{i,2} \in \mathcal{P}(X)} \frac{1}{1 + \sum_{j \in [k]} e^{F_j^{(t)}(X) - 0.4\Phi_j^{(t)}}} \right] \\ &\leq \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y \neq i} \mathbb{1}_{v_{i,1} \text{ or } v_{i,2} \in \mathcal{P}(X)} O(\Upsilon) \right] \leq O\left(\frac{s\Upsilon}{k}\right) \end{aligned} \tag{D.11}$$

¹⁷We cite (D.8) to reduce redundancy in the proofs. The only difference is that this time we have $t \geq T_{0,i}$ instead of $t \geq T_0$, and according to Claim D.11, this requires us to change the constant 0.89 to $\frac{1}{\text{polylog}(k)} = \tilde{\Theta}(1)$. We thank Bobby He for noticing this minor typo in the earlier version of this paper.

Above, the last inequality uses $F_j^{(t)}(X) \geq 0$ and Claim D.18 which says $e^{0.4\Phi_i^{(t)}} \leq O(k\Upsilon)$. Also, for obvious reason

$$\mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\left(k \mathbb{1}_{y=i} \cdot \tilde{O}(\sigma_P P) + \mathbb{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbb{1}_{\{v_{i,1}, v_{i,2}\} \in \mathcal{P}(X)} \right) \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \leq O\left(\frac{s}{k}\right)$$

Together, we arrive at the conclusion that

$$\Lambda_i^{(t+1)} \geq \Lambda_i^{(t)} + \Omega(\eta) \left(\mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y=i} \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right] - O\left(\frac{s\Upsilon}{k} + \frac{s}{k^2} \frac{N_s}{N}\right) \right)$$

Using $\Lambda_i^{(t)} \leq \tilde{O}(1)$ from Induction Hypothesis C.3g we immediately finish the proof of Claim D.17a.

As for the first part of Claim D.17b, note that when $v_{i,1} \notin P(x)$ and $v_{i,2} \notin P(x)$, we have $\mathbf{logit}_i(F^{(t)}, X) \leq O(\frac{1}{k})$; but if one of them belongs to $P(X)$ the probability is $\frac{s}{k}$ and we have (D.11). They together (and using $s\Upsilon \leq 1$) imply the first part of Claim D.17b.

As for the second part of Claim D.17b, it is similar to the first part once we take into account $\Pr[v_{j,\ell} \in \mathcal{P}(X)] = \frac{s}{k}$. \square

D.3 Tensor Power Method Bound

In this subsection we establish a lemma for comparing the growth speed of two sequences of updates of the form $x_{t+1} \leftarrow x_t + \eta C_t x_t^{q-1}$. This should be reminiscent of the classical analysis of the growth of eigenvalues on the (incremental) tensor power method of degree q .

Lemma D.19. *Let $q \geq 3$ be a constant and $x_0, y_0 = o(1)$. Let $\{x_t, y_t\}_{t \geq 0}$ be two positive sequences updated as*

- $x_{t+1} \geq x_t + \eta C_t x_t^{q-1}$ for some $C_t = \Theta(1)$, and
- $y_{t+1} \leq y_t + \eta S C_t y_t^{q-1}$ for some constant $S = \Theta(1)$.

Suppose $x_0 \geq y_0 S^{\frac{1}{q-2}} \left(1 + \frac{1}{\text{polylog}(k)} \right)$, then we must have for every $A = O(1)$, let T_x be the first iteration such that $x_t \geq A$, then

$$y_{T_x} \leq O(y_0 \cdot \text{polylog}(k))$$

We first establish a claim before proving Lemma D.19.

Claim D.20. *Consider an increasing sequence $x_t \geq 0$ defined as $x_{t+1} = x_t + \eta C_t x_t^{q-1}$ for some $C_t = \Theta(1)$, then we have for every $A > x_0$, every $\delta \in (0, 1)$, and every $\eta \in (0, 1)$:*

$$\sum_{t \geq 0, x_t \leq A} \eta C_t \geq \left[\frac{\delta(1+\delta)^{-1}}{(1+\delta)^{q-2} - 1} \left(1 - \left(\frac{(1+\delta)x_0}{A} \right)^{q-2} \right) - \frac{O(\eta A^{q-1})}{x_0} \frac{\log\left(\frac{A}{x_0}\right)}{\log(1+\delta)} \right] \cdot \frac{1}{x_0^{q-2}}$$

$$\sum_{t \geq 0, x_t \leq A} \eta C_t \leq \left[\frac{(1+\delta)^{q-2}}{(q-2)} + \frac{O(\eta A^{q-1})}{x_0} \frac{\log\left(\frac{A}{x_0}\right)}{\log(1+\delta)} \right] \cdot \frac{1}{x_0^{q-2}}$$

Proof of Claim D.20. For every $g = 0, 1, 2, \dots$, let \mathcal{T}_g be the first iteration such that $x_t \geq (1+\delta)^g x_0$. Let b be the smallest integer such that $(1+\delta)^b x_0 \geq A$. Suppose for notation simplicity that we replace x_t with exactly A whenever $x_t \geq A$.

By the definition of \mathcal{T}_g , we have

$$\sum_{t \in [\mathcal{T}_g, \mathcal{T}_{g+1})} \eta C_t [(1+\delta)^g x_0]^{(q-1)} \leq x_{\mathcal{T}_{g+1}} - x_{\mathcal{T}_g} \leq \delta(1+\delta)^g x_0 + O(\eta A^{q-1})$$

$$\sum_{t \in [\mathcal{T}_g, \mathcal{T}_{g+1})} \eta C_t [(1+\delta)^{g+1} x_0]^{(q-1)} \geq x_{\mathcal{T}_{g+1}} - x_{\mathcal{T}_g} \geq \delta(1+\delta)^g x_0 - O(\eta A^{q-1})$$

These imply that

$$\begin{aligned}\sum_{t \in [\mathcal{T}_g, \mathcal{T}_{g+1})} \eta C_t &\leq \frac{\delta}{(1+\delta)^{g(q-2)}} \frac{1}{x_0^{q-2}} + \frac{O(\eta A^{q-1})}{x_0^{q-1}} \\ \sum_{t \in [\mathcal{T}_g, \mathcal{T}_{g+1})} \eta C_t &\geq \frac{\delta}{(1+\delta)^{g(q-2)}(1+\delta)^{q-1}} \frac{1}{x_0^{q-2}} - \frac{O(\eta A^{q-1})}{x_0^{q-1}}\end{aligned}$$

Recall b is the smallest integer such that $(1+\delta)^b x_0 \geq A$, so we can calculate

$$\begin{aligned}\sum_{t \geq 0, x_t \leq A} \eta C_t &\leq \sum_{g=0}^{b-1} \frac{\delta}{(1+\delta)^{g(q-2)}} \frac{1}{x_0^{q-2}} + \frac{O(\eta A^{q-1})}{x_0^{q-1}} b = \frac{\delta}{1 - \frac{1}{(1+\delta)^{q-2}}} \frac{1}{x_0^{q-2}} + \frac{O(\eta A^{q-1})}{x_0^{q-1}} b \\ &= \frac{\delta(1+\delta)^{q-2}}{(1+\delta)^{q-2} - 1} \frac{1}{x_0^{q-2}} + \frac{O(\eta A^{q-1})}{x_0^{q-1}} b \leq \frac{(1+\delta)^{q-2}}{(q-2)} \frac{1}{x_0^{q-2}} + \frac{O(\eta A^{q-1})}{x_0^{q-1}} b \\ \sum_{t \geq 0, x_t \leq A} \eta C_t &\geq \sum_{g=0}^{b-2} \frac{\delta}{(1+\delta)^{g(q-2)}(1+\delta)^{q-1}} \frac{1}{x_0^{q-2}} - \frac{O(\eta A^{q-1})}{x_0^{q-1}} b \\ &\geq \frac{\delta(1+\delta)^{-1} \left(1 - \frac{1}{(1+\delta)^{(q-2)(b-1)}}\right)}{(1+\delta)^{q-2} - 1} \frac{1}{x_0^{q-2}} - \frac{O(\eta A^{q-1})}{x_0^{q-1}} b \\ &\geq \frac{\delta(1+\delta)^{-1} \left(1 - \left(\frac{(1+\delta)x_0}{A}\right)^{q-2}\right)}{(1+\delta)^{q-2} - 1} \frac{1}{x_0^{q-2}} - \frac{O(\eta A^{q-1})}{x_0^{q-1}} b\end{aligned} \quad \square$$

Proof of Lemma D.19. Let us apply Claim D.20 twice, once for the x_t sequence with C_t and threshold A , and the other time for the y_t sequence with $C'_t = SC_t$ and threshold $A' = y_0 \cdot \text{polylog}(k)$. Let T_x be the first iteration t in which $x_t \geq A$, and T_y be the first iteration t in which $y_t \geq A'$.

According to Claim D.20, we know

$$\begin{aligned}\sum_{t=0}^{T_x} \eta C_t &\leq \left[\frac{(1+\delta)^{q-2}}{(q-2)} + \frac{O(\eta A^{q-1})}{x_0} \frac{\log\left(\frac{A}{x_0}\right)}{\log(1+\delta)} \right] \cdot \frac{1}{x_0^{q-2}} \\ &\leq \frac{1+O(\delta)}{(q-2)x_0^{q-2}} + O\left(\frac{\eta \log(1/x_0)}{\delta x_0^{q-1}}\right) \\ \sum_{t=0}^{T_y} \eta C'_t = \sum_{t=0}^{T_y} \eta SC_t &\geq \left[\frac{\delta(1+\delta)^{-1}}{(1+\delta)^{q-2} - 1} \left(1 - \left(\frac{(1+\delta)y_0}{A'}\right)^{q-2}\right) - \frac{O(\eta(A')^{q-1})}{y_0} \frac{\log\left(\frac{A'}{y_0}\right)}{\log(1+\delta)} \right] \cdot \frac{1}{y_0^{q-2}} \\ &\geq \frac{1 - O\left(\delta + \frac{1}{\text{polylog}(k)}\right)}{(q-2)y_0^{q-2}} - \tilde{O}\left(\frac{\eta}{\delta}\right)\end{aligned}$$

Therefore, choosing $\delta = \frac{1}{\text{polylog}(k)}$ and $\eta \leq \frac{x_0}{\log(1/x_0) \text{polylog}(k)}$, together with the assumption $x_0 \geq y_0 S^{\frac{1}{q-2}} \left(1 + \frac{1}{\text{polylog}(k)}\right)$, we immediately have that $T_x < T_y$ so we finish the proof. \square

D.4 Main Lemmas for Proving the Induction Hypothesis

In this subsection, we begin to provide key technical lemmas that, when combined together, shall prove Induction Hypothesis C.3. (We combine the analysis in the next subsection, Section D.5.)

D.4.1 Lambda Lemma

Recall $\Lambda_i^{(t)} \stackrel{\text{def}}{=} \max_{r \in [m], \ell \in [2]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$. Our first lemma shows that $\Lambda_i^{(t)}$ cannot go above $\tilde{O}(1)$.

Lemma D.21. *Suppose Induction Hypothesis C.3 holds for all iterations $< t$ and suppose $N_s \leq \tilde{O}(k/\rho)$. Then, letting $\Phi_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \sum_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$, we have*

$$\forall i \in [k], \forall \ell \in [2] : \quad \Phi_{i,\ell}^{(t)} \leq \tilde{O}(1)$$

This implies $\Lambda_i^{(t)} \leq \tilde{O}(1)$ as well.

Proof of Lemma D.21. We make a simple observation:

- For those $r \in [m] \setminus \mathcal{M}_i^{(0)}$ and $\ell \in [2]$, recall Induction Hypothesis C.3i says $[\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \leq \tilde{O}(\sigma_0)$ and their summation does not exceed $\frac{1}{\text{polylog}(k)}$ due to our choice of m .

Therefore, to prove this claim, it suffices to slightly abuse the notation and prove $\bar{\Phi}_{i,\ell}^{(t)} \leq \tilde{O}(1)$ for

$$\bar{\Phi}_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \sum_{r \in \mathcal{M}_i^{(0)}} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ = \Phi_{i,\ell}^{(t)} \pm \frac{1}{\text{poly}(k)}$$

Recall gradient descent update gives

$$\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} [\langle -\nabla_{w_{i,r}} L(F^{(t)}; X, y), v_{i,\ell} \rangle]$$

when taking the positive part, we know there exists $\Delta_{i,r,\ell}^{(t)} \in [0, 1]$ such that

$$[\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle]^+ = [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ + \eta \Delta_{i,r,\ell}^{(t)} \mathbb{E}_{(X,y) \sim \mathcal{Z}} [\langle -\nabla_{w_{i,r}} L(F^{(t)}; X, y), v_{i,\ell} \rangle]$$

Now in each iteration t , for every pair (i, ℓ) , we define a special subset $\mathcal{Z}_{s,i,\ell}$ of the single-view data as

$$\mathcal{Z}_{s,i,\ell} \stackrel{\text{def}}{=} \left\{ (X, y) \in \mathcal{Z}_s \mid y = i \wedge \widehat{\ell}(X) = 3 - \ell \right\}$$

For analysis purpose, from iteration $t = 0$ onwards, we define two sequences of quantities

- define $A_{i,\ell}^{(0)} \stackrel{\text{def}}{=} \sum_{r \in \mathcal{M}_i^{(0)}} [\langle w_{i,r}^{(0)}, v_{i,\ell} \rangle]^+$ and $B_{i,\ell}^{(0)} = 0$.
- when $t \geq 0$, define

$$\begin{aligned} A_{i,\ell}^{(t+1)} &\stackrel{\text{def}}{=} A_{i,\ell}^{(t)} + \eta \sum_{r \in \mathcal{M}_i^{(0)}} \Delta_{i,r,\ell}^{(t)} \mathbb{E}_{(X,y) \sim \mathcal{Z}} [\mathbb{1}_{(X,y) \notin \mathcal{Z}_{s,i,\ell}} \cdot \langle -\nabla_{w_{i,r}} L(F^{(t)}; X, y), v_{i,\ell} \rangle] \\ B_{i,\ell}^{(t+1)} &\stackrel{\text{def}}{=} B_{i,\ell}^{(t)} + \eta \sum_{r \in \mathcal{M}_i^{(0)}} \Delta_{i,r,\ell}^{(t)} \mathbb{E}_{(X,y) \sim \mathcal{Z}} [\mathbb{1}_{(X,y) \in \mathcal{Z}_{s,i,\ell}} \cdot \langle -\nabla_{w_{i,r}} L(F^{(t)}; X, y), v_{i,\ell} \rangle] \end{aligned}$$

For obvious reason, we have $\bar{\Phi}_{i,\ell}^{(t)} = A_{i,\ell}^{(t)} + B_{i,\ell}^{(t)}$.

Bound the B sequence. Applying Claim D.7, we have

$$B_{i,\ell}^{(t+1)} = B_{i,\ell}^{(t)} + \eta \sum_{r \in \mathcal{M}_i^{(0)}} \Delta_{i,r,\ell}^{(t)} \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbb{1}_{(X,y) \in \mathcal{Z}_{s,i,\ell}} \cdot (V_{i,r,\ell}(X) \pm (\mathcal{E}_1 + \mathcal{E}_3)) \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right]$$

but according to the definition of $\mathcal{Z}_{s,i,\ell}$ (which implies $\widehat{\ell}(X) = 3 - \ell$), we must have

$$0 \leq V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widehat{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) z_p \leq O(\rho) \cdot \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widehat{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle)$$

This means, we must have

$$\begin{aligned} & |B_{i,\ell}^{(t+1)} - B_{i,\ell}^{(t)}| \\ & \leq O\left(\frac{\eta\rho N_s}{N}\right) \sum_{r \in \mathcal{M}_i^{(0)}} \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbb{1}_{(X,y) \in \mathcal{Z}_{s,i,\ell}} \left(1 - \mathbf{logit}_i(F^{(t)}, X)\right) \left(\mathcal{E}_1 + \mathcal{E}_3 + \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle)\right) \right] \end{aligned}$$

Applying Claim D.12a (and noticing $T_0 \leq O(\frac{N}{\eta})$), we conclude that (using our parameter assumption)

$$\forall t \geq 0: \quad |B_{i,\ell}^{(t)}| \leq \tilde{O}\left(\frac{\rho N_s}{k}\right) < \frac{1}{\text{polylog}(k)}$$

Bound the A sequence. So far we have derived that $\Phi_{i,\ell}^{(t)} = A_{i,\ell}^{(t)} \pm \frac{1}{\text{polylog}(k)}$ (because $|B_{i,\ell}^{(t)}|, |\Phi_{i,\ell}^{(t)} - \bar{\Phi}_{i,\ell}^{(t)}| \leq \frac{1}{\text{polylog}(k)}$). So, it suffices to bound the A sequence. Let us denote by

$$\Phi^{(t)} \stackrel{\text{def}}{=} \max_{i \in [k], \ell \in [2]} \Phi_{i,\ell}^{(t)}$$

in the remainder of the proof. Suppose we are now at some iteration $t \geq T_0$, let

$$(i, \ell) = \arg \max_{i \in [k], \ell \in [2]} \{A_{i,\ell}^{(t)}\}$$

Applying Claim D.7 and Claim D.8, and using $V_{i,r,\ell}(X) \in [0, 1]$, we have

$$A_{i,\ell}^{(t+1)} \leq A_{i,\ell}^{(t)} + O(\eta) \left(\mathbb{E}_{(X,y) \sim \mathcal{Z}} [\mathbb{1}_{y=i} \mathbb{1}_{(X,y) \notin \mathcal{Z}_{s,i,\ell}} (1 - \mathbf{logit}_y(F^{(t)}, X))] + \tilde{O}(\sigma_p P) \right)$$

Observe that, whenever $A_{i,\ell}^{(t)} > \text{polylog}(k)$, we also have $\Phi^{(t)} \geq \text{polylog}(k)$ and thus

- for every $(X, y) \in \mathcal{Z}_m$ with $y = i$, recall from Claim D.9 that

$$F_j^{(t)}(X) = \sum_{\ell \in [2]} \left(\Phi_{j,\ell}^{(t)} \times \mathbb{1}_{v_{j,\ell} \in \mathcal{V}(X)} \left(\sum_{p \in \mathcal{P}_{v_{j,\ell}}(X)} z_p \right) \right) \pm O\left(\frac{1}{\text{polylog}(k)}\right)$$

By our choice of the distribution, this implies

- $F_j^{(t)}(X) \leq 0.8001\Phi^{(t)}$ for $j \neq i$, and
- $F_i^{(t)}(X) \geq 0.9999\Phi^{(t)}$ because (i, ℓ) is the argmax of $A_{i,\ell}^{(t)}$ and $A_{i,\ell}^{(t)}$ is close to $\Phi_{i,\ell}^{(t)}$.

- for every $(X, y) \in \mathcal{Z}_s$ with $y = i$ and $\hat{\ell}(X) = \ell$, we can also use Claim D.9 to derive

$$- F_j^{(t)}(X) \leq 0.8001\Phi^{(t)} \text{ for } j \neq i.$$

However, on the i -th output, for every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, using $\langle w_{i,r}^{(t)}, x_p \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p + \langle w_{i,r}^{(t)}, \xi_p \rangle \pm \tilde{O}(\sigma_0 \gamma k)$ from Induction Hypothesis C.3d and using $\langle w_{i,r}^{(t)}, \xi_p \rangle \geq -\frac{1}{\text{polylog}(k)}$ from Claim D.13a, and $\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \geq 1$ from Def. 3.1, we also have

$$- F_i^{(t)}(X) \geq \Phi_{i,\ell}^{(t)} - 0.0001 \geq 0.9999\Phi^{(t)}.$$

In both cases, we have $1 - \mathbf{logit}_i(F^{(t)}, X) = e^{-\Omega(\log^5 k)}$ which is negligible. Therefore, we derived that

$$\text{whenever } \max_{i \in [k], \ell \in [2]} \{A_{i,\ell}^{(t)}\} \geq \text{polylog}(k) \quad \max_{i \in [k], \ell \in [2]} \{A_{i,\ell}^{(t+1)}\} \leq \max_{i \in [k], \ell \in [2]} \{A_{i,\ell}^{(t)}\} + \tilde{O}(\eta \cdot e^{-\Omega(\log^5 k)} + \eta \sigma_p P) .$$

This finishes the proof that all $A_{i,\ell}^{(t)} \leq \tilde{O}(1)$. □

D.4.2 Off-Diagonal Correlations are Small

Our previous subsection upper bounds the “diagonal” correlations $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle$, and in this subsection we bound the “off-diagonal” correlations $\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle$ for $i \neq j$.

Lemma D.22. *Suppose Parameter D.1 holds and suppose Induction Hypothesis C.3 holds for all iterations $< t$. Then,*

$$\forall i \in [k], \forall r \in [m], \forall j \in [k] \setminus \{i\}: \quad |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| \leq \tilde{O}(\sigma_0)$$

Proof of Lemma D.22. Let us denote by $R_i^{(t)} \stackrel{\text{def}}{=} \max_{r \in [m], j \in [k] \setminus \{i\}} |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle|$.

By Claim D.7 and Claim D.8,

$$\begin{aligned} |\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| &\leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbf{1}_{y=i} \left(\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X) \right) \left(1 - \mathbf{logit}_i \left(F^{(t)}, X \right) \right) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X) \right) \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \end{aligned}$$

Stage 1. In the first stage, namely when $t \leq T_{0,i}$, we have $\mathbf{logit}_i \left(F^{(t)}, X \right) \leq O\left(\frac{1}{k}\right)$ (see Claim D.4), have $\mathcal{E}_{2,i,r}(X) \leq \tilde{O}(\gamma(\Lambda_i^{(t)})^{q-1})$, and have $\mathcal{E}_{4,i,r}(X) \leq \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)}$, so

$$|\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| \leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + \tilde{O}\left(\frac{\eta}{k}\right) \left(\gamma(\Lambda_i^{(t)})^{q-1} + (\sigma_0^{q-1})\gamma s + \tilde{O}((\sigma_0 \gamma k)^{q-1}) \gamma P + (\sigma_0)^{q-1} \frac{s}{k} \right) \quad (\text{D.12})$$

Recall in the first stage $\Lambda_i^{(t+1)} = \Lambda_i^{(t)} + \Theta\left(\frac{\eta}{k}\right) \widetilde{\text{ReLU}}'(\Lambda_i^{(t)})$ (see Claim D.10) and therefore

$$\sum_{t \leq T_{0,i}} \eta(\Lambda_i^{(t)})^{q-1} \leq \tilde{O}(k)$$

Hence, together with $T_{0,i} \leq T_0 = \Theta\left(\frac{k}{\eta \sigma_0^{q-2}}\right)$ (see Claim D.11), as long as

$$\gamma = \tilde{O}(\sigma_0), \quad \gamma = \tilde{O}(1/s), \quad (\gamma k)^{q-1} \gamma P = \tilde{O}(1) \quad (\text{D.13})$$

(Note that all of them are satisfied by Parameter D.1). Plugging them back to (D.12), we have we have

$$R_i^{(t)} \leq R_i^{(0)} + \tilde{O}(\sigma_0) + \tilde{O}\left(\frac{\eta}{k} T_0\right) \left((\sigma_0^{q-1})\gamma s + \tilde{O}((\sigma_0 \gamma k)^{q-1}) \gamma P + (\sigma_0)^{q-1} \frac{s}{k} \right) \leq \tilde{O}(\sigma_0)$$

for every $t \leq T_{0,i}$.

Stage 2. In the second stage, namely when $t \in [T_{0,i}, T_0]$, we have the naive upper bounds $\mathcal{E}_{2,i,r}(X) \leq O(\gamma)$ and $\mathcal{E}_{4,i,r}(X) \leq \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)}$, so

$$\begin{aligned} |\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| &\leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbf{1}_{y=i} \left(O(\gamma) + \mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right) \left(1 - \mathbf{logit}_i \left(F^{(t)}, X \right) \right) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right) \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \end{aligned}$$

Recall that for every $(X, y) \in \mathcal{Z}_s$ and $i \neq y$, it satisfies $\mathbf{logit}_i(F^{(t)}, X) = O\left(\frac{1}{k}\right) (1 - \mathbf{logit}_y(F^{(t)}, X))$ (see Claim D.4). Therefore,

$$\begin{aligned} |\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| &\leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbf{1}_{y=i} \left(O(\gamma) + \mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right) \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \\ &\quad + O\left(\frac{\eta N_s}{kN}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[k \mathbf{1}_{y=i} \left(O(\gamma) + \mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \end{aligned}$$

$$+ \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \Big] \quad (\text{D.14})$$

Next, one can naively verify

$$\mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbf{1}_{y=i} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right] \leq \tilde{O}\left(\frac{s}{k^2}\right) \quad \text{and} \quad \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbf{1}_{y \neq i} \mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \right] \leq \tilde{O}\left(\frac{s}{k}\right)$$

Furthermore, by Claim D.17, we have

$$\begin{aligned} \forall t \in [T_{0,i}, T_0]: \quad & \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbf{1}_{y \neq i} \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \leq O\left(\frac{1}{k}\right) \\ \forall t \in [T_{0,i}, T_0]: \quad & \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbf{1}_{y \neq i} \mathbf{1}_{v_{j,\ell} \in \mathcal{P}(X)} \mathbf{logit}_i \left(F^{(t)}, X \right) \right] \leq O\left(\frac{s}{k^2}\right) \\ & \sum_{t=T_{0,i}}^{T_0} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbf{1}_{y=i} \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \leq O\left(\frac{s}{k} T_0 \Upsilon\right) + \tilde{O}\left(\frac{1}{\eta}\right) \end{aligned}$$

Putting these back to (D.14), we immediately conclude (using $N_s \ll N$) that

$$\begin{aligned} R_i^{(t)} &\leq R_i^{(T_{0,i})} + \tilde{O}\left(\eta T_0 \cdot \frac{s}{k^2} \sigma_0^{q-1}\right) \\ &+ \eta \left(O\left(\frac{s}{k} T_0 \Upsilon\right) + \tilde{O}\left(\frac{1}{\eta}\right) \right) \tilde{O}\left(\gamma + (\sigma_0^{q-1}) \gamma s + \tilde{O}\left((\sigma_0 \gamma k)^{q-1}\right) \gamma P + (\sigma_0)^{q-1} \frac{s}{k}\right) \end{aligned}$$

Hence, recalling that $T_0 = \Theta\left(\frac{k}{\eta \sigma_0^{q-2}}\right)$ (see Claim D.11), as long as (D.13) together with $s \leq k$ and¹⁸

$$\frac{s}{k} \frac{k \Upsilon \gamma}{\sigma_0^{q-2}} = \tilde{O}(\sigma_0) \quad , \quad s \Upsilon \leq O(1) \quad (\text{D.15})$$

Then, we also have $R_i^{(t)} \leq \tilde{O}(\sigma_0)$ for every $t \in [T_{0,i}, T_i]$.

Stage 3. At the third stage, namely when $t \geq T_0$, we can continue from (D.14) (and apply the naive bound $\mathbf{1}_{y=i} \leq 1$ and $\mathbf{1}_{v_{j,\ell} \in \mathcal{V}(X)} \leq 1$) to derive that

$$\begin{aligned} |\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| &\leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\left(O(\gamma) + \mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \\ &+ O\left(\frac{\eta N_s}{kN}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[k \mathbf{1}_{y=i} \left(O(\gamma) + \mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \\ &+ \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0)^{q-1} \right) \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \Big] \quad (\text{D.16}) \end{aligned}$$

If we denote by

$$\begin{aligned} S^{(t)} &\stackrel{\text{def}}{=} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right] \\ G_i^{(t)} &\stackrel{\text{def}}{=} \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbf{1}_{i=y} \cdot \left(1 - \mathbf{logit}_y \left(F^{(t)}, X \right) \right) \right] \end{aligned}$$

then we can simplify (D.16) into

$$R_i^{(t+1)} \leq R_i^{(t)} + \eta \left(S^{(t)} + \frac{N_s}{N} G_i^{(t)} + \frac{N_s}{kN} \sum_{i' \in [k]} G_{i'}^{(t)} \right) \tilde{O}\left(\gamma + \sigma_0^{q-1} \gamma s + (\sigma_0 \gamma k)^{q-1} \gamma P + (\sigma_0)^{q-1}\right) \quad . \quad (\text{D.17})$$

¹⁸Note that we have required these parameter choices in Parameter D.1. This is the tightest place for our required upper bounds on γ and on s .

Recall

$$\text{Claim D.14} \implies \sum_{t \geq T_0} S^{(t)} \leq \tilde{O}\left(\frac{k}{\eta}\right)$$

$$\text{Claim D.12} \implies \forall i' \in [k]: \sum_{t \geq T_0} G_i^{(t)} \leq \tilde{O}\left(\frac{N}{\eta k \rho^{q-1}}\right)$$

Putting them back to (D.17), we know that as long as¹⁹

$$k\gamma = \tilde{O}(\sigma_0), \quad \gamma = \tilde{O}(\sigma_0^{q-1}), \quad \frac{N_s}{k\rho^{q-1}} \leq \tilde{O}\left(\frac{1}{\sigma_0^{q-2}}\right) \quad (\text{D.18})$$

it satisfies $R_i^{(t)} \leq \tilde{O}(\sigma_0)$ for all $t \geq T_0$. \square

D.4.3 View Lottery Winning

Recall $\Lambda_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \max_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$. Also recall

$$\mathcal{M} \stackrel{\text{def}}{=} \left\{ (i, \ell^*) \in [k] \times [2] \mid \Lambda_{i,\ell^*}^{(0)} \geq \Lambda_{i,3-\ell^*}^{(0)} \left(\frac{S_{i,3-\ell^*}}{S_{i,\ell^*}} \right)^{\frac{1}{q-2}} + \frac{1}{\text{polylog}(k)} \right\}$$

in which $S_{i,\ell} = \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y=i} \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p^q \right]$.

Our next lemma shows that when $(i, \ell^*) \in \mathcal{M}$, we have $\Lambda_{i,3-\ell^*}^{(t)} \leq \tilde{O}(\sigma_0)$. (In other words, view ℓ wins the lottery and view $3 - \ell$ is negligible, on label i .)

Lemma D.23. *Suppose Parameter D.1 holds and suppose Induction Hypothesis C.3 holds for all iterations $< t$. Then,*

$$\forall i, \ell^* \in \mathcal{M}: \quad \Lambda_{i,3-\ell^*}^{(t)} = \max_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,3-\ell^*} \rangle]^+ \leq \tilde{O}(\sigma_0)$$

Proof of Lemma D.23. By Claim D.7 and Claim D.8,

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &= \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbb{1}_{y=i} (V_{i,r,\ell}(X) \pm O(\mathcal{E}_1 + \mathcal{E}_3)) \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right. \\ &\quad \left. \pm O(1) \cdot \mathbb{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbb{1}_{v_{i,\ell} \in \mathcal{P}(X)} V_{i,r,\ell}(X) \right) \mathbf{logit}_i(F^{(t)}, X) \right] \end{aligned} \quad (\text{D.19})$$

where

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) z_p$$

Stage 1. In the first stage, namely when $t \leq T_{0,i}$, we have $\Lambda_i^{(t)} \leq \Lambda_{\bar{\varnothing}} \ll \tilde{O}(\frac{1}{m_0})$ (see Claim D.11).

This implies $\mathbf{logit}_i(F^{(t)}, X) \leq O(\frac{1}{k})$ from Claim D.4. Also, we have $\Pr[v_{i,\ell} \in \mathcal{P}(X) \mid y \neq i] = \frac{s}{k}$.

Thus, we can simplify (D.19) as and thus

$$\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbb{1}_{y=i} V_{i,r,\ell}(X) \left(1 - O\left(\frac{1}{k}\right) \right) \pm O\left(\frac{1}{k}\right) \mathbb{1}_{i \neq y} \mathbb{1}_{v_{i,\ell} \in \mathcal{P}(X)} V_{i,r,\ell}(X) \pm O\left(\frac{\mathcal{E}_1 + \mathcal{E}_3}{k}\right) \right]$$

Since $N_s \ll N$, we can ignore single-view data and only focus on $(X, y) \in \mathcal{Z}_m$. Using Induction Hypothesis C.3, we know for $(X, y) \in \mathcal{Z}_m$,

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \pm \tilde{o}(\sigma_0) \right) z_p$$

¹⁹Note that we have required these parameter choices in Parameter D.1. This is the tightest place for our required upper bounds on γ and on N_s .

and furthermore since we are in stage 1, it satisfies $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \leq O(\Lambda_{\varnothing}^-) \ll \varrho$ (see Claim D.11) so

$$V_{i,r,\ell}(X) = \frac{1}{\varrho^{q-1}} ([\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+)^{q-1} \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p^q \pm \tilde{O}(\sigma_0)$$

Therefore,

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &= \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \frac{\eta}{\varrho^{q-1}} ([\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+)^{q-1} \left(\left(1 - O\left(\frac{1}{k}\right)\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{i=y} \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p^q \right] \pm O\left(\frac{s}{k^2}\right) \right) \\ &\quad \pm O\left(\frac{\mathcal{E}_1 + \mathcal{E}_3}{k}\right) \cdot \eta \\ &= \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \frac{\eta}{\varrho^{q-1}} ([\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+)^{q-1} \left(1 - O\left(\frac{1}{\text{polylog}(k)}\right)\right) S_{i,\ell} \pm \tilde{o}\left(\frac{\sigma_0}{k}\eta\right) \end{aligned} \quad (\text{D.20})$$

This means, if we take $r^* = \arg \max_{r \in [m]} \{\langle w_{i,r}^{(0)}, v_{i,\ell^*} \rangle\}$ and an arbitrary $r \in [m]$, we can define

- $x_t = \langle w_{i,r^*}^{(t)}, v_{i,\ell^*} \rangle \cdot (S_{i,\ell^*} / \varrho^{q-1})^{\frac{1}{q-2}}$
- $y_t = \max \{\langle w_{i,r}^{(t)}, v_{i,3-\ell^*} \rangle, \sigma_0\} \cdot (S_{i,3-\ell^*} / \varrho^{q-1})^{\frac{1}{q-2}}$

Then by (D.20) we know

$$x_{t+1} \geq x_t + \eta C_t x_t^{q-1} \quad \text{and} \quad y_{t+1} \leq y_t + \eta S C_t y_t^{q-1}$$

for some $C_t = 1 - O\left(\frac{1}{\text{polylog}(k)}\right) \in [0.9, 1]$ (where the value in the big O notion can vary per iteration) and constant $S = \frac{S_{i,3-\ell^*}}{S_{i,\ell^*}} \left(1 + \frac{1}{\text{polylog}(k)}\right)$ that does not depend on t .

Now, since $\Lambda_{i,\ell^*}^{(0)} \geq \Lambda_{i,3-\ell^*}^{(0)} \left(\frac{S_{i,3-\ell^*}}{S_{i,\ell^*}}\right)^{\frac{1}{q-2}} + \frac{1}{\text{polylog}(k)}$ implies $x_0 \geq y_0 S^{\frac{1}{q-2}} \left(1 + \frac{1}{\text{polylog}(k)}\right)$ and we can apply Lemma D.19 to derive that

- when $\langle w_{i,r^*}^{(t)}, v_{i,\ell^*} \rangle$ reaches $\tilde{\Omega}(1)$, which necessarily is an iteration $t \geq T_{0,i}$, we still have that

$$y_t \leq \tilde{O}(y_0) \implies \langle w_{i,r}^{(t)}, v_{i,3-\ell^*} \rangle \leq \tilde{O}(\sigma_0)$$

(This uses $|\langle w_{i,r}^{(0)}, v_{i,3-\ell^*} \rangle| \leq \tilde{O}(\sigma_0)$.) Therefore, we finished the proof that for every $t \leq T_{0,i}$, $\Lambda_{i,3-\ell}^{(t)} \leq \tilde{O}(\sigma_0)$.

Stage 2. In the second stage, namely when $t \in [T_{0,i}, T_i]$, let us denote $\ell = 3 - \ell^*$ for abbreviation. Suppose we want to prove $\Lambda_{i,3-\ell^*}^{(t+1)} \leq \tilde{O}(\sigma_0)$. By Claim D.7 and Claim D.8 again (but this time only using the upper bound part),

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &= \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbb{1}_{y=i} (V_{i,r,\ell}(X) \pm O(\mathcal{E}_1 + \mathcal{E}_3)) \left(1 - \mathbf{logit}_i(F^{(t)}, X)\right) \right. \\ &\quad \left. \pm \mathbb{1}_{y \neq i} (\mathcal{E}_1 + \mathcal{E}_3 + \mathbb{1}_{v_{i,\ell} \in \mathcal{P}(X)} V_{i,r,\ell}(X)) \mathbf{logit}_i(F^{(t)}, X) \right] \end{aligned} \quad (\text{D.21})$$

where

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) z_p$$

- Using Induction Hypothesis C.3 and $\Lambda_{i,3-\ell^*}^{(t)} \leq \tilde{O}(\sigma_0)$, we know for $(X, y) \in \mathcal{Z}_m$, and for

$(X, y) \in \mathcal{Z}_s$ but $y \neq i$,

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \pm \tilde{o}(\sigma_0) \right) z_p \leq \tilde{O}(\sigma_0^{q-1})$$

- Otherwise for $(X, y) \sim \mathcal{Z}_s$ and $y = i$, we can use $\ell = 3 - \ell^*$ and Induction Hypothesis C.3e to derive

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}' \left(\langle w_{i,r}^{(t)}, x_p \rangle \right) z_p \leq \tilde{O}(\sigma_0^{q-1}) \quad (\text{D.22})$$

- For $(X, y) \sim \mathcal{Z}_s$ and $y \neq i$, we have $\mathbf{logit}_i(F^{(t)}, X) \leq O\left(\frac{1}{k}\right)(1 - \mathbf{logit}_y(F^{(t)}, X))$ (see Claim D.4)

Putting these back to (D.21), we have

$$\begin{aligned} |\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle| &\leq |\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle| + O(\eta) \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbf{1}_{y=i} \left(\tilde{O}(\sigma_0^{q-1}) + O(\mathcal{E}_1 + \mathcal{E}_3) \right) \left(1 - \mathbf{logit}_y(F^{(t)}, X) \right) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} \tilde{O}(\sigma_0^{q-1}) \mathbf{logit}_i(F^{(t)}, X) \right] \\ &\quad + O\left(\frac{\eta N_s}{N}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbf{1}_{y=i} \tilde{O}(\sigma_0^{q-1}) \cdot \left(1 - \mathbf{logit}_y(F^{(t)}, X) \right) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} \frac{\mathcal{E}_1 + \mathcal{E}_3 + \tilde{O}(\sigma_0^{q-1})}{k} \left(1 - \mathbf{logit}_y(F^{(t)}, X) \right) \right] \end{aligned} \quad (\text{D.23})$$

Since $N_s \ll N$, we can ignore single-view data and only focus on $(X, y) \in \mathcal{Z}_m$. Applying Claim D.17, we can conclude that

$$|\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle| \leq |\langle w_{i,r}^{(T_0,i)}, v_{i,\ell} \rangle| + \eta O\left(\frac{s}{k} T_0 \Upsilon + \frac{T_0}{k}\right) \cdot \tilde{O}(\sigma_0^{q-1}) + \tilde{O}(\sigma_0)$$

Recall that $T_0 = \Theta\left(\frac{k}{\eta \sigma_0^{q-2}}\right)$ from Claim D.11, and recall we have $s\Upsilon \leq O(1)$ (from (D.15)), we finish the proof that $|\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle| \leq \tilde{O}(\sigma_0)$. This means, $\Lambda_{i,3-\ell^*}^{(t)} \leq \tilde{O}(\sigma_0)$ for every $t \in [T_0, T_i]$.

Stage 3. At the third stage, namely when $t \geq T_0$, we can continue from (D.23) but this time we do not ignore single-view data, and apply the naive bound $\mathbf{logit}_i(F^{(t)}, X) \leq 1 - \mathbf{logit}_y(F^{(t)}, X)$ for $(X, y) \in \mathcal{Z}_m$ and $i \neq y$. Recall again we denote $\ell = 3 - \ell^*$ for notational simplicity. If we denote by

$$\begin{aligned} S^{(t)} &\stackrel{\text{def}}{=} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [1 - \mathbf{logit}_y(F^{(t)}, X)] \\ G_i^{(t)} &\stackrel{\text{def}}{=} \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} [\mathbf{1}_{i=y} \cdot (1 - \mathbf{logit}_y(F^{(t)}, X))] \end{aligned}$$

then we have

$$\Lambda_{i,3-\ell^*}^{(t+1)} \leq \Lambda_{i,3-\ell^*}^{(t)} + \tilde{O}\left(\eta S^{(t)} \sigma_0^{q-1}\right) + O\left(\frac{\eta N_s}{N}\right) \cdot \left(G_i^{(t)} + \frac{\sum_{j \in [k]} G_j^{(t)}}{k}\right) \cdot \tilde{O}(\sigma_0^{q-1}) \quad (\text{D.24})$$

Recall

$$\text{Claim D.14} \implies \sum_{t \geq T_0} S^{(t)} \leq \tilde{O}\left(\frac{k}{\eta}\right)$$

$$\text{Claim D.12} \implies \forall i' \in [k]: \sum_{t \geq T_0} G_{i'}^{(t)} \leq \tilde{O}\left(\frac{N}{\eta k \rho^{q-1}}\right)$$

Putting them back to (D.24), we know that as long as $\frac{N_s}{k \rho^{q-1}} \leq \tilde{O}\left(\frac{1}{\sigma_0^{q-2}}\right)$ (already satisfied in (D.18))

and additionally²⁰

$$k\sigma_0^{q-1} = \tilde{O}(\sigma_0)$$

we have $\Lambda_{i,3-\ell^*}^{(t+1)} \leq \tilde{O}(\sigma_0)$. □

D.4.4 Neuron Lottery Winning: $\mathcal{M}_i^{(0)}$

In this subsection, we prove that the neurons outside $\mathcal{M}_i^{(0)}$ for each index $i \in [k]$ is negligible. (They did not win the lottery so only the neurons in $\mathcal{M}_i^{(0)}$ count.)

Lemma D.24. *Suppose Parameter D.1 holds and suppose Induction Hypothesis C.3 holds for all iterations $< t$. Then,*

$$\forall i \in [k], \forall \ell \in [2], \forall r \in [m] \setminus \mathcal{M}_i^{(0)} : \quad \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \leq \tilde{O}(\sigma_0)$$

Proof of Lemma D.24. The proof is nearly identical to that of Lemma D.23.

Stage 1. In the first stage, namely when $t \leq T_{0,i}$, the proof is nearly identical to stage 1 of the proof of Lemma D.23. At a high level, this time we instead compare two sequences $x_t = \langle w_{i,r^*}^{(t)}, v_{i,\ell} \rangle$ and $y_t = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle$ for $r^* = \arg \max_{r \in [m]} \langle w_{i,r}^{(0)}, v_{i,\ell} \rangle$ and $r \in \mathcal{M}_{i,\ell}^{(0)}$. We skip the details for the sake of cleanness.

Stage 2. In the second stage, namely when $t \in [T_{0,i}, T_0]$, we can also completely reuse the stage 2 of the proof of Lemma D.23. The only slight difference is that in order to derive (D.22), this time we instead use Induction Hypothesis C.3f.

Stage 3. In the third stage, namely when $t \geq T_0$, we can also completely reuse the stage 3 of the proof of Lemma D.23. The only slight difference is that in order to derive (D.22), this time we instead use Induction Hypothesis C.3f. □

D.4.5 Noise Correlation is Small

In this subsection, we prove that the neurons correlate negligibly with the random noise, that is $\langle w_{i,r}^{(t)}, \xi_p \rangle$ is small, except for those single-view data on the lottery winning views. (Recall single-view data are learned through *memorization*, so the learner network can correlate with their noise ξ_p significantly.)

Lemma D.25. *Suppose Parameter D.1 holds and suppose Induction Hypothesis C.3 holds for all iterations $< t$.²¹ For every $\ell \in [2]$, for every $r \in [m]$, for every $(X, y) \in \mathcal{Z}_m$ and $i \in [k]$, or for every $(X, y) \in \mathcal{Z}_s$ and $i \in [k] \setminus \{y\}$:*

(a) *For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, we have: $\langle w_{i,r}^{(t)}, \xi_p \rangle \leq \tilde{o}(\sigma_0)$.*

(b) *For every $p \in \mathcal{P}(X) \setminus (\mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X))$, we have: $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$.*

(c) *For every $p \in [P] \setminus \mathcal{P}(X)$, we have: $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0 \gamma k)$.*

In addition, for every $(X, y) \in \mathcal{Z}_s$, every $i \in [k]$, every $r \in [m]$, every $\ell \in [2]$,

(d) *For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, if $(i, 3 - \ell) \in \mathcal{M}$ we have: $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$.*

²⁰This is the place in our proof that we require $\sigma_0^{q-2} \leq \frac{1}{k}$. For simplicity, we have chosen $\sigma_0^{q-2} \leq \frac{1}{k}$ in Parameter D.1.

²¹In particular, we need to ensure $(\sigma_0)^{q-2} \leq \rho^{q-1}$ and $N \geq \tilde{\Omega}(\frac{k^5}{\sigma_0^{q-1}})$ in the proof of this lemma.

(e) For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, if $r \in [m] \setminus \mathcal{M}_i^{(0)}$ we have: $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$.

We prove Lemma D.25 after we establish the following claim.

Claim D.26. For every $(X, y) \in \mathcal{Z}$, $i \in [k]$, $r \in [m]$, and $p \in [P]$, suppose it satisfies $|\langle w_{i,r}^{(t)}, x_p \rangle| \leq A$ for every $t < t_0$ where t_0 is any iteration $t_0 \leq T$. Then,

- If $(X, y) \in \mathcal{Z}_m$ then $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}\left(\frac{kA^{q-1}}{N\sigma_0^{q-2}} + \frac{k^5A^{q-1}}{s^2N}\right) + \frac{\eta T}{\sqrt{d}}$
- If $(X, y) \in \mathcal{Z}_s$ then $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}\left(\frac{kA^{q-1}}{N\sigma_0^{q-2}} + \frac{A^{q-1}}{\rho^{q-1}}\right) + \frac{\eta T}{\sqrt{d}}$

We first prove Lemma D.25 using Claim D.26 (which is trivial), and then we prove Claim D.26.

Proof of Lemma D.25.

- In the case of Lemma D.25a, since we have $|\langle w_{i,r}^{(t')}, x_p \rangle| \leq \tilde{O}(1)$ for every $t' < t$ according to our Induction Hypothesis C.3a and C.3g, by applying Claim D.26 we immediately have $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ once we plug in $A = \tilde{O}(1)$, $N \geq \tilde{\omega}\left(\frac{k}{\sigma_0^{q-1}}\right)$, and $N \geq \tilde{\omega}\left(\frac{k^5}{\sigma_0}\right)$.
- In the case of Lemma D.25b, since we have $|\langle w_{i,r}^{(t')}, x_p \rangle| \leq \tilde{O}(\sigma_0)$ for every $t' < t$ according to our Induction Hypothesis C.3b, by applying Claim D.26 we immediately have $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ once we plug in $A = \tilde{O}(\sigma_0)$, $N \geq k^5$ and $(\sigma_0)^{q-2} \leq \rho^{q-1}$.
- In the case of Lemma D.25c, since we have $|\langle w_{i,r}^{(t')}, x_p \rangle| \leq \tilde{O}(\sigma_0\gamma k)$ for every $t' < t$ according to our Induction Hypothesis C.3c, by applying Claim D.26 we immediately have $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0\gamma k)$ once we plug in $A = \tilde{O}(\sigma_0\gamma k)$, $N \geq k^5$ and $(\sigma_0\gamma k)^{q-2} \leq \rho^{q-1}$.
- In the case of Lemma D.25d, since we have $|\langle w_{i,r}^{(t')}, x_p \rangle| \leq \tilde{O}(\sigma_0)$ for every $t' < t$ according to our Induction Hypothesis C.3e, by applying Claim D.26 we immediately have $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ once we plug in $A = \tilde{O}(\sigma_0)$, $N \geq k^5$ and $(\sigma_0)^{q-2} \leq \rho^{q-1}$.
- In the case of Lemma D.25e, since we have $|\langle w_{i,r}^{(t')}, x_p \rangle| \leq \tilde{O}(\sigma_0)$ for every $t' < t$ according to our Induction Hypothesis C.3f, by applying Claim D.26 we immediately have $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ once we plug in $A = \tilde{O}(\sigma_0)$, $N \geq k^5$ and $(\sigma_0)^{q-2} \leq \rho^{q-1}$.

□

Proof of Claim D.26. Recall from our earlier calculation (see (D.6)) that for every $(X, y) \in \mathcal{Z}$ and $p \in [P]$, if $y = i$ then

$$\langle w_{i,r}^{(t+1)}, \xi_p \rangle = \langle w_{i,r}^{(t)}, \xi_p \rangle + \tilde{\Theta}\left(\frac{\eta}{N}\right) \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) \left(1 - \mathbf{logit}_i(F^{(t)}, X)\right) \pm \frac{\eta}{\sqrt{d}}$$

for similar reason, if $y \neq i$, then

$$\langle w_{i,r}^{(t+1)}, \xi_p \rangle = \langle w_{i,r}^{(t)}, \xi_p \rangle - \tilde{\Theta}\left(\frac{\eta}{N}\right) \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) \mathbf{logit}_i(F^{(t)}, X) \pm \frac{\eta}{\sqrt{d}}$$

Note by our assumption, we have $\widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) \leq \tilde{O}(A^{q-1})$. For obvious reason, when $t = T_0 = \Theta\left(\frac{k}{\eta\sigma_0^{q-2}}\right)$ (see Claim D.11)

$$\left|\langle w_{i,r}^{(t)}, \xi_p \rangle\right| \leq \tilde{O}\left(\frac{\eta}{N}A^{q-1}T_0\right) + \frac{\eta T_0}{\sqrt{d}} \leq \tilde{O}\left(\frac{kA^{q-1}}{N\sigma_0^{q-2}}\right) + \frac{\eta T_0}{\sqrt{d}}$$

Case 1: multi-view data. We first consider $(X, y) \in \mathcal{Z}_m$. Using Claim D.16 and Claim D.14, we have

$$y = i \implies \sum_{t=T_0}^T 1 - \mathbf{logit}_y(F^{(t)}, X) \leq \tilde{O}\left(\frac{k^4}{s^2}\right) \cdot \sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[1 - \mathbf{logit}_y(F^{(t)}, X) \right] \leq \tilde{O}\left(\frac{k^5}{s^2 \eta}\right)$$

$$y \neq i \implies \sum_{t=T_0}^T \mathbf{logit}_i(F^{(t)}, X) \leq \sum_{t=T_0}^T 1 - \mathbf{logit}_y(F^{(t)}, X) \leq \tilde{O}\left(\frac{k^5}{s^2 \eta}\right)$$

Combining this with the bound at $t = T_0$, we have

$$\left| \langle w_{i,r}^{(t)}, \xi_p \rangle \right| \leq \tilde{O}\left(\frac{kA^{q-1}}{N\sigma_0^{q-2}} + \frac{k^5 A^{q-1}}{s^2 N}\right) + \frac{\eta T}{\sqrt{d}}$$

Case 2: single-view data. Let us now consider $(X, y) \in \mathcal{Z}_s$. Recall from Claim D.12b that

$$\sum_{t=T_0}^T \left(1 - \mathbf{logit}_y(F^{(t)}, X) \right) \leq \tilde{O}\left(\frac{N}{\eta \rho^{q-1}}\right)$$

so using the same analysis, we have

$$\left| \langle w_{i,r}^{(t)}, \xi_p \rangle \right| \leq \tilde{O}\left(\frac{kA^{q-1}}{N\sigma_0^{q-2}} + \frac{A^{q-1}}{\rho^{q-1}}\right) + \frac{\eta T}{\sqrt{d}} \quad \square$$

D.4.6 Diagonal Correlations are Nearly Non-Negative

Lemma D.27. *Suppose Parameter D.1 holds and suppose Induction Hypothesis C.3 holds for all iterations $< t$. Then,*

$$\forall i \in [k], \forall r \in [m], \forall \ell \in [2]: \quad \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq -\tilde{O}(\sigma_0).$$

Proof of Lemma D.27. Let us consider any iteration t so that $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \leq -\tilde{\Omega}(\sigma_0)$. We start from this iteration to see how negative the next iterations can be. Without loss of generality we consider the case when $\langle w_{i,r}^{(t')}, v_{i,\ell} \rangle \leq -\tilde{\Omega}(\sigma_0)$ holds for every $t' \geq t$.

By Claim D.7 and Claim D.8,

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbf{1}_{y=i} \left(V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P) \right) \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right. \\ &\quad \left. - \mathbf{1}_{y \neq i} \left(\mathcal{E}_1 + \mathcal{E}_3 + \mathbf{1}_{v_{i,\ell} \in \mathcal{P}(X)} V_{i,r,\ell}(X) \right) \mathbf{logit}_i(F^{(t)}, X) \right] \end{aligned}$$

Recall $V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) z_p$. Since $V_{i,r,\ell}(X) \geq 0$ we can ignore it in the first occurrence corresponding to $\mathbf{1}_{y=i}$. Also, applying Induction Hypothesis C.3, we know that as long as $(X, y) \in \mathcal{Z}$ and $y \neq i$, it satisfies

$$V_{i,r,\ell}(X) = \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \pm \tilde{o}(\sigma_0)) z_p = 0$$

because we have assumed $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \leq -\tilde{\Omega}(\sigma_0)$. Therefore,

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left[\mathbf{1}_{y=i} \tilde{O}(\sigma_p P) \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} (\mathcal{E}_1 + \mathcal{E}_3) \mathbf{logit}_i(F^{(t)}, X) \right] \end{aligned}$$

We first consider every $t \leq T_0 \stackrel{\text{def}}{=} \Theta\left(\frac{k}{\eta\sigma_0^{q-2}}\right)$ (recall Claim D.11). Using Claim D.4 we have $\mathbf{logit}_i(F^{(t)}, X) = O\left(\frac{1}{k}\right)$. This implies

$$\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq -\tilde{O}(\sigma_0) - O\left(\frac{\eta T_0}{k}\right)(\mathcal{E}_1 + \mathcal{E}_3) = -\tilde{O}(\sigma_0) - O\left(\frac{1}{\sigma_0^{q-2}}\right)(\mathcal{E}_1 + \mathcal{E}_3) \geq -\tilde{O}(\sigma_0)$$

(Above, the last inequality uses our earlier parameter choices, see (D.12).)

As for $t \geq T_0$, we combining this with $\mathbf{logit}_i(F^{(t)}, X) \leq 1 - \mathbf{logit}_y(F^{(t)}, X)$ for $i \neq y$ and $(X, y) \in \mathcal{Z}_m$, and $\mathbf{logit}_i(F^{(t)}, X) \leq O\left(\frac{1}{k}\right)(1 - \mathbf{logit}_y(F^{(t)}, X))$ for $i \neq y$ and $(X, y) \in \mathcal{Z}_s$, we have

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[O(\mathcal{E}_1 + \mathcal{E}_3) \left(1 - \mathbf{logit}_y(F^{(t)}, X)\right) \right] \\ &\quad - O\left(\frac{\eta N_s}{N}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbf{1}_{y=i} \tilde{O}(\sigma_p P) \left(1 - \mathbf{logit}_y(F^{(t)}, X)\right) \right. \\ &\quad \left. + \mathbf{1}_{y \neq i} \frac{\mathcal{E}_1 + \mathcal{E}_3}{k} \left(1 - \mathbf{logit}_y(F^{(t)}, X)\right) \right] \end{aligned}$$

Finally, recall

$$\text{Claim D.14} \implies \sum_{t \geq T_0} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [1 - \mathbf{logit}_y(F^{(t)}, X)] \leq \tilde{O}\left(\frac{k}{\eta}\right)$$

$$\text{Claim D.12} \implies \forall i' \in [k]: \sum_{t \geq T_0} \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} [\mathbf{1}_{i=y} \cdot (1 - \mathbf{logit}_y(F^{(t)}, X))] \leq \tilde{O}\left(\frac{N}{\eta k \rho^{q-1}}\right)$$

Therefore, we know for every $t \in [T_0, T]$:

$$\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \langle w_{i,r}^{(T_0)}, v_{i,\ell} \rangle - \tilde{O}(k \cdot (\mathcal{E}_1 + \mathcal{E}_3)) \stackrel{\textcircled{1}}{\geq} \langle w_{i,r}^{(T_0)}, v_{i,\ell} \rangle - \tilde{O}(\sigma_0) \geq -\tilde{O}(\sigma_0)$$

(Above, the inequality $\textcircled{1}$ uses our earlier parameter choices, see (D.18).) □

D.5 Putting All Together

We are now ready to restate Theorem D.2 and prove it.

Theorem D.2. *Under Parameter D.1, for any $m \in [\tilde{\Omega}(1), \tilde{O}(\frac{1}{\sigma_0})]$ and sufficiently small $\eta \leq \frac{1}{\text{poly}(k)}$, our Induction Hypothesis C.3 holds for all iterations $t = 0, 1, \dots, T$.*

Proof of Theorem D.2. At iteration t , we first calculate

$$\forall p \in P_{v_{j,\ell}}(X): \quad \langle w_{i,r}^{(t)}, x_p \rangle = \langle w_{i,r}^{(t)}, v_{j,\ell} \rangle z_p + \sum_{v' \in \mathcal{V}} \alpha_{p,v'} \langle w_{i,r}^{(t)}, v' \rangle + \langle w_{i,r}^{(t)}, \xi_p \rangle \quad (\text{D.25})$$

$$\forall p \in [P] \setminus P(X): \quad \langle w_{i,r}^{(t)}, x_p \rangle = \sum_{v' \in \mathcal{V}} \alpha_{p,v'} \langle w_{i,r}^{(t)}, v' \rangle + \langle w_{i,r}^{(t)}, \xi_p \rangle \quad (\text{D.26})$$

It is easy to verify Induction Hypothesis C.3 holds at iteration $t = 0$ (merely some simple high probability bounds on Gaussian random variables). Suppose Induction Hypothesis C.3 holds for all iterations $< t$. We have established several lemmas

$$\text{Lemma D.22} \implies \forall i \in [k] \forall r \in [m] \forall j \in [k] \setminus \{i\}: \quad |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| \leq \tilde{O}(\sigma_0) \quad (\text{D.27})$$

$$\text{Lemma D.21+Lemma D.27} \implies \forall i \in [k] \forall r \in [m] \forall \ell \in [2]: \quad \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \in [-\tilde{O}(\sigma_0), \tilde{O}(1)] \quad (\text{D.28})$$

$$\text{Lemma D.23+Lemma D.27} \implies \forall (i, \ell^*) \in \mathcal{M} \forall r \in [m]: \quad |\langle w_{i,r}^{(t)}, v_{i,3-\ell^*} \rangle| \leq \tilde{O}(\sigma_0) \quad (\text{D.29})$$

$$\text{Lemma D.27} \implies \forall i \in [k], \forall r \in [m], \forall \ell \in [2]: \quad \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq -\tilde{O}(\sigma_0) \quad (\text{D.30})$$

$$\text{Lemma D.24+Lemma D.27} \implies \forall i \in [k], \forall \ell \in [2], \forall r \in [m] \setminus \mathcal{M}_i^{(0)} : \quad |\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle| \leq \tilde{O}(\sigma_0) \quad (\text{D.31})$$

- To prove C.3a, it suffices to plug (D.27), (D.28) into (D.25), use $\alpha_{p,v'} \in [0, \gamma]$, use $|\mathcal{V}| = 2k$, and use $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ from Lemma D.25a.
- To prove C.3b, it suffices to plug (D.27), (D.28) into (D.25), use $\alpha_{p,v'} \in [0, \gamma]$, use $|\mathcal{V}| = 2k$, and use $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ from Lemma D.25b.
- To prove C.3c, it suffices to plug (D.27), (D.28) into (D.26), use $\alpha_{p,v'} \in [0, \gamma]$, use $|\mathcal{V}| = 2k$, and use $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0 \gamma k)$ from Lemma D.25c.
- To prove C.3d, it suffices to plug (D.27), (D.28) into (D.25), use $\alpha_{p,v'} \in [0, \gamma]$, use $|\mathcal{V}| = 2k$.
- To prove C.3e, it suffices to plug (D.27), (D.29) into (D.25), use $\alpha_{p,v'} \in [0, \gamma]$, use $|\mathcal{V}| = 2k$, and use $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ from Lemma D.25d.
- To prove C.3f, it suffices to plug (D.27), (D.31) into (D.25), use $\alpha_{p,v'} \in [0, \gamma]$, use $|\mathcal{V}| = 2k$, and use $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$ from Lemma D.25e.
- To prove C.3g, it suffices to note that (D.28) exactly implies $\Lambda_i^{(t)} \leq \tilde{O}(1)$, and note that Claim D.10 (which says as long as $\Lambda_i^{(t)} \leq O(1/m_0)$, then it must grow by $\Lambda_i^{(t+1)} \geq \Lambda_i^{(t)} + \Theta(\frac{\eta}{k}) \widetilde{\text{ReLU}}'(\Lambda_i^{(t)})$) implies $\Lambda_i^{(t)} \geq \Omega(\Lambda_i^{(0)}) \geq \tilde{\Omega}(\sigma_0)$.
- To prove C.3h, it suffices to invoke (D.30).
- To prove C.3i, it suffices to invoke (D.31).

□

E Single Model and Ensemble: Theorem Statements

We can now state the general version of the main theorem for single model, as below:

Theorem 1 (single model, restated). *For sufficiently large $k > 0$, every $m \in [\text{polylog}(k), \frac{1}{s\sigma_0^q \text{polylog}(k)}]$, every $\eta \leq \frac{1}{\text{poly}(k)}$, after $T = \frac{\text{poly}(k)}{\eta}$ many iterations, when Parameter D.1 is satisfied, with probability at least $1 - e^{-\Omega(\log^2 k)}$:*

- (training accuracy is perfect) for every $(X, y) \in \mathcal{Z}$:

$$\forall i \neq y: F_y^{(T)}(X) \geq F_i^{(T)}(X) + \Omega(\log k).$$

- (multi-view testing is good) for every $i, j \in [k]$ we have $\tilde{O}(1) \geq \Phi_i^{(T)} \geq 0.4\Phi_j^{(T)} + \Omega(\log k)$, and thus

$$\Pr_{(X,y) \in \mathcal{D}_m} \left[F_y^{(T)}(X) \geq \max_{j \neq y} F_j^{(T)}(X) + \Omega(\log k) \right] \geq 1 - e^{-\Omega(\log^2 k)}$$

- (single-view testing is bad) for every $(i, \ell) \in \mathcal{M}$ we have $\Phi_{i,3-\ell}^{(T)} \leq \tilde{O}(\sigma_0 m) \ll \frac{1}{\text{polylog}(k)}$, and since $|\mathcal{M}| \geq k(1 - o(1))$, we have²²

$$\Pr_{(X,y) \in \mathcal{D}_s} \left[F_y^{(T)}(X) \geq \max_{j \neq y} F_j^{(T)}(X) - \frac{1}{\text{polylog}(k)} \right] \leq \frac{1}{2}(1 + o(1))$$

²²Note that we have assumed for simplicity that there are 2 views with equal probability, and this is why the testing accuracy is close to $\frac{1}{2}$; in more general settings, as we stated in Section 4, this accuracy may be some other constant μ .

We also state the general version of the main theorem for ensemble model, as below:

Theorem 2 (ensemble accuracy, restated). *In the same setting as above, suppose $\{F^{[w]}\}_{w \in [K]}$ are K independently randomly trained models with $m \in [\log^{\Omega(1)}(k), \log^{O(1)} k]$ for $T = \frac{\text{poly}(k)}{\eta}$ iterations each. Let us define $G(X) = \frac{1}{K} \sum_w F^{[w]}(X)$.*

- (training is perfect) same as the single model;
- (multi-view testing is good) same as the single model;
- (single-view testing is good) when $K \geq \text{polylog}(k)$, ensemble model satisfies

$$\Pr_{(X,y) \sim \mathcal{D}_s} \left[G_y(X) \geq \max_{i \in [k] \setminus \{y\}} G_i(X) + \frac{1}{\text{polylog}(k)} \right] \geq 1 - e^{-\Omega(\log^2 k)}$$

E.1 Proof of Theorem 1

Since Theorem D.2 implies the induction hypothesis holds for every $t \leq T$, we have according to Claim D.14 and Claim D.12b that

$$\begin{aligned} \sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [1 - \mathbf{logit}_y(F^{(t)}, X)] &\leq \tilde{O}\left(\frac{k}{\eta}\right) \\ \sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} (1 - \mathbf{logit}_y(F^{(t)}, X)) &\leq \tilde{O}\left(\frac{N}{\eta \rho^{q-1}}\right) \end{aligned} \quad (\text{E.1})$$

Also recall that our training objective is

$$L(F^{(t)}) = \mathbb{E}_{(X,y) \sim \mathcal{Z}} [-\log \mathbf{logit}_y(F^{(t)}, X)]$$

Now, since for every data,

- if $\mathbf{logit}_y(F^{(t)}, X) \geq \frac{1}{2}$ then $-\log \mathbf{logit}_y(F^{(t)}, X) \leq O(1 - \mathbf{logit}_y(F^{(t)}, X))$;
- if $\mathbf{logit}_y(F^{(t)}, X) \leq \frac{1}{2}$, this cannot happen for too many tuples (X, y, t) thanks to (E.1), and when this happens we have a naive bound $-\log \mathbf{logit}_y(F^{(t)}, X) \in [0, \tilde{O}(1)]$ using Claim D.4.

Therefore, we can safely conclude using (E.1) that, when $T \geq \text{poly}(k)/\eta$,

$$\frac{1}{T} \sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}} [-\log \mathbf{logit}_y(F^{(t)}, X)] \leq \frac{1}{\text{poly}(k)}$$

On the other hand, since we are using full gradient descent and the objective function is $O(1)$ -Lipschitz continuous, it means the objective value is monotonically non-increasing. In other words, we have

$$\mathbb{E}_{(X,y) \sim \mathcal{Z}} (1 - \mathbf{logit}_y(F^{(T)}, X)) \leq \mathbb{E}_{(X,y) \sim \mathcal{Z}} [-\log \mathbf{logit}_y(F^{(T)}, X)] \leq \frac{1}{\text{poly}(k)}$$

also for the last iteration T . This immediately implies that the training accuracy is perfect.

As for the multi-view test accuracy, we recall from Claim D.16 that $0.4\Phi_i^{(T)} - \Phi_j^{(T)} \leq -\Omega(\log k)$ for every $i \neq j$. This combined with the function approximation Claim D.9 shows that with high probability $F_y^{(T)}(X) \geq \max_{j \neq y} F_j^{(T)}(X) + \Omega(\log k)$ for every $(X, y) \in \mathcal{D}_m$.

As for the single-view test accuracy, whenever $(i, \ell) \in \mathcal{M}$, using Lemma D.23 we have $\Lambda_{i,3-\ell}^{(T)} \leq \tilde{O}(\sigma_0)$ so $\Phi_{i,3-\ell}^{(T)} \leq \tilde{O}(\sigma_0 m)$.

Now, for every single-view data $(X, y) \in \mathcal{D}_s$ with $y = i$, we know that with half probability $\hat{\ell}(X) = 3 - \ell$. When this happens, according to Claim D.9, we have $F_y^{(T)}(X) \leq O(\rho) + \frac{1}{\text{polylog}(k)}$ (using Def. 3.1 for the single-view distribution).

For every other $j \neq y$, whenever suppose $\ell' = \arg \max_{\ell' \in [2]} \{\Phi_{j,\ell'}^{(T)}\}$, we have $\Phi_{j,\ell'}^{(T)} \geq \Omega(\log k)$ by a few lines above. This means, as long as $v_{j,\ell'} \in \mathcal{V}(X)$ (which happens with probability s/k for

every j), invoking Claim D.9 again, that $F_j^{(T)}(X) \geq \tilde{\Omega}(\Gamma)$. In other words, when this happens for some $j \in [k] \setminus \{i\}$ (which happens with probability at least $1 - e^{-\Omega(\log^2 k)}$), we have

$$F_y^{(T)}(X) \leq \max_{j \neq y} F_j^{(T)}(X) - \frac{1}{\text{polylog}(k)} \quad (\text{E.2})$$

To sum up, we have shown for every $(i, \ell) \in \mathcal{M}$, for every $(X, y) \in \mathcal{D}_s$ with $y = i$, we know that with probability at least $\frac{1}{2}(1 - o(1))$, inequality (E.2) holds. Since the size $|\mathcal{M}| \geq k(1 - o(1))$ (see Proposition C.2), we finish the proof. \blacksquare

E.2 Proof of Theorem 2

Recall from Proposition C.2 that for every $i \in [k]$, $\ell \in [2]$ and every model $F^{[w]}$, the probability for (i, ℓ) to be included in the set $\mathcal{M}^{[w]}$ (defined by model $F^{[w]}$) is at least $m^{-O(1)}$. When this happens, we also have $\Phi_{i,\ell}^{(T)} \geq \Omega(\log k)$ for this model (because $\Phi_i^{(T)} \geq \Omega(\log k)$ while $\Phi_{i,3-\ell}^{(T)} \ll 1$). Let us denote it as $\Phi_{i,\ell}^{[w]}$.

Now, for every $(X, y) \in \mathcal{D}_s$, with probability at least $1 - e^{-\Omega(\log^2 k)}$, letting $\ell = \hat{\ell}(X)$, we have (see Claim D.9)

$$\begin{aligned} \text{for every } F^{[w]} \text{ with } \Phi_{y,\ell}^{[w]} \geq \Omega(\log k) &\implies F_y^{[w]}(X) \geq \Phi_{y,\ell}^{[w]} - \frac{1}{\text{polylog}(k)} \geq \Omega(\log k) \\ \text{for every } F^{[w]} \text{ with } i \neq y &\implies F_i^{[w]}(X) \leq \Gamma(\Phi_{i,1} + \Phi_{i,2}) + \frac{1}{\text{polylog}(k)} \leq O(\Gamma) \end{aligned}$$

Therefore, once we have $K \geq m^{\Omega(1)}$ models in the ensemble, and suppose $\Gamma \leq \frac{1}{m^{\Omega(1)}}$, after taking average, we shall have $G_y(X) \geq G_i(X) + \frac{1}{\text{polylog}(k)}$ for every $i \neq [k]$. \blacksquare

F Knowledge Distillation: Theorem Statement

In this section we show how knowledge distillation (both for ensemble and for self-distillation) can improve the final generalization accuracy. For every i , let us define the *truncated scaled logit* as (for $\tau = \frac{1}{\log^2 k}$):

$$\mathbf{logit}_i^\tau(F, X) = \frac{e^{\min\{\tau^2 F_i(X), 1\}/\tau}}{\sum_{j \in [k]} e^{\min\{\tau^2 F_j(X), 1\}/\tau}}$$

This logit function should be reminiscent of the logit function with temperature used by the seminal knowledge distillation paper by [42]; we use the truncation function instead which is easier to analyze.

F.1 Using Ensemble for Knowledge Distillation

Suppose $\{F^{[i]}\}_{i \in [K]}$ are $K = \tilde{\Theta}(1)$ independently trained models of F for $T = O(\frac{\text{poly}(k)}{\eta})$ iterations (i.e., the same setting as Theorem 1). Let us define their ensemble

$$G(X) = \frac{\Xi}{K} \sum_i F^{[i]}(X) \quad \text{for some } \Xi = \tilde{\Theta}(1) \quad (\text{F.1})$$

Recall from (4.3) that we train a new network F from random initialization, and at every

iteration t , we update each weight $w_{i,r}$ by:

$$w_{i,r}^{(t+1)} = w_{i,r}^{(t)} - \eta \nabla_{w_{i,r}} L(F^{(t)}) - \eta' \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left(\left(\mathbf{logit}_i^\tau(F^{(t)}, X) - \mathbf{logit}_i^\tau(G, X) \right)^- \nabla_{w_{i,r}} F_i^{(t)}(X) \right) \quad (4.3) \text{ restated}$$

Let $F^{(t)}$ be the resulted network obtained by distilling G using algorithm (4.3) at iteration t . We have the following theorem:

Theorem 3 (ensemble distillation, restated). *For sufficiently large $k > 0$, for every $m \in [\log^{\Omega(1)}(k), \log^{O(1)} k]$, every $\eta \leq \frac{1}{\text{poly}(k)}$, setting $\eta' = \eta \text{poly}(k)$, after $T = \frac{\text{poly}(k)}{\eta}$ many iterations, when Parameter G.2 is satisfied, with probability at least $1 - e^{-\Omega(\log^2 k)}$, for at least 90% of the iterations $t \leq T$:*

- (training accuracy is perfect) for every $(X, y) \in \mathcal{Z}$:

$$\forall i \neq y: F_y^{(t)}(X) \geq F_i^{(t)}(X) + \Omega(\log k).$$

- (multi-view testing is good) for every $i, j \in [k]$ we have $\tilde{O}(1) \geq \Phi_i^{(t)} \geq 0.4\Phi_j^{(t)} + \Omega(\log k)$, and thus

$$\Pr_{(X,y) \in \mathcal{D}_m} \left[F_y^{(t)}(X) \geq \max_{j \neq y} F_j^{(t)}(X) + \Omega(\log k) \right] \geq 1 - e^{-\Omega(\log^2 k)}$$

- (single-view testing is good) for every $i \in [k]$ and $\ell \in [2]$ we have $\Phi_{i,\ell}^{(t)} \geq \Omega(\log k)$ and thus

$$\Pr_{(X,y) \in \mathcal{D}_s} \left[F_y^{(T)}(X) \geq \max_{j \neq y} F_j^{(T)}(X) + \Omega(\log k) \right] \leq 1 - e^{-\Omega(\log^2 k)}$$

Our proof to Theorem 3 is in the next Section G.

F.2 Self-Distillation: Using a Single Model to Distill Itself

Recall in the self-distillation case, we made an additional assumption for simplicity:

Assumption 4.1 (balanced \mathcal{D}_m , restated). *In Def. 3.1, for multi-view data (X, y) , we additionally assume that the marginal distributions of $\sum_{p \in \mathcal{P}_v(X)} z_p^q \in [1, O(1)]$ for $v \in \{v_{y,1}, v_{y,2}\}$.*

Let G be a single model trained in the same way as Theorem 1. At the end of training, we scale it up by a small factor $G \leftarrow \log^4 k \cdot G$. Define a “lottery winning” set

$$\mathcal{M}_G \stackrel{\text{def}}{=} \left\{ (i, \ell^*) \in [k] \times [2] \mid \Lambda_{i,\ell^*}^{(0)} \geq \Lambda_{i,3-\ell^*}^{(0)} \left(1 + \frac{2}{\log^2(m)} \right) \right\} \quad (\text{F.2})$$

which only depends on $\Lambda_{i,\ell^*}^{(0)}$ which in terms depends on G ’s random initialization. (Note \mathcal{M}_G is provably a subset of \mathcal{M} defined in (C.2).)

As for F , we break its update into two stages.

1. (Learn.) In the first stage, in the same way as Theorem 1, we start from random initialization and update

$$w_{i,r}^{(t+1)} = w_{i,r}^{(t)} - \eta \nabla_{w_{i,r}} L(F^{(t)}) \text{ for } T = \frac{\text{poly}(k)}{\eta} \text{ iterations}$$

We let \mathcal{M}_F be the “lottery winning” set of network F at the end of stage 1 defined in the same way as (F.2) (which now depends only on F ’s random initialization).

2. (Distill.) In the second stage, for another $T' = \frac{\text{poly}(k)}{\eta}$ iterations, we update

$$w_{i,r}^{(t+1)} = w_{i,r}^{(t)} - \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} \left(\left(\mathbf{logit}_i^\tau(F, X) - \mathbf{logit}_i^\tau(G, X) \right)^- \nabla_{w_{i,r}} F_i^{(t)}(X) \right) \quad (4.4) \text{ restated}$$

Theorem 4 (self distillation, restated). *Suppose the data satisfies Assumption 4.1. For sufficiently large $k > 0$, for every $m \in [\log^{\Omega(1)}(k), k]$, every $\eta \leq \frac{1}{\text{poly}(k)}$, setting $T = \frac{\text{poly}(k)}{\eta}$ and $T' = \frac{\text{poly}(k)}{\eta}$, when Parameter D.1 is satisfied, with probability at least $1 - e^{-\Omega(\log^2 k)}$:*

- (training accuracy is perfect) for every $(X, y) \in \mathcal{Z}$:

$$\forall i \neq y: F_y^{(T+T')}(X) \geq F_i^{(T+T')}(X) + \Omega(\log k).$$

- (multi-view testing is good) for every $i, j \in [k]$ we have $\tilde{O}(1) \geq \Phi_i^{(T+T')} \geq 0.4\Phi_j^{(T+T')} + \Omega(\log k)$, and thus

$$\Pr_{(X,y) \in \mathcal{D}_m} \left[F_y^{(T+T')}(X) \geq \max_{j \neq y} F_j^{(T+T')}(X) + \Omega(\log k) \right] \geq 1 - e^{-\Omega(\log^2 k)}$$

- (single-view testing is better) for every $(i, \ell) \in \mathcal{M}_F \cup \mathcal{M}_G$ we have $\Phi_{i,\ell}^{(T+T')} \geq \Omega(\frac{1}{\log k})$, and since $|\mathcal{M}_F \cup \mathcal{M}_G| \geq 1.5k(1 - o(1))$, we have

$$\Pr_{(X,y) \in \mathcal{D}_s} \left[F_y^{(T+T')}(X) \geq \max_{j \neq y} F_j^{(T+T')}(X) + \Omega(\log k) \right] \geq \frac{3}{4}(1 - o(1))$$

The proof of Theorem 4 is quite easy once the reader is familiar with the proofs of Theorem 1 and Theorem 3. We include it at the end of the next Section G.

G Knowledge Distillation Proof for Ensemble

Our proof structure of Theorem 3 is the same as that for the single model case, but is a lot simpler thanks to our special choice of the truncated distillation function.

Specifically, we maintain the following set of simpler induction hypothesis.

Induction Hypothesis G.1. *For every $\ell \in [2]$, for every $r \in [m]$, for every $(X, y) \in \mathcal{Z}$ and $i \in [k]$,*

(a) *For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, we have: $\langle w_{i,r}^{(t)}, x_p \rangle = \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle z_p \pm \tilde{o}(\sigma_0)$.*

(b) *For every $p \in \mathcal{P}(X) \setminus (\mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X))$, we have: $|\langle w_{i,r}^{(t)}, x_p \rangle| \leq \tilde{O}(\sigma_0)$.*

(c) *For every $p \in [P] \setminus \mathcal{P}(X)$, we have: $|\langle w_{i,r}^{(t)}, x_p \rangle| \leq \tilde{O}(\sigma_0 \gamma k)$.*

Moreover, we have for every $i \in [k]$, every $\ell \in [2]$,

(g) $\Phi_{i,\ell}^{(t)} \geq \Omega(\sigma_0)$ and $\Phi_{i,\ell}^{(t)} \leq \tilde{O}(1)$.

(h) for every $r \in [m]$, it holds that $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq -\tilde{O}(\sigma_0)$.

(Recall $\Phi_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \sum_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$.)

Parameter G.2. The parameter range for our proofs in this section to hold is the same as Parameter D.1, except that

- $N \geq \eta T \cdot \text{poly}(k)$ and $\eta T \geq \text{poly}(k)$. (Instead of $\eta T \geq N \cdot \text{poly}(k)$.)

Explanation: single models need a longer training time because they need to memorize single-view data; instead, here ensemble distillation can truly learn all the training data so the training time T can be shorter.²³

²³Our result also holds for longer T , at the expense of adding an additional lemma. We choose to assume $T \geq N \cdot \text{poly}(k)$ for simplicity.

- $m = \text{polylog}(k)$.

Explanation: we do not need the model to have too much over-parameterization.

- $\eta' = \eta \text{poly}(k)$.

Theorem G.3. *Under Parameter G.2, for any $m = \text{polylog}(k)$ and sufficiently small $\eta \leq \frac{1}{\text{poly}(k)}$ and $\eta' = \eta \text{poly}(k)$, our Induction Hypothesis G.1 holds for all iterations $t = 0, 1, \dots, T$.*

This entire section is devoted to proving Theorem G.3, and we shall explain in the end of this section how Theorem G.3 implies Theorem 3.

Disclaimer. To make this paper more concise, in the rest of this section we highlight the key technical claims/lemmas that we need to prove Induction Hypothesis G.1. Some of the proofs we give in this section are more “sketched” because we assume the readers are already familiar with our proof languages used in Section D.

G.1 Gradient Calculations and Function Approximation

Claim G.4. *There exists some parameter $\Xi = \text{polylog}(k)$ in (F.1) so that for every $(X, y) \in \mathcal{Z}_m$,*

$$\forall i \in [k]: \quad \mathbf{logit}_i^\tau(G, X) = \begin{cases} \frac{1}{s(X)} - k^{-\Omega(\log k)}, & \text{if } v_{i,1} \text{ or } v_{i,2} \text{ is in } \mathcal{V}(X); \\ k^{-\Omega(\log k)}, & \text{if neither } v_{i,1} \text{ nor } v_{i,2} \text{ is in } \mathcal{V}(X); \end{cases}$$

where $s(X)$ is the number of indices $i \in [k]$ such that $v_{i,1}$ or $v_{i,2}$ is in $\mathcal{V}(X)$. (Recall with high probability $s(X) = \Theta(s)$.) And, for every $(X, y) \in \mathcal{Z}_s$,

$$\forall i \in [k]: \quad \mathbf{logit}_i^\tau(G, X) = \begin{cases} 1 - k^{-\Omega(\log k)}, & \text{if } i = y; \\ k^{-\Omega(\log k)}, & \text{if } i \neq y; \end{cases}$$

Proof of Claim G.4. Using the same analysis as the proof of Theorem 2, we know after ensemble, for every $i \in [k]$ and $\ell \in [2]$,

$$\frac{1}{K} \sum_{w \in [K]} \Phi_{i,\ell}^{[w]} \geq \frac{1}{\text{polylog}(k)}$$

and therefore there exists some scale-up factor $\Xi = \text{polylog}(k)$ for (F.1) so that for every $(X, y) \in \mathcal{Z}_m$, every $i \in [k]$,

- $G_i(X) \geq \log^4 k$ when either $v_{i,1}$ or $v_{i,2}$ is in $\mathcal{V}(X)$;
- $G_i(X) \leq 1$ when neither $v_{i,1}$ nor $v_{i,2}$ is in $\mathcal{V}(X)$.

and at the same time, for every $(X, y) \in \mathcal{Z}_s$, every $i \in [k]$,

- $G_i(X) \geq \log^4 k$ when $i = y$;
- $G_i(X) \leq 1$ when $i \neq y$.

Plugging this into the threshold logit function (4.2) finishes the proof. □

Note that our update rule (4.3) is not precisely the gradient of a function (due to our truncation to the negative part for simpler analysis). However, in the remainder of the proof, slightly abusing notation, let us denote by

$$\nabla_{w_{i,r}} \tilde{L}(F; X, y) \stackrel{\text{def}}{=} \nabla_{w_{i,r}} L(F; X, y) - \frac{\eta'}{\eta} (\mathbf{logit}_i^\tau(F, X) - \mathbf{logit}_i^\tau(G, X))^- \nabla_{w_{i,r}} F_i(X)$$

so that when

Fact G.5. Given data point $(X, y) \in \mathcal{D}$, for every $i \in [k]$, $r \in [m]$, up to a negligible additive error $\frac{1}{k^{\Omega(\log k)}}$, we have

$$\begin{aligned} -\nabla_{w_{i,r}} \tilde{L}(F; X, y) &= \left(\mathbb{1}_{y=i} - \mathbf{logit}_i(F, X) + \frac{\eta'}{\eta} \mathbb{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(\frac{1}{s(X)} - \mathbf{logit}_i^r(F, X) \right)^+ \right) \nabla_{w_{i,r}} F_i(X) \\ &\quad \text{when } (X, y) \in \mathcal{Z}_m \\ -\nabla_{w_{i,r}} \tilde{L}(F; X, y) &= \left(\mathbb{1}_{y=i} - \mathbf{logit}_i(F, X) + \frac{\eta'}{\eta} \mathbb{1}_{y=i} (1 - \mathbf{logit}_i^r(F, X))^+ \right) \nabla_{w_{i,r}} F_i(X) \\ &\quad \text{when } (X, y) \in \mathcal{Z}_s \end{aligned}$$

where recall $\nabla_{w_{i,r}} F_i(X) = \sum_{p \in [P]} \widetilde{\text{ReLU}}'(\langle w_{i,r}, x_p \rangle) x_p$

Because $\frac{1}{k^{\Omega(\log k)}}$ is negligible, for proof simplicity, we ignore it in the rest of the proof.

We also summarize a simple calculation that is analogous to Claim D.7 and Claim D.8.

Claim G.6 (gradient, c.f. Claim D.7 and D.8). For every $t \leq T$, for every $(X, y) \in \mathcal{Z}$, every $i \in [k]$, $r \in [m]$ and $\ell \in [2]$, we have:

- If $v_{i,1}, v_{i,2} \in \mathcal{V}(X)$ then $\langle \nabla_{w_{i,r}} F_i^{(t)}(X), v_{i,\ell} \rangle \geq (V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P))$
- $\langle \nabla_{w_{i,r}} F_i^{(t)}(X), v_{i,\ell} \rangle \leq (\mathbb{1}_{v_{i,\ell} \in \mathcal{V}(X)} V_{i,r,\ell}(X) + \mathcal{E}_1 + \mathcal{E}_3)$
- for every $j \in [k] \setminus \{i\}$, $\left| \langle -\nabla_{w_{i,r}} F_i^{(t)}(X), v_{j,\ell} \rangle \right| \leq (\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X))$

Notation. Throughout the remainder of the proof, let us use $v_{i,1}, v_{i,2} \in \mathcal{V}(X)$ to denote that at least one of $v_{i,1}, v_{i,2}$ is in \mathcal{V} . This simplifies our notations.

Recall

$$\Phi_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \sum_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+ \quad \text{and} \quad \Phi_i^{(t)} \stackrel{\text{def}}{=} \sum_{\ell \in [2]} \Phi_{i,\ell}^{(t)}$$

and this time we have

Claim G.7 (function approximation, c.f. Claim D.9). Under the new Induction Hypothesis G.1, let us define $Z_{i,\ell}^{(t)}(X) := \mathbb{1}_{v_{i,\ell} \in \mathcal{V}(X)} \left(\sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p \right)$, we have: for every t , every $i \in [k]$, every $(X, y) \in \mathcal{Z}$ (or for every new sample $(X, y) \sim \mathcal{D}$, with probability at least $1 - e^{-\Omega(\log^2 k)}$):

$$F_i^{(t)}(X) = \sum_{\ell \in [2]} \left(\Phi_{i,\ell}^{(t)} \times Z_{i,\ell}^{(t)}(X) \right) \pm \tilde{O}(\sigma_0 \cdot m) = \sum_{\ell \in [2]} \left(\Phi_{i,\ell}^{(t)} \times Z_{i,\ell}^{(t)}(X) \right) \pm O\left(\frac{1}{\text{polylog}(k)}\right)$$

G.2 Useful Claims as Consequences of the Induction Hypothesis

Recall we had five useful claims in Section D.2 for the proof of the single model case. This time, we only have three and they are also easier than their counterparts in Section D.2.

G.2.1 Lambda Growth

Claim G.8 (growth, c.f. Claim D.10). Suppose Induction Hypothesis G.1 holds at iteration t , then for every $i \in [k]$, $\ell \in [2]$, suppose $\Phi_{i,\ell}^{(t)} \leq \frac{1}{\tau}$, then it satisfies

$$\Phi_{i,\ell}^{(t+1)} \geq \Phi_{i,\ell}^{(t)} + \tilde{\Omega}\left(\frac{\eta'}{k}\right) \widetilde{\text{ReLU}}'(\Phi_{i,\ell}^{(t)})$$

Proof of Claim D.10. Using the same calculation as (D.5), but this time substituting the new gradient formula in Fact G.5, we have

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle - O(\eta + \eta' \frac{N_s}{N}) \\ &\quad + \Omega(\eta') \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P) \right) \left(\frac{1}{s(X)} - \mathbf{logit}_i^\tau(F^{(t)}, X) \right)^+ \right] \end{aligned}$$

Let us now consider $r \in [m]$ to be the $\arg \max_{r \in [m]} \{ \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \}$, so we have $\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq \tilde{\Omega}(\Phi_{i,\ell}^{(t)})$.

Following a similar analysis as before, we can derive that as long as $\mathbb{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} = 1$, we have $V_{i,r,\ell}(X) \geq \Omega(1) \cdot \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle) \geq \tilde{\Omega}(\widetilde{\text{ReLU}}'(\Phi_{i,\ell}^{(t)}))$.

Now, since $\Phi_{i,\ell}^{(t)} \leq \frac{1}{\tau}$, we know that as long as $v_{i,\ell} \in \mathcal{V}(X)$ and $v_{i,3-\ell} \notin \mathcal{V}(X)$ (which happens for $\Theta(\frac{s}{k})$ fraction of the multi-view training data), it satisfies (see Claim G.7):

$$F_i^{(t)}(X) \leq \Phi_{i,\ell}^{(t)} \times Z_{i,\ell}^{(t)}(X) - \tilde{O}(\sigma_0 m) = O\left(\frac{1}{\tau}\right) - \tilde{O}(\sigma_0 m) \leq O\left(\frac{1}{\tau}\right).$$

When this happens, we know $\mathbf{logit}_i^\tau(F^{(t)}, X) \leq O(\frac{1}{k})$. This implies, after summing over $r \in [m]$,

$$\Phi_{i,\ell}^{(t+1)} \geq \Phi_{i,\ell}^{(t)} + \tilde{\Omega}\left(\frac{\eta'}{k}\right) \widetilde{\text{ReLU}}'(\Phi_{i,\ell}^{(t)})$$

□

Now we can define T_0 as follows.

Claim G.9. Define iteration threshold $T_0 \stackrel{\text{def}}{=} \tilde{\Theta}\left(\frac{k}{\eta' \sigma_0^{q-2}}\right)$, then

- for every $i \in [k], \ell \in [2]$ and $t \geq T_0$, it satisfies $\Phi_{i,\ell}^{(t)} \geq \frac{1}{2\tau}$

G.2.2 Single-View Error Till the End

Claim G.10 (single-view after T_0). Suppose Induction Hypothesis G.1 holds for all iterations $< t$ and $t \geq T_0$. For every single-view data $(X, y) \in \mathcal{Z}_s$ (or any $(X, y) \in \mathcal{D}_s$ but with probability $1 - e^{-\Omega(\log^2 k)}$), we have

$$F_y^{(t)}(X) \geq \max_{i \in [k] \setminus \{y\}} F_i^{(t)}(X) + \Omega(\log k) \quad \text{and} \quad 1 - \mathbf{logit}_y(F^{(t)}, X) \leq \frac{1}{\text{poly}(k)}$$

Proof. This is because for single-view data $(X, y) \in \mathcal{Z}_s$, it satisfies $Z_{i,\ell}^{(t)}(X) \leq \Gamma$ as long as $i \neq y$. As a result, applying Claim G.9, we must have $F_y^{(t)}(X) \geq \Omega(\frac{1}{\tau}) \geq \Omega(\log k)$ but $F_i^{(t)}(X) \leq O(1)$ for $i \neq y$ (using $\Gamma < \frac{1}{\text{polylog}(k)}$). (Similar for $(X, y) \in \mathcal{D}_s$.) □

G.2.3 Multi-View Error Till the End

Claim G.11 (multi till the end, c.f. Claim D.14). Suppose Induction Hypothesis G.1 holds for every iteration $t < T$, then

- $\eta \sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y=i} (1 - \mathbf{logit}_y(F^{(t)}, X)) \right] \leq \tilde{O}(1)$
- $\frac{\eta' N_s}{N} \sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbb{1}_{y=i} (1 - \mathbf{logit}_i^\tau(F^{(t)}, X))^+ \right] \leq \tilde{O}(1)$
- $\eta' \sum_{t=T_0}^T \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(\frac{1}{s(X)} - \mathbf{logit}_i^\tau(F^{(t)}, X) \right)^+ \right] \leq \tilde{O}(1)$

Proof of Claim G.11. By Fact G.5 and Claim G.6 again (similar to the calculation in the proof of Claim G.8), we have

$$\begin{aligned}
\langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &\geq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle - \frac{1}{\text{poly}(k)} \\
&\quad + \Omega(\eta) \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y=i} \left(V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P) \right) \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right] \\
&\quad + \Omega\left(\frac{\eta' N_s}{N}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbb{1}_{y=i} \left(V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P) \right) \left(1 - \mathbf{logit}_i^\tau(F^{(t)}, X) \right)^+ \right] \\
&\quad + \Omega(\eta') \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(V_{i,r,\ell}(X) - \tilde{O}(\sigma_p P) \right) \left(\frac{1}{s(X)} - \mathbf{logit}_i^\tau(F^{(t)}, X) \right)^+ \right]
\end{aligned}$$

In the above formula, we can ignore the single-view data for the $(1 - \mathbf{logit}_i(F^{(t)}, X))$ term because they are extremely small (see Claim G.10).

Now, if we take $r = \arg \max_{r \in [m]} \{\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle\}$, we must have (whenever $v_{i,1}, v_{i,2} \in \mathcal{V}(X)$) $V_{i,r,\ell}(X) \geq \Omega(1) \cdot \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle) \geq \tilde{\Omega}(1)$. Therefore, when summing up over all possible $r \in [m]$, we have

$$\begin{aligned}
\Phi_{i,\ell}^{(t+1)} &\geq \Phi_{i,\ell}^{(t)} - \frac{\eta}{\text{poly}(k)} \\
&\quad + \tilde{\Omega}(\eta) \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y=i} \left(1 - \mathbf{logit}_i(F^{(t)}, X) \right) \right] \\
&\quad + \tilde{\Omega}\left(\frac{\eta' N_s}{N}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\mathbb{1}_{y=i} \left(1 - \mathbf{logit}_i^\tau(F^{(t)}, X) \right)^+ \right] \\
&\quad + \tilde{\Omega}(\eta') \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(\frac{1}{s(X)} - \mathbf{logit}_i^\tau(F^{(t)}, X) \right)^+ \right]
\end{aligned}$$

After telescoping, and using $\Phi^{(t)} \leq \tilde{O}(1)$, we finish the proof. \square

G.3 Main Lemmas for Proving the Induction Hypothesis

In this subsection, we provide key technical lemmas that, when combined together, shall prove that Induction Hypothesis G.1 holds for every iteration (and thus prove Theorem G.3).²⁴

G.3.1 Correlation Growth

Lemma G.12 (c.f. Lemma D.21). *Suppose Parameter G.2 holds and suppose Induction Hypothesis G.1 holds for all iterations $< t$. Then, letting $\Phi_{i,\ell}^{(t)} \stackrel{\text{def}}{=} \sum_{r \in [m]} [\langle w_{i,r}^{(t)}, v_{i,\ell} \rangle]^+$, we have*

$$\forall i \in [k], \forall \ell \in [2] : \quad \Phi_{i,\ell}^{(t)} \leq \tilde{O}(1)$$

Proof of Lemma G.12. Let us denote by $\Phi^{(t)} = \max_{i \in [k], \ell \in [2]} \Phi_{i,\ell}^{(t)}$. Suppose t is some iteration so that $\Phi^{(t)} \geq \frac{10}{\tau^2}$ but $\Phi^{(t)} \leq \tilde{O}(1)$. We wish to prove that if we continue from iteration t for at most T iterations, then $\Phi^{(t')} \leq \tilde{O}(1)$ for every $t' \in [t, T]$.

Without loss of generality, we assume that $\Phi^{(t)} \geq \frac{10}{\tau^2}$ always holds from iteration t onwards (because otherwise we can start with the next iteration t' so that $\Phi^{(t)}$ goes above $\frac{10}{\tau^2}$).

²⁴We only sketch the proofs to these technical lemmas, but ignore the last step of putting them together to prove Theorem G.3 because it is trivial (but anyways almost identical to that in Section D.5).

Let $(i, \ell) = \arg \max_{i \in [k], \ell \in [2]} \Phi_{i, \ell}^{(t)}$. Then, for every $(X, y) \in \mathcal{Z}_m$, (G.7) tells us that $F_i^{(t)}(X) \geq \frac{1}{\tau^2}$ for every i such that $v_{i,1}, v_{i,2} \in \mathcal{V}(X)$. Therefore,

$$\mathbb{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(\frac{1}{s(X)} - \mathbf{logit}_i^r(F, X) \right)^+ \leq \frac{1}{k^{\Omega(\log k)}}$$

Also, for every $(X, y) \in \mathcal{Z}_s$, for similar reason we have

$$(1 - \mathbf{logit}_i^r(F, X))^+ \leq \frac{1}{k^{\Omega(\log k)}}$$

Therefore, at this iteration t , up to negligible $\frac{1}{k^{\Omega(\log k)}}$ terms, we have according to Claim G.7:

$$\begin{aligned} -\nabla_{w_{i,r}} \tilde{L}(F; X, y) &= (1 - \mathbf{logit}_i(F, X)) \nabla_{w_{i,r}} F_i(X) && \text{when } i = y \\ -\nabla_{w_{i,r}} \tilde{L}(F; X, y) &= -\mathbf{logit}_i(F, X) \nabla_{w_{i,r}} F_i(X) && \text{when } i \neq y \end{aligned}$$

This is already identical to what we had in the single-model case (without distillation).

This time, we can calculate (using Claim G.6)

$$\begin{aligned} \langle w_{i,r}^{(t+1)}, v_{i,\ell} \rangle &= \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} [\langle -\nabla_{w_{i,r}} \tilde{L}(F^{(t)}; X, y), v_{i,\ell} \rangle] \\ &\leq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + \eta \mathbb{E}_{(X,y) \sim \mathcal{Z}} [\mathbb{1}_{y=i} (V_{i,r,\ell}(X) + \mathcal{E}_1 + \mathcal{E}_3) (1 - \mathbf{logit}_i(F, X))] \\ &\leq \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle + O(\eta) \mathbb{E}_{(X,y) \sim \mathcal{Z}} [\mathbb{1}_{y=i} (1 - \mathbf{logit}_y(F, X))] \end{aligned}$$

- For every $(X, y) \in \mathcal{Z}_s$, since $\Phi_{i,\ell}^{(t)} \geq \frac{10}{\tau^2}$, it satisfies (similar to Claim G.10) that $F_y^{(t)}(X) \geq \Omega(\frac{1}{\tau^2}) \gg \Omega(\log^2 k)$ but $F_j^{(t)}(X) \leq O(1)$ for $j \neq y$ (using $\Gamma < \frac{1}{\text{polylog}(k)}$). This implies

$$(1 - \mathbf{logit}_y(F, X)) \leq \frac{1}{k^{\Omega(\log k)}}$$

- For every $(X, y) \in \mathcal{Z}_m$ with $y = i$, recall from Claim G.7 that

$$F_j^{(t)}(X) = \sum_{\ell \in [2]} \left(\Phi_{j,\ell}^{(t)} \times \mathbb{1}_{v_{j,\ell} \in \mathcal{V}(X)} \left(\sum_{p \in \mathcal{P}_{v_{j,\ell}}(X)} z_p \right) \right) \pm O\left(\frac{1}{\text{polylog}(k)}\right)$$

By our choice of the distribution, this implies

- $F_j^{(t)}(X) \leq 0.8001\Phi^{(t)}$ for $j \neq i$, and
- $F_i^{(t)}(X) \geq 0.9999\Phi^{(t)}$ because (i, ℓ) is the argmax of $\Phi_{i,\ell}^{(t)}$.

This again means

$$(1 - \mathbf{logit}_y(F, X)) \leq \frac{1}{k^{\Omega(\log k)}}$$

Together, and summing up over all $r \in [m]$, we have

$$\Phi^{(t+1)} \leq \Phi^{(t)} + \frac{\eta m}{k^{\Omega(\log k)}}$$

so if we continue this for T iterations we still have $\Phi^{(T)} \leq \tilde{O}(1)$. □

G.3.2 Off-Diagonal Correlations are Small

Lemma G.13 (c.f. Lemma D.22). *Suppose Parameter G.2 holds and suppose Induction Hypothesis G.1 holds for all iterations $< t$. Then,*

$$\forall i \in [k], \forall r \in [m], \forall j \in [k] \setminus \{i\}: \quad |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| \leq \tilde{O}(\sigma_0)$$

Proof of Lemma G.13. We separately treat $t \leq T_0$ and $t \geq T_0$. (This should be reminiscent of the three-stage proof in the single model case.)

Consider $t \leq T_0$. By Fact G.5 and Claim G.6

$$|\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| \leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + O\left(\eta + \frac{\eta' N_s}{N}\right) \\ + O(\eta') \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\left(\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X) \right) \mathbf{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(\frac{1}{s(X)} - \mathbf{logit}_i^r(F, X) \right)^+ \right]$$

Using the property that $v_{i,1}, v_{i,2} \in \mathcal{V}(X)$ with probability $\Theta(\frac{s}{k})$ over a sample $(X, y) \in \mathcal{Z}_m$, and using the trivial bound $\left(\frac{1}{s(X)} - \mathbf{logit}_i^r(F, X)\right)^+ \leq \Theta(\frac{1}{s})$, we immediately have

$$|\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| \leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + O\left(\eta + \frac{\eta' N_s}{N}\right) \\ + O\left(\frac{\eta'}{k}\right) \left(\gamma + (\sigma_0^{q-1})\gamma s + \tilde{O}((\sigma_0 \gamma k)^{q-1}) \gamma P + (\sigma_0)^{q-1} \frac{s}{k} \right)$$

Finally, using $T_0 = \tilde{\Theta}(\frac{k}{\eta' \sigma_0^{q-2}})$, $N_s \leq \frac{N}{\text{poly}(k)}$, and using $\eta \leq \frac{\eta'}{\text{poly}(k)}$, together with the same parameter choices as before, we conclude that $|\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| \leq \tilde{O}(\sigma_0)$ for every $t \leq T_0$.

Consider $t > T_0$. By Fact G.5 and Claim G.6 again

$$|\langle w_{i,r}^{(t+1)}, v_{j,\ell} \rangle| \leq |\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| + O\left(\frac{\eta}{\text{poly}(k)}\right) \\ + O(\eta) \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\left(\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X) \right) \mathbf{1}_{y=i} (1 - \mathbf{logit}_i(F, X)) \right] \\ + O(\eta) \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\left(\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X) \right) \mathbf{1}_{y \neq i} (\mathbf{logit}_i(F, X)) \right] \\ + O\left(\frac{\eta' N_s}{N}\right) \mathbb{E}_{(X,y) \sim \mathcal{Z}_s} \left[\left(\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X) \right) \mathbf{1}_{y=i} (1 - \mathbf{logit}_i^r(F, X))^+ \right] \\ + O(\eta') \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\left(\mathcal{E}_{2,i,r}(X) + \mathcal{E}_1 + \mathcal{E}_3 + \mathcal{E}_{4,j,\ell}(X) \right) \mathbf{1}_{v_{i,1}, v_{i,2} \in \mathcal{V}(X)} \left(\frac{1}{s(X)} - \mathbf{logit}_i^r(F, X) \right)^+ \right]$$

In the above formula, we can ignore the single-view data for the $(1 - \mathbf{logit}_y(F^{(t)}, X))$ and $\mathbf{logit}_i(F^{(t)}, X)$ (for $i \neq y$) terms because they are extremely small (see Claim G.10).

Now, applying the naive upper bounds $\mathcal{E}_{2,i,r}(X) \leq \gamma$ and $\mathcal{E}_{4,j,\ell}(X) \leq \tilde{O}(\sigma_0^{q-1})$, and telescoping for all $t \geq T_0$ and applying Claim G.11, we immediately have $|\langle w_{i,r}^{(t)}, v_{j,\ell} \rangle| \leq |\langle w_{i,r}^{(T_0)}, v_{j,\ell} \rangle| + \tilde{O}(\sigma_0)$. \square

G.3.3 Noise Correlation is Small

Lemma G.14 (c.f. Lemma D.25). *Suppose Parameter G.2 holds and suppose Induction Hypothesis G.1 holds for all iterations $< t$. For every $\ell \in [2]$, for every $r \in [m]$, for every $(X, y) \in \mathcal{Z}$ and $i \in [k]$:*

- (a) *For every $p \in \mathcal{P}_{v_{i,\ell}}(X)$, we have: $\langle w_{i,r}^{(t)}, \xi_p \rangle \leq \tilde{o}(\sigma_0)$.*
- (b) *For every $p \in \mathcal{P}(X) \setminus (\mathcal{P}_{v_{i,1}}(X) \cup \mathcal{P}_{v_{i,2}}(X))$, we have: $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0)$.*
- (c) *For every $p \in [P] \setminus \mathcal{P}(X)$, we have: $|\langle w_{i,r}^{(t)}, \xi_p \rangle| \leq \tilde{O}(\sigma_0 \gamma k)$.*

Proof. From a similar calculation (see (D.6)) we have for every $(X, y) \in \mathcal{Z}$ and $p \in [P]$,

$$\langle w_{i,r}^{(t+1)}, \xi_p \rangle = \langle w_{i,r}^{(t)}, \xi_p \rangle \pm \tilde{\Theta} \left(\frac{\eta'}{N} \right) \widetilde{\text{ReLU}}'(\langle w_{i,r}^{(t)}, x_p \rangle) \pm \frac{\eta}{\sqrt{d}}$$

After telescoping and using $N \gg T$ from Parameter G.2, we immediately finish the proof. (If one instead wishes to consider the case of $T \geq N$, she has to do a more careful calculation here. We skip it to keep this paper concise. \square)

G.3.4 Diagonal Correlations are Nearly Non-Negative

Lemma G.15 (c.f. Lemma D.27). *Suppose Parameter G.2 holds and suppose Induction Hypothesis G.1 holds for all iterations $< t$. Then,*

$$\forall i \in [k], \forall r \in [m], \forall \ell \in [2]: \quad \langle w_{i,r}^{(t)}, v_{i,\ell} \rangle \geq -\tilde{O}(\sigma_0) .$$

The proof of Lemma G.15 is almost identical to Lemma D.27 so we skip here.

G.4 Proof of Theorem 3

First of all, applying Claim G.11 and $T \geq \frac{\text{poly}(k)}{\eta}$, we know there are at most 90% of the iterations $t \leq T_0$ satisfying

$$\mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[(1 - \mathbf{logit}_y(F^{(t)}, X)) \right] \leq \frac{1}{\text{poly}(k)}$$

Applying Claim D.16, we immediately have the test accuracy result for multi-view data.

Applying Claim G.10 (which uses $\Phi_{i,\ell}^{(t)} \geq \Omega(\log k)$), we immediately have the test accuracy result for single-view data. \blacksquare

G.5 Proof of Theorem 4

We assume the readers are now familiar with the proofs of the single model Theorem 1 and of the ensemble distill Theorem 3. They easily imply Theorem 4 for reasons we explain below.

Stage 1 of F . Recall Theorem 1 (and Lemma D.21) imply that, at the end of the stage 1 for training a network F , the quantity $\Phi_i^{(T)} \in [\Omega(\log k), \tilde{O}(1)]$ for every $i \in [k]$; thus, if $(i, \ell) \in \mathcal{M}$, we must have $\Phi_{i,\ell}^{(T)} \geq \Omega(\log k)$. At the end of stage 1, also recall for every $(i, \ell) \in \mathcal{M}$, for any single-view data $(X, y) \in \mathcal{D}_s$ with $y = i$ and $\hat{\ell}(X) = \ell$, with high probability F predicts correctly on (X, y) . Let us remind the readers from (C.2) that

$$\mathcal{M} \stackrel{\text{def}}{=} \left\{ (i, \ell^*) \in [k] \times [2] \mid \Lambda_{i,\ell^*}^{(0)} \geq \Lambda_{i,3-\ell^*}^{(0)} \left(\frac{S_{i,3-\ell^*}}{S_{i,\ell^*}} \right)^{\frac{1}{q-2}} \left(1 + \frac{1}{\log^2(m)} \right) \right\}$$

where $S_{i,\ell} \stackrel{\text{def}}{=} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} \left[\mathbb{1}_{y=i} \sum_{p \in \mathcal{P}_{v_{i,\ell}}(X)} z_p^q \right]$. Since in this self-distillation theorem, we have Assumption 4.1 which says the distribution of $\sum_{p \in \mathcal{P}_v(X)} z_p^q$ for $v \in \{v_{y,1}, v_{y,2}\}$ are the same over multi-view data, by standard concentration, we know with high probability $S_{i,1} = S_{i,2} \left(1 \pm \frac{1}{2 \log^2 k} \right)$ for every $i \in [k]$. This means, we can alternatively define

$$\mathcal{M}_F \stackrel{\text{def}}{=} \left\{ (i, \ell^*) \in [k] \times [2] \mid \Lambda_{i,\ell^*}^{(0)} \geq \Lambda_{i,3-\ell^*}^{(0)} \left(1 + \frac{2}{\log^2(m)} \right) \right\}$$

(which is a subset of \mathcal{M}) and all the statements about \mathcal{M} also apply to \mathcal{M}_F .

Base Model G. For a similar reason, if $(i, \ell) \in \mathcal{M}_G$ for the distill model G , then we have that the quantity $\Phi_{i,\ell}^{(T)} \geq \text{polylog}(k)$ for network G (recall G is scaled up by a $\text{polylog}(k)$ factor). Using a similar analysis to Claim G.4, we can derive that:

- for every $(X, y) \in \mathcal{Z}_m$, we have

$$\forall (i, \ell) \in \mathcal{M}_G: \quad \mathbf{logit}_i^\tau(G, X) \begin{cases} \geq \frac{1}{s(X)} - k^{-\Omega(\log k)}, & \text{if } v_{i,\ell} \text{ is in } \mathcal{V}(X); \\ = k^{-\Omega(\log k)}, & \text{if } v_{i,\ell} \text{ is not in } \mathcal{V}(X); \end{cases} \quad (\text{G.1})$$

$$\forall i \in [k]: \quad \mathbf{logit}_i^\tau(G, X) \begin{cases} \leq \frac{1}{s'(X)} + k^{-\Omega(\log k)}, & \text{if } v_{i,1} \text{ or } v_{i,2} \text{ is in } \mathcal{V}(X); \\ = k^{-\Omega(\log k)}, & \text{if neither } v_{i,1} \text{ or } v_{i,2} \text{ is not in } \mathcal{V}(X); \end{cases} \quad (\text{G.2})$$

where recall $s(X)$ is the number of indices $i \in [k]$ such that $v_{i,1}$ or $v_{i,2}$ is in $\mathcal{V}(X)$, and we newly define $s'(X)$ as the number of indices $i \in [k]$ such that $(i, \ell) \in \mathcal{M}_G$ and $v_{i,\ell} \in \mathcal{V}(X)$ for some $\ell \in [2]$. One can derive using concentration that with high probability $\frac{s}{2} \leq s'(X) \leq s(X) \leq 3s$ for all multi-view training data.²⁵

Stage 2 of F. Similar to the proof of Theorem 3, we can ignore single-view data's contribution to the gradient updates (since they are negligible) and only focus on multi-view data. Using a similar gradient calculation to Fact G.5, we know that, up to some small error,

- the quantity $\Phi_{i,\ell} = \sum_{r \in [m]} [\langle w_{i,r}, v_{i,\ell} \rangle]^+$ never decreases during stage 2.
- the quantity $\Phi_i = \sum_{r \in [m], \ell \in [2]} [\langle w_{i,r}, v_{i,\ell} \rangle]^+$ no longer changes during stage 2, when it reaches $\Phi_i \geq \frac{2}{\tau^2}$.²⁶

Recall after stage 1, we have $\Phi_i \geq \Omega(\log k)$ for network F ; but since at the beginning of stage 2 we have scaled up F by a factor of $\log^4 k$, this means $\Phi_i \geq \frac{2}{\tau^2}$ is already satisfied at the beginning of stage 2 (up to small error), so it does not change during stage 2. As a result, at the end of stage 2, network F should give the same (nearly perfect) accuracy on multi-view data as claimed in Theorem 1.

Furthermore, through a similar analysis to Claim G.8 and Claim G.9 (and combining with (G.1)), we know that for when $(i, \ell) \in \mathcal{M}_G$, the quantity $\Phi_{i,\ell}$ must increase to at least $\frac{1}{2\tau} \geq \Omega(\log^2 k)$. This allows us to conclude that, when $(i, \ell) \in \mathcal{M}_G$, at the end of stage 2, for those single-view data $(X, y) \in \mathcal{D}_s$ with $y = i$ and $\widehat{\ell}(X) = \ell$, with high probability $F_y(X) \geq \Omega(\log^2 k) \gg F_j(X)$ for $j \neq y$, so F predicts correctly on (X, y) .

At the same time, for every $(i, \ell) \in \mathcal{M}_F$, we have $\Phi_{i,\ell} \geq \Omega(\log k)$ is already satisfied at the end of stage 1, so at the end of stage 2 it must also satisfy $\Phi_{i,\ell} \geq \Omega(\log^5 k)$ (the extra factors are due to scale-up). Thus, F also predicts correctly on those single-view data $(X, y) \in \mathcal{D}_s$ with $y = i$ and $\widehat{\ell}(X) = \ell$.

Finally, using $|\mathcal{M}_F| \geq k(1 - o(1))$ and $|\mathcal{M}_G| \geq k(1 - o(1))$, together with the fact that they are totally independent random sets, we obtain $|\mathcal{M}_F \cup \mathcal{M}_G| \geq \frac{3}{2}k(1 - o(1))$. This means, learned model $F^{(T+T')}$ at the end of stage 2 through self-distillation has an accuracy of $\geq \frac{3}{4}(1 - o(1))$ over single-view data. ■

²⁵We remark here that \mathcal{M}_G does not depend on the randomness of the training set, so the lower bound $\frac{s}{10} \leq s'(X)$ can be derived trivially using $|\mathcal{M}_G| \geq k(1 - o(1))$.

²⁶Whenever $\Phi_i \geq \frac{2}{\tau^2}$ for network F , one can verify that $F_i(X) \geq \frac{1}{\tau^2}$ for every multi-view data $(X, y) \in \mathcal{Z}_m$ with $v_{i,1}, v_{i,2} \in \mathcal{V}(X)$. This means $\mathbf{logit}_i^\tau(F, X) \geq \frac{1}{s(X)} - k^{-\Omega(\log k)}$ for every multi-view data (X, y) with $v_{i,1}, v_{i,2} \in \mathcal{V}(X)$. When this happens, combining with (G.2), we have $[\mathbf{logit}_i^\tau(G, X) - 10\mathbf{logit}_i^\tau(F, X)]^+ = 0$ so (up to small error) there is no gradient and $\langle w_{i,r}, v_{i,\ell} \rangle$ stays unchanged.

H Simple Probability Lemmas

We first state a simple proposition that directly implies Fact C.1.

Proposition H.1. *Given m i.i.d. standard Gaussian random variables $g_1, \dots, g_m \sim \mathcal{N}(0, 1)$, with probability at least $1 - \delta$, we have that except for at most $O(\log(1/\delta))$ indices $i \in [m]$, we have*

$$g_i \leq \max_{j \in [m]} \{g_j\} \cdot \left(1 - \Omega\left(\frac{1}{\log(m/\log(1/\delta))}\right)\right)$$

Proof of Proposition H.1. Recall for every $x > 0$

$$\left(\frac{1}{x} - \frac{1}{x^3}\right) \frac{e^{-x^2/2}}{\sqrt{2\pi}} \leq \Pr_g[g > x] \leq \frac{1}{x} \frac{e^{-x^2/2}}{\sqrt{2\pi}} \quad (\text{H.1})$$

The probability for one of them to exceed

$$\Pr[\max_i g_i > x] = 1 - (1 - \Pr_g[g > x])^m$$

Let us choose x^* so that $\Pr[\max_i g_i > x^*] = 1 - \delta/2$. By the asymptotic bound above, it is easy to derive that $x^* = \Theta(\sqrt{\log(m/\log(1/\delta))})$ and $\Pr[g > x^*] = \Theta(\frac{\log(1/\delta)}{m})$.

Now, consider $x = x^* - \frac{1}{x^*} = x^*(1 - \frac{1}{(x^*)^2})$. By the asymptotic bound above, it is not hard to see

$$\Pr[g > x] \leq O(1) \cdot \Pr[g > x^*] \leq O\left(\frac{\log(1/\delta)}{m}\right)$$

By Chernoff bound, we know with probability at least $1 - \delta$, it satisfies that

$$\Pr\left[\sum_{i=1}^m \mathbf{1}_{g_i > x} \geq \Omega(\log(1/\delta))\right] \leq \delta/2$$

□

We next state a proposition that shall be used to prove Proposition C.2.

Proposition H.2. *Consider two sequences of i.i.d. Gaussian, $g_1, \dots, g_m \sim \mathcal{N}(0, 1)$ and $h_1, \dots, h_m \sim \mathcal{N}(0, \sigma^2)$. Then,*

- when $\sigma > 1$, we have $\Pr[\max_{i \in [m]} g_i > \max_{i \in [m]} h_i] \geq \Omega\left(\frac{1}{\sigma}\right) \frac{1}{m\sigma^2 - 1}$.
- when $\sigma \leq 1$, given $\tau > 0$, we have $\Pr[\max_{i \in [m]} g_i = (1 \pm O(\tau)) \max_{i \in [m]} h_i] \leq O(\tau \log m + \frac{1}{\text{poly}(m)})$.

Proof of Proposition H.2. The case of $\sigma \leq 1$ is trivial and is simply by symmetry.

The case of $\sigma > 1$. For any threshold $x > 0$, we have

$$\Pr[\max_i g_i > x] = 1 - (1 - \Pr_{g \sim \mathcal{N}(0,1)}[g > x])^m$$

$$\Pr[\max_i h_i < x] = (1 - \Pr_{h \sim \mathcal{N}(0, \sigma^2)}[h > x])^m$$

Let $x^* > \sigma$ be a threshold satisfying $\frac{\sigma}{x^*} \frac{e^{-(x^*)^2/2\sigma^2}}{\sqrt{2\pi}} = \frac{1}{m}$. By the earlier Gaussian tail bound (H.1), it is easy to verify

$$\Pr[\max_i h_i < x^*] = (1 - \Pr_{h \sim \mathcal{N}(0, \sigma^2)}[h > x^*])^m \geq \Omega(1)$$

Using the earlier Gaussian tail bound (H.1), one can also verify that

$$\Pr_{g \sim \mathcal{N}(0,1)} [g > x^*] \geq \frac{1}{2x^*} \frac{e^{-(x^*)^2/2}}{\sqrt{2\pi}} \geq \frac{1}{2\sigma} \left(\frac{1}{m}\right)^{\sigma^2}$$

Therefore,

$$\Pr[\max_i g_i > x^*] = 1 - (1 - \Pr_{g \sim \mathcal{N}(0,1)} [g > x^*])^m \geq 1 - (1 - \Omega(m) \cdot \frac{1}{\sigma m \sigma^2})$$

Combining both, we have

$$\Pr[\max_i h_i < x^* < \max_i g_i] = \Pr[\max_i h_i < x^*] \cdot \Pr[\max_i g_i > x^*] \geq \frac{\Omega(1)}{\sigma m \sigma^2 - 1}$$

The case of $\sigma \leq 1$. Let us first generate h and then generate g . Since $\sigma \leq 1$, we have with probability at least $1 - \frac{1}{\text{poly}(m)}$, it satisfies that $0 < \max_{i \in [m]} h_i \leq O(\sqrt{\log m})$. When this happens, denoting by $z = \max_{i \in [m]} h_i$, we can apply a known anti-concentration result for the maximum of Gaussian variables [21, Theorem 3]:

$$\Pr \left[\max_{i \in [m]} g_i \in [z(1 - \tau), z(1 + \tau)] \right] \leq O(z\tau) \cdot O(\mathbb{E}[\max_{i \in [m]} g_i]) \leq O(z\tau \sqrt{\log m})$$

This finishes the proof. □

Let us restate Proposition C.2 for the readers' convenience. Suppose we denote by $S_{i,\ell} \stackrel{\text{def}}{=} \mathbb{E}_{(X,y) \sim \mathcal{Z}_m} [\mathbb{1}_{y=i} \sum_{p \in P_{v_{i,\ell}}(X)} z_p^q]$. Then, define

$$\mathcal{M} \stackrel{\text{def}}{=} \left\{ (i, \ell^*) \in [k] \times [2] \mid \Lambda_{i,\ell^*}^{(0)} \geq \Lambda_{i,3-\ell^*}^{(0)} \left(\frac{S_{i,3-\ell^*}}{S_{i,\ell^*}} \right)^{\frac{1}{q-2}} \left(1 + \frac{1}{\log^2(m)} \right) \right\}$$

Proposition C.2. *Suppose $m \leq \text{poly}(k)$. We have the following properties about \mathcal{M} .*

- For every $i \in [k]$, at most one of $(i, 1)$ or $(i, 2)$ is in \mathcal{M} (obvious).
- For every $i \in [k]$, suppose $S_{i,\ell} \geq S_{i,3-\ell}$, then

- $\Pr [(i, 3 - \ell) \in \mathcal{M}] \geq m^{-O(1)}$.
- $\Pr [(i, \ell) \in \mathcal{M} \text{ or } (i, 3 - \ell) \in \mathcal{M}] \geq 1 - o(1)$

Proof of Proposition C.2. We only prove the second item since the first one is trivial. Suppose $S_{i,\ell} \geq S_{i,3-\ell}$.

By our assumption on the data distribution, it is easy to verify $S_{i,1}/S_{i,2} > 0$ is a constant for every $i \in [k]$. Therefore, $(i, 3 - \ell) \in \mathcal{M}$ with probability at least $\frac{1}{m^{O(1)}}$ following the first item of Proposition H.2.

Finally, if neither (i, ℓ) or $(i, 3 - \ell)$ is in \mathcal{M} , it necessarily implies

$$\Lambda_{i,\ell}^{(0)} = \Lambda_{i,3-\ell}^{(0)} \left(\frac{S_{i,3-\ell}}{S_{i,\ell}} \right)^{\frac{1}{q-2}} \left(1 \pm O\left(\frac{1}{\log^2 m}\right) \right)$$

but according to the second item of Proposition H.2, this happens with probability at most $\frac{1}{\log m}$. □

References

- [1] Monther Alhamdoosh and Dianhui Wang. Fast decorrelated neural network ensembles with random weights. *Information Sciences*, 264:104–117, 2014.

- [2] Zeyuan Allen-Zhu and Yuanzhi Li. What Can ResNet Learn Efficiently, Going Beyond Kernels? In *NeurIPS*, 2019. Full version available at <http://arxiv.org/abs/1905.10337>.
- [3] Zeyuan Allen-Zhu and Yuanzhi Li. Can SGD Learn Recurrent Neural Networks with Provable Generalization? In *NeurIPS*, 2019. Full version available at <http://arxiv.org/abs/1902.01028>.
- [4] Zeyuan Allen-Zhu and Yuanzhi Li. Backward feature correction: How deep learning performs deep learning. *arXiv preprint arXiv:2001.04413*, 2020.
- [5] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training recurrent neural networks. In *NeurIPS*, 2019. Full version available at <http://arxiv.org/abs/1810.12065>.
- [6] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. In *ICML*, 2019. Full version available at <http://arxiv.org/abs/1811.03962>.
- [7] Jose M Alvarez, Yann LeCun, Theo Gevers, and Antonio M Lopez. Semantic road segmentation via multi-scale ensembles of learned features. In *European Conference on Computer Vision*, pages 586–595. Springer, 2012.
- [8] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *arXiv preprint arXiv:1904.11955*, 2019.
- [9] Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *CoRR*, abs/1901.08584, 2019. URL <http://arxiv.org/abs/1901.08584>.
- [10] Ainesh Bakshi, Rajesh Jayaram, and David P Woodruff. Learning two layer rectified neural networks in polynomial time. *arXiv preprint arXiv:1811.01885*, 2018.
- [11] Verónica Bolón-Canedo and Amparo Alonso-Betanzos. Ensembles for feature selection: A review and future trends. *Information Fusion*, 52:1–12, 2019.
- [12] Digvijay Boob and Guanghui Lan. Theoretical properties of the global optimizer of two layer neural network. *arXiv preprint arXiv:1710.11241*, 2017.
- [13] Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- [14] Gavin Brown, Jeremy Wyatt, Rachel Harris, and Xin Yao. Diversity creation methods: a survey and categorisation. *Information Fusion*, 6(1):5–20, 2005.
- [15] Gavin Brown, Jeremy L Wyatt, and Peter Tiño. Managing diversity in regression ensembles. *Journal of machine learning research*, 6(Sep):1621–1650, 2005.
- [16] Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with gaussian inputs. *arXiv preprint arXiv:1702.07966*, 2017.
- [17] Robert Bryll, Ricardo Gutierrez-Osuna, and Francis Quek. Attribute bagging: improving accuracy of classifier ensembles by using random feature subsets. *Pattern recognition*, 36(6):1291–1302, 2003.
- [18] Jie Cai, Jiawei Luo, Shulin Wang, and Sheng Yang. Feature selection in machine learning: A new perspective. *Neurocomputing*, 300:70–79, 2018.
- [19] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep neural networks. In *Advances in Neural Information Processing Systems*, pages 10835–10845, 2019.
- [20] Yevgen Chebotar and Austin Waters. Distilling knowledge from ensembles of neural networks for speech recognition. In *Interspeech*, pages 3439–3443, 2016.
- [21] Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. Comparison and anti-concentration bounds for maxima of gaussian random vectors. *Probability Theory and Related Fields*, 162(1):47–70, 2015.
- [22] Jia Cui, Brian Kingsbury, Bhuvana Ramabhadran, George Saon, Tom Sercu, Kartik Audhkhasi, Abhinav Sethy, Markus Nussbaum-Thom, and Andrew Rosenberg. Knowledge distillation across ensembles of multilingual models for low-resource languages. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4825–4829. IEEE, 2017.
- [23] Amit Daniely. Sgd learns the conjugate kernel class of the network. In *Advances in Neural Information Processing Systems*, pages 2422–2430, 2017.
- [24] Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The

- power of initialization and a dual view on expressivity. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2253–2261, 2016.
- [25] Thomas G Dietterich. Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, pages 1–15. Springer, 2000.
- [26] Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, November 2018.
- [27] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- [28] David A Freedman et al. Bootstrapping regression models. *The Annals of Statistics*, 9(6):1218–1228, 1981.
- [29] Markus Freitag, Yaser Al-Onaizan, and Baskaran Sankaran. Ensemble distillation for neural machine translation. *arXiv preprint arXiv:1702.01802*, 2017.
- [30] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- [31] Yoav Freund, Robert Schapire, and Naoki Abe. A short introduction to boosting. *Journal-Japanese Society For Artificial Intelligence*, 14(771-780):1612, 1999.
- [32] Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). *The annals of statistics*, 28(2):337–407, 2000.
- [33] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- [34] Takashi Fukuda, Masayuki Suzuki, Gakuto Kurata, Samuel Thomas, Jia Cui, and Bhuvana Ramabhadran. Efficient knowledge distillation from an ensemble of teachers. In *Interspeech*, pages 3697–3701, 2017.
- [35] Tommaso Furlanello, Zachary C Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again neural networks. *arXiv preprint arXiv:1805.04770*, 2018.
- [36] Mikel Galar, Alberto Fernandez, Edurne Barrenechea, Humberto Bustince, and Francisco Herrera. A review on ensembles for the class imbalance problem: bagging-, boosting-, and hybrid-based approaches. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 42(4):463–484, 2011.
- [37] Rong Ge, Jason D Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. *arXiv preprint arXiv:1711.00501*, 2017.
- [38] Rong Ge, Rohith Kuditipudi, Zhize Li, and Xiang Wang. Learning two-layer neural networks with symmetric inputs. *arXiv preprint arXiv:1810.06793*, 2018.
- [39] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers neural networks in high dimension. *arXiv preprint arXiv:1904.12191*, 2019.
- [40] Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. *arXiv preprint arXiv:1909.05989*, 2019.
- [41] Lars Kai Hansen and Peter Salamon. Neural network ensembles. *IEEE transactions on pattern analysis and machine intelligence*, 12(10):993–1001, 1990.
- [42] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [43] Tin Kam Ho. The random subspace method for constructing decision forests. *IEEE transactions on pattern analysis and machine intelligence*, 20(8):832–844, 1998.
- [44] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [45] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, pages 586–594, 2016.

- [46] Josef Kittler, Mohamad Hatef, Robert PW Duin, and Jiri Matas. On combining classifiers. *IEEE transactions on pattern analysis and machine intelligence*, 20(3):226–239, 1998.
- [47] Ron Kohavi, George H John, et al. Wrappers for feature subset selection. *Artificial intelligence*, 97(1-2):273–324, 1997.
- [48] J Zico Kolter and Marcus A Maloof. Dynamic weighted majority: An ensemble method for drifting concepts. *Journal of Machine Learning Research*, 8(Dec):2755–2790, 2007.
- [49] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- [50] Anders Krogh and Jesper Vedelsby. Neural network ensembles, cross validation, and active learning. *Advances in neural information processing systems*, 7:231–238, 1994.
- [51] Ludmila I Kuncheva. *Combining pattern classifiers: methods and algorithms*. John Wiley & Sons, 2014.
- [52] Xu Lan, Xiatian Zhu, and Shaogang Gong. Knowledge distillation by on-the-fly native ensemble. In *Advances in neural information processing systems*, pages 7517–7527, 2018.
- [53] Yuanzhi Li and Zehao Dou. When can wasserstein gans minimize wasserstein distance? *arXiv preprint arXiv:2003.04033*, 2020.
- [54] Yuanzhi Li and Yingyu Liang. Provable alternating gradient descent for non-negative matrix factorization with strong correlations. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2062–2070. JMLR. org, 2017.
- [55] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, 2018.
- [56] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. In *Advances in Neural Information Processing Systems*, pages 597–607. <http://arxiv.org/abs/1705.09886>, 2017.
- [57] Yuanzhi Li, Yingyu Liang, and Andrej Risteski. Recovery guarantee of non-negative matrix factorization via alternating updates. In *Advances in neural information processing systems*, pages 4987–4995, 2016.
- [58] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *COLT*, 2018.
- [59] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial large learning rate in training neural networks. *arXiv preprint arXiv:1907.04595*, 2019.
- [60] Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang. Learning over-parametrized two-layer relu neural networks beyond ntk. *arXiv preprint arXiv:2007.04596*, 2020.
- [61] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Improving multi-task deep neural networks via knowledge distillation for natural language understanding. *arXiv preprint arXiv:1904.09482*, 2019.
- [62] Pankaj Mehta, Marin Bukov, Ching-Hao Wang, Alexandre GR Day, Clint Richardson, Charles K Fisher, and David J Schwab. A high-bias, low-variance introduction to machine learning for physicists. *Physics reports*, 810:1–124, 2019.
- [63] Hossein Mobahi, Mehrdad Farajtabar, and Peter L Bartlett. Self-distillation amplifies regularization in hilbert space. *arXiv preprint arXiv:2002.05715*, 2020.
- [64] M Arthur Munson and Rich Caruana. On feature selection, bias-variance, and bagging. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 144–159. Springer, 2009.
- [65] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2(11):e7, 2017. doi: 10.23915/distill.00007. <https://distill.pub/2017/feature-visualization>.
- [66] Luiz S Oliveira, Robert Sabourin, Flávio Bortolozzi, and Ching Y Suen. Feature selection for ensembles: A hierarchical multi-objective genetic algorithm approach. In *Seventh International Conference on Document Analysis and Recognition, 2003. Proceedings.*, pages 676–680. Citeseer, 2003.
- [67] David Opitz and Richard Maclin. Popular ensemble methods: An empirical study. *Journal of artificial intelligence research*, 11:169–198, 1999.

- [68] David W Opitz. Feature selection for ensembles. In *AAAI*, pages 379–384, 1999.
- [69] Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *arXiv preprint arXiv:1902.04674*, 2019.
- [70] Michael P Perrone and Leon N Cooper. When networks disagree: Ensemble methods for hybrid neural networks. Technical report, BROWN UNIV PROVIDENCE RI INST FOR BRAIN AND NEURAL SYSTEMS, 1992.
- [71] Robi Polikar. Ensemble based systems in decision making. *IEEE Circuits and systems magazine*, 6(3): 21–45, 2006.
- [72] Juan José Rodríguez, Ludmila I Kuncheva, and Carlos J Alonso. Rotation forest: A new classifier ensemble method. *IEEE transactions on pattern analysis and machine intelligence*, 28(10):1619–1630, 2006.
- [73] Lior Rokach. Ensemble-based classifiers. *Artificial intelligence review*, 33(1-2):1–39, 2010.
- [74] Lior Rokach. *Pattern classification using ensemble methods*, volume 75. World Scientific, 2010.
- [75] Lior Rokach and Oded Z Maimon. *Data mining with decision trees: theory and applications*, volume 69. World scientific, 2008.
- [76] Robert E Schapire, Yoav Freund, Peter Bartlett, Wee Sun Lee, et al. Boosting the margin: A new explanation for the effectiveness of voting methods. *The annals of statistics*, 26(5):1651–1686, 1998.
- [77] Vaishaal Shankar, Alex Fang, Wenshuo Guo, Sara Fridovich-Keil, Ludwig Schmidt, Jonathan Ragan-Kelley, and Benjamin Recht. Neural kernels without tangents. *arXiv preprint arXiv:2003.02237*, 2020.
- [78] Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *arXiv preprint arXiv:1707.04926*, 2017.
- [79] Daniel Soudry and Yair Carmon. No bad local minima: Data independent training error guarantees for multilayer neural networks. *arXiv preprint arXiv:1605.08361*, 2016.
- [80] Yuandong Tian. An analytical formula of population gradient for two-layered relu network and its applications in convergence and critical point analysis. *arXiv preprint arXiv:1703.00560*, 2017.
- [81] Alexey Tsymbal, Mykola Pechenizkiy, and Pádraig Cunningham. Diversity in search strategies for ensemble feature selection. *Information fusion*, 6(1):83–98, 2005.
- [82] Giorgio Valentini. An experimental bias-variance analysis of svm ensembles based on resampling techniques. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 35(6):1252–1271, 2005.
- [83] Giorgio Valentini and Thomas G Dietterich. Bias-variance analysis of support vector machines for the development of svm-based ensemble methods. *Journal of Machine Learning Research*, 5(Jul):725–775, 2004.
- [84] Santosh Vempala and John Wilmes. Polynomial convergence of gradient descent for training one-hidden-layer neural networks. *arXiv preprint arXiv:1805.02677*, 2018.
- [85] Bo Xie, Yingyu Liang, and Le Song. Diversity leads to generalization in neural networks. *arXiv preprint Arxiv:1611.03131*, 2016.
- [86] Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. *arXiv preprint arXiv:1902.04760*, 2019.
- [87] Gilad Yehudai and Ohad Shamir. On the power and limitations of random features for understanding neural networks. *arXiv preprint arXiv:1904.00687*, 2019.
- [88] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [89] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *ICCV*, pages 3713–3722, 2019.
- [90] Kai Zhong, Zhao Song, Prateek Jain, Peter L Bartlett, and Inderjit S Dhillon. Recovery guarantees for one-hidden-layer neural networks. *arXiv preprint arXiv:1706.03175*, 2017.

- [91] Zhi-Hua Zhou, Jianxin Wu, and Wei Tang. Ensembling neural networks: many could be better than all. *Artificial intelligence*, 137(1-2):239–263, 2002.
- [92] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.