
Improved Generalization Bounds for Transfer Learning via Neural Collapse

Tomer Galanti¹ András György² Marcus Hutter²

Abstract

Using representations learned by large, pretrained models, also called foundation models, in new tasks with fewer data has been successful in a wide range of machine learning problems. Recently, Galanti et al. (2022) introduced a theoretical framework for studying this transfer learning setting for classification. Their analysis is based on the recently observed phenomenon that the features learned by overparameterized deep classification networks show an interesting clustering property, called neural collapse (Papayan et al., 2020). A cornerstone of their analysis demonstrates that neural collapse generalizes from the source classes to new target classes. However, this analysis is limited as it relies on several unrealistic assumptions. In this work, we provide an improved theoretical analysis significantly relaxing these modeling assumptions.

1. Introduction

Transfer learning is a prominent approach for dealing with overfitting (see, e.g., Caruana, 1995; Bengio, 2012; Yosinski et al., 2014). A popular approach for transfer learning suggests to pretrain a large neural network on a large-scale source task (e.g., ResNet-50 on ImageNet ILSVRC, Rusakovsky et al., 2015), and then to train a relatively small network on top of the penultimate layer of the pretrained model, using the available data in the target task. In fact, due to the impressive success of transfer learning, large pretrained models that can be effectively adapted to a wide variety of tasks (Brown et al., 2020; Ramesh et al., 2021) have recently been characterized as ‘foundation models’ (Bommasani et al., 2021), emphasizing their adaptive nature.

While foundation models are intended to be generic and widely adaptive to downstream tasks, when some specifics

¹MIT, Cambridge, MA, USA ²Deepmind, London, UK. Correspondence to: Tomer Galanti <galanti@mit.edu>, András György <agyorgy@deepmind.com>.

First Workshop of Pre-training: Perspectives, Pitfalls, and Paths Forward at ICML 2022, Baltimore, Maryland, USA, 2022. Copyright 2022 by the author(s).

of the target tasks are known, often special-purpose algorithms are designed to utilize this information. This is the case, for example, in few-shot learning, where it is known in advance that the target problems come with a very small training set (Vinyals et al., 2016; Ravi & Larochelle, 2017; Finn et al., 2017; Lee et al., 2019). Even though these approaches significantly improved the state of the art for a long time, surprisingly, recent works have demonstrated that predictors trained on top of pretrained models can achieve better performance on few-shot learning benchmarks (Tian et al., 2020; Dhillon et al., 2020; Galanti et al., 2022).

Complementing the empirical results, Galanti et al. (2022) also studied this approach theoretically, and provided an explanation for its success based on the recently discovered phenomenon of neural collapse (Papayan et al., 2020). Informally, neural collapse (see Section 3) identifies training dynamics of deep networks for standard classification tasks, where the features (the output of the penultimate layer) associated with training samples belonging to the same class concentrate around their class feature mean. Galanti et al. (2022) showed that this property generalizes to new data points and new classes (e.g., the target classes), when the model is trained on several classes with many samples for each. In addition, they showed that in the presence of neural collapse, training a linear classifier on top of the learned feature map can generalize well, even if trained with few samples only. However, their analysis demonstrating that neural collapse generalizes to new classes relies on some hard-to-justify assumptions regarding the feature maps, making their results less relevant in practical situations.

In this paper we provide a stronger theoretical analysis for this problem without relying on this kind of assumptions.

2. Problem Setup

We consider a transfer learning setting, where a model is pretrained on some source task and is adapted to solve many downstream tasks. To model this problem, we assume that a final downstream task is a k -class classification problem (a *target* problem), and the auxiliary task where the feature representation is learned on is an l -class classification problem, called the *source* problem. Formally, a target task is defined by a distribution P over samples $(x, y) \in \mathcal{X} \times \mathcal{Y}_k$, where $\mathcal{X} \subset \mathbb{R}^d$ is the instance space, and $\mathcal{Y}_k = [k] := \{1, \dots, k\}$

is a label space. For a pair (x, y) with distribution P , we denote by P_c the class-conditional distribution of x given $y = c$ (i.e., $P_c(\cdot) := \mathbb{P}[x \in \cdot \mid y = c]$).

For a given target task, our goal is to learn a classifier $h : \mathcal{X} \rightarrow \mathcal{Y}_k$ that minimizes the target test error

$$L_P(h) := \mathbb{E}_{(x,y) \sim P} [\mathbb{I}[h(x) \neq y]], \quad (1)$$

where $\mathbb{I}[E]$ denotes the indicator function of an event E (equals 1 if the event holds and 0 otherwise). As usual, P is unknown and the algorithm is provided with some training data $S = \{(x_i, y_i)\}_{i=1}^{nk}$. In this work we assume that P is a balanced distribution (i.e., $P[y = c] = 1/k$) with class-conditionals $\{P_c\}_{c=1}^k$ and that the data S is also a balanced set over the k classes, that is, $S = \cup_{c=1}^k \{(x_{ci}, c)\}_{i=1}^n$, where the sets $S_c = \{x_{ci}\}_{i=1}^n \sim P_c^n$ are drawn independently.

When n is small, training a classifier on S may not be sufficient to achieve reasonable performance. To facilitate finding a good solution, we aim to find a classifier of the form $h = g \circ f$, where $f \in \mathcal{F} \subset \{f' : \mathbb{R}^d \rightarrow \mathbb{R}^p\}$ is a feature map and $g \in \mathcal{G} \subset \{g' : \mathbb{R}^p \rightarrow \mathbb{R}^k\}$ is a classifier used on the feature space \mathbb{R}^p . The idea is that the feature map f can be learned on some other problem where more data is available, and then g is trained to solve the hopefully simpler classification problem of finding y based on $f(x)$, instead of x . That is, g is actually a function of the modified training data $\{(f(x_i), y_i)\}_{i=1}^{nk}$; to emphasize this dependence, we use the notation $g_{f,S}$ for the trained classifier based on the features, and $h_{f,S} = g_{f,S} \circ f$ for the full classifier (to be specified in Section 4).

We assume that the source task helping to find f is a single l -class classification problem over the same sample space \mathcal{X} , given by a distribution \tilde{P} , and here we are interested in finding a classifier $\tilde{h} : \mathcal{X} \rightarrow \mathbb{R}^l$ of the form $\tilde{h} = \tilde{g} \circ f$, where $\tilde{g} \in \tilde{\mathcal{G}} \subset \{g' : \mathbb{R}^p \rightarrow \mathbb{R}^l\}$ is a classifier over the feature space $f(\mathcal{X}) := \{f(x) : x \in \mathcal{X}\}$. Given a training dataset $\tilde{S} = \{(\tilde{x}_i, \tilde{y}_i)\}_{i=1}^{ml}$, the classifier \tilde{h} is trained on \tilde{S} , with the goal of minimizing the source test error, $L_{\tilde{P}}(\tilde{h}_{\tilde{S}})$.

Similarly to the target task, we assume that \tilde{P} is balanced (i.e., $\tilde{P}[y = c] = 1/l$) with class-conditionals $\{\tilde{P}_c\}_{c=1}^l$, and that the dataset $\tilde{S} = \cup_{c=1}^l \{(\tilde{x}_{ci}, c)\}_{i=1}^m$ is also balanced with $\tilde{S}_c = \{\tilde{x}_{ci}\}_{i=1}^m \sim \tilde{P}_c^m$. Finally, we assume that the classes in the source and target tasks are selected randomly: that is, the sets of class-conditionals $\tilde{\mathcal{P}} = \{\tilde{P}_c\}_{c=1}^l$ and $\mathcal{P} = \{P_c\}_{c=1}^k$ for the source and target tasks are sampled i.i.d. from a distribution \mathcal{D} over a set of class-conditional distributions \mathcal{E} . This completes the definition of the distribution of P and \tilde{P} (which are random themselves).

In a typical setting, the classifier \tilde{h} is a deep neural network, f is the representation in the last internal layer of the network (i.e., the penultimate, a.k.a. *embedding* layer), and \tilde{g} , the last layer of the network, is a linear map; similarly, g in

the target problem is often taken to be linear. The learned feature map f is called a *foundation model* (Bommasani et al., 2021) when it can be effectively used in a wide range of target tasks. In particular, its effectiveness in dealing with downstream tasks (for any f) can be measured by its expected *transfer error*,

$$\mathcal{L}_{\mathcal{D}}(f) := \mathbb{E}_P \mathbb{E}_S [L_P(h_{f,S})], \quad (2)$$

where the expectation is taken over the random choice of P (i.e., averaging over randomly selected target tasks) and the training data S for the target class.

Notice that while the feature map f is evaluated on the distribution of target tasks P determined by \mathcal{D} , the training of f in a foundation model, as described above, is fully agnostic of this target. In the rest of the paper we analyze this setting, and provide an explanation through the recent concept of neural collapse.

Notation. $\|\cdot\|$ denotes the Euclidean norm for vectors and the spectral norm for matrices. For a distribution Q over $\mathcal{X} \subset \mathbb{R}^d$ and $u : \mathcal{X} \rightarrow \mathbb{R}^p$, the mean and variance of $u(x)$ for $x \sim Q$ are denoted by $\mu_u(Q) := \mathbb{E}_{x \sim Q}[u(x)]$ and by $\text{Var}_u(Q) := \mathbb{E}_{x \sim Q}[\|u(x) - \mu_u(Q)\|^2]$. For a finite set $A = \{a_i\}_{i=1}^n$, we denote $\text{Avg}_{i=1}^n[a_i] := \frac{1}{n} \sum_{i=1}^n a_i$ and by $U[A]$ the uniform distribution over A .

3. Neural Collapse

Neural collapse (NC) is a recently discovered phenomenon in deep learning (Papayan et al., 2020): it has been observed that during the training of deep networks for standard classification tasks, the features (the output of the penultimate, a.k.a. embedding layer) associated with training samples belonging to the same class concentrate around the mean feature value for the same class, also satisfying some additional conditions (such as the mean feature vectors being orthogonal to each other).

From these we focus on the class-features variance collapse (called NC1 by Papayan et al., 2020), and used a simplified version introduced by Galanti et al. (2022): For a feature map $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and two distributions Q_1, Q_2 over $\mathcal{X} \subset \mathbb{R}^d$, we define their *class-distance normalized variance* (CDNV) as

$$V_f(Q_1, Q_2) := \frac{\text{Var}_f(Q_1)}{\|\mu_f(Q_1) - \mu_f(Q_2)\|^2}$$

(the definition can be extended to finite sets $S_1, S_2 \subset \mathcal{X}$ by defining $V_f(S_1, S_2) = V_f(U[S_1], U[S_2])$). This quantity essentially measures to what extent the feature vectors of samples from Q_1 and Q_2 are clustered in space. We say that a training process satisfies *neural collapse* if

$$\lim_{t \rightarrow \infty} \text{Avg}_{i \neq j \in [l]} [V_{f^{(t)}}(\tilde{S}_i, \tilde{S}_j)] = 0, \quad (3)$$

where $f^{(t)}$ is the penultimate layer of a neural network $h^{(t)} = g^{(t)} \circ f^{(t)}$ that is trained to fit \tilde{S} for t steps. It has

been observed (Papayan et al., 2020; Galanti et al., 2022) that for several problems the left-hand side of (3) indeed reduces to some small value during training.

4. Generalization Bounds for Transfer Learning

Galanti et al. (2022) introduced generalization bounds for transfer learning in the setting of Section 2. Their analysis takes three main steps; first, they show that if the CDNV on the training data, $\text{Avg}_{i \neq j \in [l]} [V_f(\tilde{S}_i, \tilde{S}_j)]$, is small at the end of training, then we expect that the CDNV over unseen samples, $\text{Avg}_{i \neq j \in [l]} [V_f(\tilde{P}_i, \tilde{P}_j)]$, would also be small for the same classes (see Lemma 4.1 below). As a second step, they bound the expected CDNV between two new target classes, $\mathbb{E}_{P_c \neq P_{c'}} [V_f(P_c, P_{c'})]$, using the CDNV between the source classes, $\text{Avg}_{i \neq j \in [l]} [V_f(\tilde{P}_i, \tilde{P}_j)]$. Finally, they prove that $\mathcal{L}_{\mathcal{D}}(f) \stackrel{\times}{\leq} (k + k/n) \cdot \mathbb{E}_{P_c \neq P_{c'}} [V_f(P_c, P_{c'})]$, with $h_{f,S} := \arg \min_{c \in [C]} \|f(x_j) - \mu_f(S_c)\|$ being the *nearest class-center (NCC) classifier*.

However, there are multiple limitations to their arguments. First, the term $\mathbb{E}_{P_c \neq P_{c'}} [V_f(P_c, P_{c'})]$ may be very large even when $\mathcal{L}_{\mathcal{D}}(f)$ is very small. For example, if there is at least one anomalous pair of class-conditionals $P_1, P_2 \in \mathcal{E}$ for which $\mu_f(P_1) = \mu_f(P_2)$, then, $\mathbb{E}_{P_c \neq P_{c'}} [V_f(P_c, P_{c'})] = \infty$ (as long as $\mathcal{D}[P_1], \mathcal{D}[P_2] > 0$) even if for all other pairs $P'_1, P'_2 \in \mathcal{E}$, $V_f(P'_1, P'_2)$ is tiny. In fact, if the support \mathcal{E} of \mathcal{D} is infinite and f is bounded, then clearly $\inf_{P_1 \neq P_2 \in \mathcal{E}} \|\mu_f(P_1) - \mu_f(P_2)\| = 0$ and $\mathbb{E}_{P_c \neq P_{c'}} [V_f(P_c, P_{c'})] = \infty$. Second, the bound of Galanti et al. (2022) on $\mathbb{E}_{P_c \neq P_{c'}} [V_f(P_c, P_{c'})]$ scales with $(\inf_{f \in \mathcal{F}} \inf_{P_1 \neq P_2 \in \mathcal{E}} \|\mu_f(P_1) - \mu_f(P_2)\|)^{-1}$, which may be very large even if \mathcal{E} is finite (e.g., if \mathcal{F} contains a constant function).

In this paper we aim to circumvent these issues and provide an upper bound on the error of h which does not suffer from these limitations. Instead of bounding the term $\mathbb{E}_{P_c \neq P_{c'}} [V_f(P_c, P_{c'})]$, we bound $\mathcal{L}_{\mathcal{D}}(f)$ in terms of the averaged margin error of the NCC classifier between pairs of source classes. As a next step, in Lemma 4.3, we show that each of these error terms can be bounded using the CDNV between pairs of source classes. Finally, in Theorem 4.4 we combine Lemmas 4.1-4.3 and obtain an upper bound on the transfer error of f in terms of $\text{Avg}_{i \neq j \in [l]} [V_f(\tilde{S}_i, \tilde{S}_j)]$ which is typically small (see Section 3).

4.1. Neural Collapse Generalizes to New Samples

We start by recalling Proposition 1 of Galanti et al. (2022). This proposition provides an upper bound on $V_f(\tilde{P}_i, \tilde{P}_j)$, decomposed into $V_f(\tilde{S}_i, \tilde{S}_j)$ and additional generalization gap terms bounding the difference between expectations

and empirical averages for f and its variants, where f is the output of the learning algorithm with access to \tilde{S} .

For any $\delta \in (0, 1)$, we let $\epsilon_1^c(\delta)$ and $\epsilon_2^c(\delta)$ be the smallest positive values such that with probability at least $1 - \delta$, the learning algorithm returns a function $f \in \mathcal{F}$ that satisfies

$$\begin{aligned} \left| \mathbb{E}_{x \sim \tilde{P}_c} [f(x)] - \text{Avg}_{x \in \tilde{S}_c} [f(x)] \right| &\leq \epsilon_1^c(\delta); \\ \left| \mathbb{E}_{x \sim \tilde{P}_c} [\|f(x)\|^2] - \text{Avg}_{x \in \tilde{S}_c} [\|f(x)\|^2] \right| &\leq \epsilon_2^c(\delta). \end{aligned}$$

These quantities are typically bounded using Rademacher complexities (Bartlett & Mendelson, 2002) related to \mathcal{F} , scaling usually as $\mathcal{O}_m(1/\sqrt{m})$ (as also shown by Galanti et al., 2022). Next, we present their bound on the CDNV of the source distributions.

Lemma 4.1 (Galanti et al., 2022, Proposition 1). *Fix two source classes, i and j with distributions \tilde{P}_i and \tilde{P}_j , and let $\delta \in (0, 1)$. Let $\tilde{S}_c \sim \tilde{P}_c^m$ for $c \in \{i, j\}$. Let*

$$\begin{aligned} A_{ij}(\delta) &:= \frac{\epsilon_1^i(\delta/4) + \epsilon_1^j(\delta)}{\|\mu_f(\tilde{P}_i) - \mu_f(\tilde{P}_j)\|} \\ B_{ij}(\delta) &:= \frac{\text{Avg}_{c \in \{i, j\}} [\epsilon_2^c(\delta) + 2\|\mu_f(\tilde{P}_c)\| \cdot \epsilon_1^c(\delta) + \epsilon_1^c(\delta)^2]}{\|\mu_f(\tilde{S}_i) - \mu_f(\tilde{S}_j)\|^2}. \end{aligned}$$

Then, with probability at least $1 - \delta$ over the selection of \tilde{S}_i, \tilde{S}_j , we have $V_f(\tilde{P}_i, \tilde{P}_j) \leq (V_f(\tilde{S}_i, \tilde{S}_j) + B_{ij}(\delta/4)) \cdot (1 + A_{ij}(\delta/4))^2$.

The bound in the lemma has several terms. The first one is the empirical CDNV, $V_f(\tilde{S}_i, \tilde{S}_j)$, which is assumed to be implicitly minimized by the training algorithm (see Section 3). The rest of the terms are proportional to the generalization gaps $\epsilon_1^c(\delta/4)$ and $\epsilon_2^c(\delta/4)$ — as discussed above, typically we expect these terms to scale as $\mathcal{O}_m(1/\sqrt{m})$. As discussed in (Galanti et al., 2022), the denominators in $A_{ij}(\delta/4)$ and $B_{ij}(\delta/4)$ are expected to be lower bounded by positive constants under reasonable conditions.

Therefore, for a pair $i \neq j$, if the CDNV $V_f(\tilde{S}_i, \tilde{S}_j)$ is small and m is large, then we expect $V_f(\tilde{P}_i, \tilde{P}_j)$ to be small, as well. That is, if neural collapse emerges in the training data of two source classes, we should also expect it to emerge in unseen samples of the same classes (for large m). For an extensive empirical validation of this argument, see the works of Galanti et al. (2022); Galanti & Galanti (2022).

4.2. Bounding the Transfer Error

In this section we bound the transfer error of f using the averaged margin error of the NCC classifier between pairs of source classes. In the analysis, similarly to Galanti et al. (2022), we treat class-conditional distributions as data points on which the feature map f is trained. We apply standard techniques to derive generalization bounds to new data points, which, in this case, are class-conditional distribu-

tions. Recall that, in line with this view, we have already assumed that the source and target class-conditionals $\{\tilde{P}_c\}_{c=1}^l$ and $\{P_c\}_{c=1}^k$ are i.i.d. samples from \mathcal{D} .

For simplicity, we focus on the case where \mathcal{F} , the class of feature maps is a set of depth- q ReLU neural networks of the form $f(\cdot) = W^q \sigma(W^{q-1} \dots \sigma(W^1 \cdot)) : \mathbb{R}^d \rightarrow \mathbb{R}^p$, where $\sigma(x) := \max\{0, x\}$ is the element-wise ReLU function, and $W^i \in \mathbb{R}^{d_{i+1} \times d_i}$ for $i \in [q]$, where $d_1 = d$ and $d_{q+1} = p$. Our bounds depend on the spectral complexity of a network f , which is defined as $\mathcal{C}(f) := \max_{j \in [p]} \|W_j^q\| \cdot \prod_{r=1}^{q-1} \|W^r\|$. This quantity upper bounds the Lipschitz constant of f and is similar in fashion to other (slightly different) notions of spectral complexity for neural networks (Golowich et al., 2017; Bartlett et al., 2017). We also denote $B := \sup_{x \in \mathcal{X}} \|x\|$.

We start by a generic result that allows us to upper bound the transfer error of the NCC classifier (on top of the feature map f) by the margin-error over the source classes (defined below), plus some additional terms, depending on the complexity of f and the number of classes. The proof, presented in Appendix A, is partially based on Corollary 3 of Maurer & Pontil (2019) and the analysis of Bartlett et al. (2017).

Lemma 4.2. *Let $\delta \in (0, 1)$, $l \geq 2$ and let \mathcal{F} be a class of ReLU neural networks of depth q . Then, with probability at least $1 - \delta$ over the random selection of the class-conditionals $\tilde{\mathcal{P}} = \{\tilde{P}_c\}_{c=1}^l$ of the source classes, for any $f \in \mathcal{F}$ and $\Delta > 0$, we have*

$$\mathcal{L}_{\mathcal{D}}(f) \leq (k-1) \cdot \left[\text{Avg}_{i \neq j \in [l]} [L_{ij}^{2\Delta}(f)] + \delta \right] + \frac{\sqrt{q}pkB \cdot [\mathcal{C}(f)]}{\Delta\sqrt{l}} \cdot \text{polylog}(k, \delta^{-1}, [\mathcal{C}(f)], \lceil \log(\frac{\Delta}{B}) \rceil, l)$$

where $L_{ij}^{\Delta}(f) := \Pr_{\hat{x}_i, \hat{S}_i, \hat{S}_j} [\|f(\hat{x}_i) - \mu_f(\hat{S}_j)\| \leq \|f(\hat{x}_i) - \mu_f(\hat{S}_i)\| + \Delta]$, where $\hat{S}_i \sim \tilde{P}_i^n$, $\hat{S}_j \sim \tilde{P}_j^n$ and $\hat{x}_i \sim \tilde{P}_i$.

The bound above can be decomposed into several parts. The first term is the average expected Δ -margin error of a NCC classifier, when averaged over all classification problems given by the pairs of source classes with n samples each. For each pair (i, j) , the NCC classifier is defined with two datasets of size n , \hat{S}_i and \hat{S}_j (which are independent of the training datasets \tilde{S}_i and \tilde{S}_j), and its performance is evaluated over a new test sample $\hat{x}_i \sim \tilde{P}_i$. The second term, δ , is a free parameter and can be selected to be very small as the third term scales polylogarithmically with respect to $1/\delta$. The third term is proportional to the ratio between the complexity of the selected function f , captured by $\sqrt{q}pkB \cdot [\mathcal{C}(f)]$ and $\Delta\sqrt{l}$. Therefore, as long as we increase the number of source classes l , we can expect the generalization to new classes to improve.

The next lemma is an extension of Proposition 5 of Galanti et al. (2022) and relates the margin-error to the CDNV.

Lemma 4.3. *Let $\{\tilde{P}_i, \tilde{P}_j\}$ be a pair of class-conditional distributions. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$ be any feature map and denote $\mu_c = \mu_f(\tilde{P}_c)$. Let $\Delta \leq 0.12s(f, \tilde{\mathcal{P}})^{-1/2} \cdot \|\mu_f(\tilde{P}_i) - \mu_f(\tilde{P}_j)\|$. Then, we have, for some universal constant $C > 0$,*

$$\frac{1}{2}(L_{ij}^{\Delta}(f) + L_{ji}^{\Delta}(f)) \leq C \left[\frac{1}{s(f, \tilde{\mathcal{P}})} + \frac{1}{n} \right] \cdot V_f(\tilde{P}_i, \tilde{P}_j),$$

where $s(f, \tilde{\mathcal{P}}) = p$ if $\{f \circ \tilde{P}_c\}_{c=1}^l$ are spherically symmetric and $s(f, \tilde{\mathcal{P}}) = 1$ otherwise.

Informally, the expected margin error $L_{ij}^{\Delta}(f)$ of f is upper bounded by $V_f(\tilde{P}_i, \tilde{P}_j)$, and when $f \circ \tilde{P}_i$ and $f \circ \tilde{P}_j$ are spherically symmetric then by $(1/p + 1/n) \cdot V_f(\tilde{P}_i, \tilde{P}_j)$. Therefore, in case of neural collapse (i.e., when the $V_f(\tilde{P}_i, \tilde{P}_j)$ are small on average), $L_{ij}^{\Delta}(f)$ is small, explaining the success of foundation models in the low-data regime (such as in few-shot learning problems) in the presence of neural collapse. Finally, we also derive a much better bound (see Lemma B.2) which is exponentially small in $p/V_f(\tilde{P}_1, \tilde{P}_2)$ if $f \circ \tilde{P}_c$ are (assumed to be) centered Gaussian distributions as also assumed by Papyan et al. (2020).

Putting together Lemmas 4.1-4.3, we obtain the following generalization bound.

Theorem 4.4. *Let $\delta \in (0, 1)$, $l \geq 2$ and let \mathcal{F} be a class of ReLU neural networks of depth q . Then, with probability at least $1 - \delta$ over the selection of the source class-conditionals $\tilde{\mathcal{P}} = \{\tilde{P}_c\}_{c=1}^l$ and the corresponding source training data $\tilde{S}_1, \dots, \tilde{S}_l$, for any $f \in \mathcal{F}$ and $\Delta = 0.12s(f, \tilde{\mathcal{P}})^{-1/2} \cdot \min_{i \neq j \in [l]} \|\mu_f(\tilde{P}_i) - \mu_f(\tilde{P}_j)\|$, we have*

$$\mathcal{L}_{\mathcal{D}}(f) \leq \frac{k\delta}{2} + Ck \cdot \left[\frac{1}{s(f, \tilde{\mathcal{P}})} + \frac{1}{n} \right] \cdot \text{Avg}_{i \neq j \in [l]} \left[(V_f(\tilde{S}_i, \tilde{S}_j) + B_{ij}(\frac{\delta}{8l^2})) \cdot (1 + A_{ij}(\frac{\delta}{8l^2}))^2 \right] + \frac{\sqrt{q}pkB \cdot [\mathcal{C}(f)]}{\Delta\sqrt{l}} \cdot \text{polylog}(k, \delta^{-1}, [\mathcal{C}(f)], \lceil \log(\frac{\Delta}{B}) \rceil, l)$$

for some universal constant $C > 0$.

Namely, the transfer error of f is bounded by the sum of the averaged CDNV over the source training data, $V_f(\tilde{S}_i, \tilde{S}_j)$, combinations of the terms $A_{ij}(\frac{\delta}{8l^2})$ and $B_{ij}(\frac{\delta}{8l^2})$, where $\epsilon_1^c(\frac{\delta}{8l^2}) = \tilde{\mathcal{O}}\left(\frac{p\sqrt{q} \cdot B \cdot [\mathcal{C}(f)] \cdot \log(l/\delta)}{\sqrt{m}}\right)$ and $\epsilon_2^c(\frac{\delta}{8l^2}) = \tilde{\mathcal{O}}\left(\frac{p\sqrt{q} \cdot B^2 \cdot [\mathcal{C}(f)]^2 \cdot \log(l/\delta)}{\sqrt{m}}\right)$ (see Galanti et al., 2022) and a term that scales as $\tilde{\mathcal{O}}\left(\frac{\sqrt{q}pkB \cdot [\mathcal{C}(f)]}{\Delta\sqrt{l}}\right)$. Therefore, if l and m are large and we have neural collapse across the source training samples (i.e., $\text{Avg}_{i \neq j \in [l]} [V_f(\tilde{S}_i, \tilde{S}_j)]$ is small) and $[\mathcal{C}(f)]$ is bounded (as a function of m, l), then the transfer error $\mathcal{L}_{\mathcal{D}}(f)$ is also small. Interestingly, in the presence of neural collapse in the source data, the dependence on n is very weak.

5. Conclusions

Using pretrained models for transfer learning is a successful approach in the low-data regime. Recently, Galanti et al. (2022) suggested a new perspective on analysing this problem via the newly discovered phenomenon of neural collapse. However, their theoretical analysis relied on several unrealistic assumptions. In this work we provided better theoretical analysis by significantly relaxing their assumptions.

6. Acknowledgements

Tomer Galanti was supported by the Center for Minds, Brains and Machines (CBMM), funded by NSF STC award CCF-1231216 and by NSF award 213418.

References

- Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *J. Mach. Learn. Res.*, 3:463–482, 2002.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. Spectrally-normalized margin bounds for neural networks. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, pp. 6241–6250, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- Bengio, Y. Deep learning of representations for unsupervised and transfer learning. In *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*, volume 27 of *Proceedings of Machine Learning Research*, pp. 17–36, Bellevue, Washington, USA, 02 Jul 2012. PMLR.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A. S., Creel, K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N. D., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P. W., Krass, M. S., Krishna, R., Kudithipudi, R., and et al. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258, 2021.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.
- Caruana, R. Learning many related tasks at the same time with backpropagation. In *Advances in Neural Information Processing Systems*, volume 7. MIT Press, 1995.
- Dhillon, G. S., Chaudhari, P., Ravichandran, A., and Soatto, S. A baseline for few-shot image classification. In *International Conference on Learning Representations*, 2020.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1126–1135. PMLR, 06–11 Aug 2017.
- Galanti, T. and Galanti, L. On the implicit bias towards minimal depth of deep neural networks, 2022.
- Galanti, T., György, A., and Hutter, M. On the role of neural collapse in transfer learning. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=SwIp410B6aQ>.
- Golowich, N., Rakhlin, A., and Shamir, O. Size-independent sample complexity of neural networks. *Information and Inference: A Journal of the IMA*, 9, 12 2017. doi: 10.1093/imaiai/iaz007.
- Lee, K., Maji, S., Ravichandran, A., and Soatto, S. Meta-learning with differentiable convex optimization. In *CVPR*, 2019.
- Maurer, A. and Pontil, M. Uniform concentration and symmetrization for weak interactions. In *COLT*, 2019.
- Papayan, V., Han, X. Y., and Donoho, D. L. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.
- Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., and Sutskever, I. Zero-shot text-to-image generation. *CoRR*, abs/2102.12092, 2021.
- Ravi, S. and Larochelle, H. Optimization as a model for few-shot learning. In *International Conference on Learning Representations*, 2017.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.

Tian, Y., Wang, Y., Krishnan, D., Tenenbaum, J. B., and Isola, P. Rethinking few-shot image classification: A good embedding is all you need? In *Proceedings of the European Conference on Computer Vision (ECCV)*, volume 12359 of *Lecture Notes in Computer Science*, pp. 266–282. Springer, 2020.

Vinyals, O., Blundell, C., Lillicrap, T., kavukcuoglu, k., and Wierstra, D. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.

Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.