

---

# Distributional Generalization: Characterizing Classifiers Beyond Test Error

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We present a new set of empirical properties of interpolating classifiers, including  
2 neural networks, kernel machines and decision trees. Informally, the output  
3 distribution of an interpolating classifier matches the distribution of true labels,  
4 when conditioned on certain subgroups of the input space. For example, if we  
5 mislabel 30% of dogs as cats in the train set of CIFAR-10, then a ResNet trained  
6 to interpolation will in fact mislabel roughly 30% of dogs as cats on the *test set*  
7 as well, while leaving other classes unaffected. These behaviors are not captured  
8 by classical generalization, which would only consider the average error over  
9 the inputs, and not *where* these errors occur. We introduce and experimentally  
10 validate a formal conjecture that specifies the subgroups for which we expect this  
11 distributional closeness. Further, we show that these properties can be seen as a  
12 new form of generalization, which advances our understanding of the implicit bias  
13 of interpolating methods.

## 14 1 Introduction

15 In learning theory, when we study how well a classifier “generalizes”, we usually consider a single  
16 metric – its test error [59]. However, there could be many different classifiers with the same test error  
17 that differ substantially in, say, the subgroups of inputs on which they make errors or in the features  
18 they use to attain this performance. Reducing classifiers to a single number misses these rich aspects  
19 of their behavior. In this work, we propose formally studying the entire *joint distribution* of classifier  
20 inputs and outputs. That is, the distribution  $(x, f(x))$  for samples from the distribution  $x \sim D$  for a  
21 classifier  $f(x)$ . This distribution reveals many structural properties of the classifier beyond test error  
22 (such as *where* the errors occur). In fact, we discover new behaviors of modern classifiers that can  
23 only be understood in this framework. As an example, consider the following experiment (Figure 1).

24 **Experiment 1.** Consider a binary classification version of CIFAR-10, where CIFAR-10 images  $x$   
25 have binary labels *Animal/Object*. Take 50K samples from this distribution as a train set, but  
26 apply the following label noise: flip the label of cats to *Object* with probability 30%. Now train  
27 a WideResNet  $f$  to 0 train error on this train set. How does the trained classifier behave on test  
28 samples? Options below:

- 29 (1) The test error is low across all classes, since there is only 3% overall label noise in the train set.  
30 (2) Test error is “spread” across the animal class. After all, the classifier is not explicitly told what a  
31 cat or a dog is, just that they are all animals.  
32 (3) The classifier misclassifies roughly 30% of test cats as “objects”, but all other animals are largely  
33 unaffected.

34 The reality is closest to option (3) as shown in Figure 1. The left panel shows the joint density of  
35 train inputs  $x$  with train labels *Object/Animal*. Since the classifier is interpolating, the classifier

36 outputs on the train set are identical to the left panel. The right panel shows the *classifier predictions*  
 37  $f(x)$  on *test inputs*  $x$ .

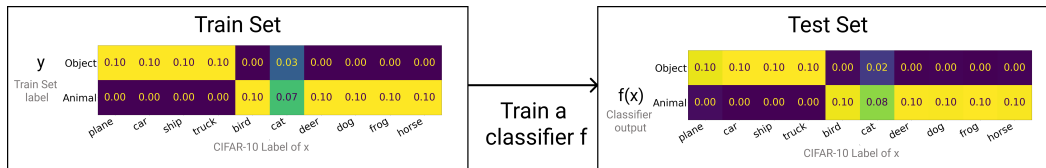


Figure 1: The setup and result of Experiment 1. The CIFAR-10 train set is labeled as either Animals or Objects, with label noise affecting only cats. A WideResNet-28-10 is then trained to 0 train error on this train set, and evaluated on the test set. Full experimental details in Appendix C.2

38 There are several notable things about this experiment. First, the error is *localized* to cats in the test  
 39 set as it was in the train set, even though no explicit cat labels were provided. The interpolating  
 40 model is thus sensitive to subgroup-structures in the distribution. Second, the *amount* of error on  
 41 the cat class is close to the noise applied on the train set. Thus, the behavior of the classifier on the  
 42 train set *generalizes* to the test set in a stronger sense than just average error. Specifically, when  
 43 *conditioned on a subgroup* (cat), the *distribution* of the true labels is close to that of the classifier  
 44 outputs. Third, this is not the behavior of the Bayes-optimal classifier, which would always output  
 45 the maximum-likelihood label instead of reproducing the noise in the distribution. The network  
 46 is thus behaving poorly from the perspective of Bayes-optimality, but behaving well in a certain  
 47 distributional sense (which we will formalize soon).

48 Now, consider a seemingly unrelated experimental observation. Take an AlexNet trained on ImageNet,  
 49 a 1000-way classification problem with 116 varieties of dogs. AlexNet only achieves 56.5% test  
 50 accuracy on ImageNet. However, it at least classifies most dogs as *some* variety of dog (with 98.4%  
 51 accuracy), though it may mistake the exact breed. In this work, we show that both of these experiments  
 52 are examples of the same underlying phenomenon. We empirically show that for an interpolating  
 53 classifier, its classification outputs are close in distribution to the true labels — even when conditioned  
 54 on many subsets of the domain. For example, in Figure 1, the distribution of  $p(f(x)|x = \text{cat})$  is close  
 55 to the true label distribution of  $p(y|x = \text{cat})$ . We propose a formal conjecture (Feature Calibration),  
 56 that predicts which subgroups of the domain can be conditioned on for the above distributional  
 57 closeness to hold.

58 These experimental behaviors could not have been captured solely by looking at average test error,  
 59 as is done in the classical theory of generalization. In fact, they are special cases of a new kind of  
 60 generalization, which we call “Distributional Generalization”.

## 61 1.1 Distributional Generalization

62 Informally, Distributional Generalization states that the outputs of classifiers  $f$  on their train sets  
 63 and test sets are close *as distributions* (as opposed to close in just error). That is, the following joint  
 64 distributions are close:

$$(x, f(x))_{x \sim \text{TestSet}} \approx (x, f(x))_{x \sim \text{TrainSet}} \quad (1)$$

65 The remainder of this paper is devoted to making the above statement precise, and empirically  
 66 checking its validity on real-world tasks. Specifically, we want to formally define the notion of  
 67 approximation ( $\approx$ ), and understand how it depends on the problem parameters (the type of classifier,  
 68 number of train samples, etc). We focus primarily on interpolating methods, where we formalize  
 69 Equation (1) through our Feature Calibration Conjecture.

## 70 1.2 Our Contributions and Organization

71 In this work, we discover new empirical properties of interpolating classifiers, which are not captured  
 72 in the classical framework of generalization. We then propose formal conjectures to characterize  
 73 these behaviors.

<sup>1</sup>These distributions also include the randomness in sampling the train and test sets, and in training the classifier, as we define more precisely in Section 3

- 74 • In Section 3, we introduce a formal “Feature Calibration” conjecture, which unifies our  
75 experimental observations. Roughly, Feature Calibration says that the outputs of classifiers  
76 match the statistics of their training distribution when conditioned on certain subgroups.
- 77 • In Section 4, we experimentally stress test our Feature Calibration conjecture across various  
78 settings in machine learning, including neural networks, kernel machines, and decision trees.  
79 This highlights the universality of our results across machine learning.
- 80 • In Section 5, we relate our results to classical generalization, by defining a new notion of  
81 Distributional Generalization which subsumes both classical generalization and our new  
82 conjectures.
- 83 • Finally, in Section 5.2 we informally discuss how Distributional Generalization can be  
84 applied even for non-interpolating methods.

85 Our results, thus, extend our understanding of the *implicit bias* of interpolating methods, and introduce  
86 a new type of generalization exhibited across many methods in machine learning.

### 87 1.3 Related Work and Significance

88 Our work has connections to, and implications for many existing research programs in deep learning.

89 **Implicit Bias and Overparameterization.** There has been a long line of recent work towards  
90 understanding overparameterized and interpolating methods, since these pose challenges for classical  
91 theories of generalization (e.g. Belkin et al. [8, 9, 10], Breiman [11], Gunasekar et al. [25], Liang  
92 and Rakhlin [36], Nakkiran et al. [43], Schapire et al. [58], Soudry et al. [62], Zhang et al. [71]). The  
93 “implicit bias” program here aims to answer: *Among all models with 0 train error, which model is*  
94 *actually produced by SGD?* Most existing work seeks to characterize the exact implicit bias of models  
95 under certain (sometimes strong) assumptions on the model, training method or the data distribution.  
96 In contrast, our conjecture applies across many different interpolating models (from neural nets to  
97 decision trees) as they would be used in practice, and thus form a sort of “universal implicit bias” of  
98 these methods. Moreover, our results place constraints on potential future theories of implicit bias,  
99 and guide us towards theories that better capture practice.

100 **Benign Overfitting.** Most prior works on interpolating classifiers attempt to explain why training  
101 to interpolation “does not harm” the the model. This has been dubbed “benign overfitting” [7] and  
102 “harmless interpolation” [40], reflecting the widely-held belief that interpolation does not harm the  
103 decision boundary of classifiers. In contrast, we find that interpolation actually does “harm” classifiers,  
104 in predictable ways: fitting the label noise on the train set causes similar noise to be reproduced at  
105 test time. Our results thus indicate that interpolation can significantly affect the decision boundary of  
106 classifiers, and should not be considered a purely “benign” effect.

107 **Classical Generalization and Scaling Limits.** Our framework of Distributional Generalization is  
108 insightful even to study classical generalization, since it reveals much more about models than just  
109 their test error. For example, statistical learning theory attempts to understand if and when models  
110 will asymptotically converge to Bayes optimal classifiers, in the limit of large data (“asymptotic  
111 consistency” [59, 65]). In deep learning, there are at least two distinct ways to scale model and data  
112 to infinity together: the *underparameterized* scaling limit, where data-size  $\gg$  model-size always, and  
113 the *overparameterized* scaling limit, where data-size  $\ll$  model-size always. The underparameterized  
114 scaling limit is well-understood: when data is essentially infinite, neural networks will converge to  
115 the Bayes-optimal classifier (provided the model-size is large enough, and the optimization is run  
116 for long enough, with enough noise to escape local minima). On the other hand, our work suggests  
117 that in the *overparameterized* scaling limit, models will *not* converge to the Bayes-optimal classifier.  
118 Specifically, our Feature Calibration Conjecture implies that in the limit of large data, interpolating  
119 models will approach a *sampler* from the distribution. That is, the limiting model  $f$  will be such that  
120 the output  $f(x)$  is a sample from  $p(y|x)$ , as opposed to the Bayes-optimal  $f^*(x) = \operatorname{argmax}_y p(y|x)$ .  
121 This claim—that overparameterized models do not converge to Bayes-optimal classifiers—is unique  
122 to our work as far as we know, and highlights the broad implications of our results.

123 **Locality and Manifold Learning.** Our intuition for the behaviors in this work is that they arise due to  
124 some form of “locality” of the trained classifiers, in an appropriate embedding space. For example, the  
125 behavior observed in Experiment 1 would be consistent with that of a 1-Nearest-Neighbor classifier  
126 in a embedding that separates the CIFAR-10 classes well. This intuition that classifiers learn good

127 embeddings is present in various forms in the literature, for example: the so-called called “manifold  
 128 hypothesis,” that natural data lie on a low-dimensional manifold [44, 61], as well as works on local  
 129 stiffness of the loss landscape [19], and works showing that overparameterized neural networks can  
 130 learn hidden low-dimensional structure in high-dimensional settings [6, 15, 21]. It is open to more  
 131 formally understand connections between our work and the above.

132 **Other Related Works.** Our conjectures also describe neural networks under label noise, which has  
 133 been empirically and theoretically studied in the past [9, 14, 45, 54, 63, 71, 72], though not formally  
 134 characterized. A full discussion of related works is in Appendix A

## 135 2 Preliminaries

136 **Notation.** We consider joint distributions  $\mathcal{D}$  on  $x \in \mathcal{X}$  and discrete  $y \in \mathcal{Y} = [k]$ . Let  $S =$   
 137  $\{(x_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$  denote a train set of  $n$  iid samples from  $\mathcal{D}$ . Let  $\mathcal{A}$  denote the training procedure  
 138 (including architecture and training algorithm for neural networks), and let  $f \leftarrow \text{Train}_{\mathcal{A}}(S)$  denote  
 139 training a classifier  $f$  on train-set  $S$  using procedure  $\mathcal{A}$ . We consider classifiers which output hard  
 140 decisions  $f : \mathcal{X} \rightarrow \mathcal{Y}$ . Let  $\text{NN}_S(x) = x_i$  denote the nearest-neighbor to  $x$  in train-set  $S$ , with  
 141 respect to a distance metric  $d$ . Our theorems will apply to any distance metric, and so we leave  
 142 this unspecified. Let  $\text{NN}_S^{(y)}(x)$  denote the nearest-neighbor estimator itself, that is,  $\text{NN}_S^{(y)}(x) := y_i$   
 143 where  $x_i = \text{NN}_S(x)$ .

144 **Experimental Setup.** Briefly, we train all classifiers to interpolation (to 0 train error). Neural  
 145 networks (MLPs and ResNets [29]) are trained with SGD. Interpolating decision trees are trained  
 146 using the growth rule from Random Forests [12]. For kernel classification, we consider kernel  
 147 regression on one-hot labels and kernel SVM, with small or 0 of regularization (which is often  
 148 optimal [60]). Full experimental details are provided in Appendix B

149 **Distributional Closeness.** We consider the following notion of closeness for two probability dis-  
 150 tributions: For two distributions  $P, Q$  over  $\mathcal{X} \times \mathcal{Y}$ , let a “test” (or “distinguisher”) be a function  
 151  $T : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$  which accepts a sample from either distribution, and is intended to classify the  
 152 sample as either from distribution  $P$  or  $Q$ . For any set  $\mathcal{C} \subseteq \{T : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]\}$  of tests, we say  
 153 distributions  $P$  and  $Q$  are “ $\varepsilon$ -indistinguishable up to  $\mathcal{C}$ -tests” if they are close with respect to all tests  
 154 in class  $\mathcal{C}$ . That is,

$$P \approx_{\varepsilon}^{\mathcal{C}} Q \iff \sup_{T \in \mathcal{C}} \left| \mathbb{E}_{(x,y) \sim P} [T(x,y)] - \mathbb{E}_{(x,y) \sim Q} [T(x,y)] \right| \leq \varepsilon \quad (2)$$

155 Total-Variation distance is equivalent to closeness in all tests, i.e.  $\mathcal{C} = \{T : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]\}$ , but we  
 156 consider closeness for restricted families of tests  $\mathcal{C}$ .  $P \approx_{\varepsilon} Q$  denotes  $\varepsilon$ -closeness in TV-distance.

## 157 3 Feature Calibration Conjecture

### 158 3.1 Distributions of Interest

159 We first define three key distributions that we will use in stating our formal conjecture. For a given  
 160 data distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$  and training procedure  $\text{Train}_{\mathcal{A}}$ , we consider the following three  
 161 distributions over  $\mathcal{X} \times \mathcal{Y}$ :

- 162 1. **Source  $\mathcal{D}$ :**  $(x, y)$  where  $x, y \sim \mathcal{D}$ .
- 163 2. **Train  $\mathcal{D}_{\text{tr}}$ :**  $(x_{\text{tr}}, f(x_{\text{tr}}))$  where  $S \sim \mathcal{D}^n$ ,  $f \leftarrow \text{Train}_{\mathcal{A}}(S)$ ,  $(x_{\text{tr}}, y_{\text{tr}}) \sim S$
- 164 3. **Test  $\mathcal{D}_{\text{te}}$ :**  $(x, f(x))$  where  $S \sim \mathcal{D}^n$ ,  $f \leftarrow \text{Train}_{\mathcal{A}}(S)$ ,  $x, y \sim \mathcal{D}$

165 The source distribution  $\mathcal{D}$  is simply the original distribution. To sample once from the **Train Dis-**  
 166 **tribution  $\mathcal{D}_{\text{tr}}$** , we first sample a train set  $S \sim \mathcal{D}^n$ , train a classifier  $f$  on it, then output  $(x_{\text{tr}}, f(x_{\text{tr}}))$   
 167 for a random *train point*  $x_{\text{tr}} \in S$ . That is,  $\mathcal{D}_{\text{tr}}$  is the distribution of input and outputs of a trained  
 168 classifier  $f$  on its train set. To sample once from the **Test Distribution  $\mathcal{D}_{\text{te}}$** , we do this same proce-  
 169 dure, but output  $(x, f(x))$  for a random *test point*  $x$ . That is, the  $\mathcal{D}_{\text{te}}$  is the distribution of input and  
 170 outputs of a trained classifier  $f$  at test time. The only difference between the Train Distribution and

171 Test Distribution is that the point  $x$  is sampled from the train set or the test set, respectively<sup>2</sup>. For  
 172 interpolating classifiers,  $f(x_{\text{tr}}) = y_{\text{tr}}$  on the train set, and so the Source and Train distributions are  
 173 equivalent:  $\mathcal{D} \equiv \mathcal{D}_{\text{tr}}$ . (Note that these definitions, crucially, involve randomness from sampling the  
 174 train set, training the classifier, and sampling a test point).

### 175 3.2 Feature Calibration

We now formally describe the Feature Calibration Conjecture. At a high level, we argue that the  
 distributions  $\mathcal{D}_{\text{te}}$  and  $\mathcal{D}$  are statistically close for interpolating classifiers if we first “coarsen” the  
 domain of  $x$  by some partition  $L : \mathcal{X} \rightarrow [M]$  in to  $M$  parts. That is, for certain partitions  $L$ , the  
 following distributions are statistically close:

$$(L(x), f(x))_{x \sim \mathcal{D}} \approx_{\varepsilon} (L(x), y)_{x \sim \mathcal{D}}$$

176 We think of  $L$  as defining subgroups over the domain— for example,  $L(x) \in \{\text{dog, cat, horse} \dots\}$ .  
 177 Then, the above statistical closeness is essentially equivalent to requiring that for all subgroups  
 178  $\ell \in [M]$ , the conditional distribution of classifier output on the subgroup— $p(f(x)|L(x) = \ell)$  — is  
 179 close to the true conditional distribution:  $p(y|L(x) = \ell)$ .

180 The crux of our conjecture lies in defining exactly which subgroups  $L$  satisfy this distributional  
 181 closeness, and quantifying the  $\varepsilon$  approximation. This is subtle, since it must depend on almost all  
 182 parameters of the problem. For example, consider a modification to Experiment 1, where we use  
 183 a fully-connected network (MLP) instead of a ResNet. An MLP cannot properly distinguish cats  
 184 even when it is actually provided the real CIFAR-10 labels, and so (informally) it has no hope of  
 185 behaving differently on cats in the setting of Experiment 1, where the cats are not labeled explicitly  
 186 (See Figure C.2 for results with MLPs). Similarly, if we train the ResNet with very few samples from  
 187 the distribution, the network will be unable to recognize cats. Thus, the allowable partitions must  
 188 depend on the classifier family and the training method, including the number of samples.

189 We conjecture that allowable partitions are those which can themselves be learnt to good test  
 190 performance with an identical training procedure, but trained with the labels of the partition  $L$  instead  
 191 of  $y$ . To formalize this, we define a *distinguishable feature*: a partition of the domain  $\mathcal{X}$  that is  
 192 learnable for a given training procedure. Thus, in Experiment 1, the partition into CIFAR-10 classes  
 193 would be a distinguishable feature for ResNets (trained with SGD with 50K or more samples), but  
 194 not for MLPs. The definition below depends on the training procedure  $\mathcal{A}$ , the data distribution  $\mathcal{D}$ ,  
 195 number of train samples  $n$ , and an approximation parameter  $\varepsilon$  (which we think of as  $\varepsilon \approx 0$ ).

**Definition 1** ( $(\varepsilon, \mathcal{A}, \mathcal{D}, n)$ -Distinguishable Feature). *For a distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$ , number of  
 samples  $n$ , training procedure  $\mathcal{A}$ , and small  $\varepsilon \geq 0$ , an  $(\varepsilon, \mathcal{A}, \mathcal{D}, n)$ -distinguishable feature is a  
 partition  $L : \mathcal{X} \rightarrow [M]$  of the domain  $\mathcal{X}$  into  $M$  parts, such that training a model using  $\mathcal{A}$  on  $n$   
 samples labeled by  $L$  works to classify  $L$  with high test accuracy. Precisely,  $L$  is a  $(\varepsilon, \mathcal{A}, \mathcal{D}, n)$ -  
 distinguishable feature if:*

$$\Pr_{\substack{S = \{(x_i, L(x_i))\}_{x_1, \dots, x_n \sim \mathcal{D}} \\ f \leftarrow \text{Train}_{\mathcal{A}}(S); x \sim \mathcal{D}}} [f(x) = L(x)] \geq 1 - \varepsilon$$

196 This definition depends only on the marginal distribution of  $\mathcal{D}$  on  $x$ , and not on the label distribution  
 197  $p_{\mathcal{D}}(y|x)$ . To recap, this definition is meant to capture a labeling of the domain  $\mathcal{X}$  that is learnable for  
 198 a given training procedure  $\mathcal{A}$ . It must depend on the architecture used by  $\mathcal{A}$  and number of samples  
 199  $n$ , since more powerful classifiers can distinguish more features. Note that there could be many  
 200 distinguishable features for a given setting  $(\varepsilon, \mathcal{A}, \mathcal{D}, n)$  — including features not implied by the class  
 201 label such as the presence of grass in a CIFAR-10 image. Our main conjecture follows.

202 **Conjecture 1** (Feature Calibration). *For all natural distributions  $\mathcal{D}$ , number of samples  $n$ , interpo-  
 203 lating training procedures  $\mathcal{A}$ , and  $\varepsilon \geq 0$ , the following distributions are statistically close for all  
 204  $(\varepsilon, \mathcal{A}, \mathcal{D}, n)$ -distinguishable features  $L$ :*

$$(L(x), f(x))_{f \leftarrow \text{Train}_{\mathcal{A}}(\mathcal{D}^n); x, y \sim \mathcal{D}} \approx_{\varepsilon} (L(x), y)_{x, y \sim \mathcal{D}} \quad (3)$$

205 or equivalently:

$$(L(x), \hat{y})_{x, \hat{y} \sim \mathcal{D}_{\text{te}}} \approx_{\varepsilon} (L(x), y)_{x, y \sim \mathcal{D}} \quad (4)$$

<sup>2</sup>Technically, these definitions require training a fresh classifier for each sample, using independent train sets.  
 For practical reasons most of our experiments train a single classifier  $f$  and evaluate it on the entire train/test set.

206 This claims that the TV distance between the LHS and RHS of Equation (4) is at most  $\varepsilon$ , where  $\varepsilon$  is the  
 207 error of the distinguishable feature (in Definition 1). We claim that this holds *for all* distinguishable  
 208 features  $L$  “automatically” – we simply train a classifier, without specifying any particular partition.  
 209 The formal statements of Definition 1 and Conjecture 1 may seem somewhat arbitrary, involving  
 210 many quantifiers over  $(\varepsilon, \mathcal{A}, \mathcal{D}, n)$ . However, we believe these statements are natural: In addition  
 211 to extensive experimental evidence in Section 4, we also prove that Conjecture 1 is formally true as  
 212 stated for 1-Nearest-Neighbor classifiers in Theorem 1.

### 213 3.3 Feature Calibration for 1-Nearest-Neighbors

214 Here we prove that the 1-Nearest-Neighbor classifier formally satisfies Conjecture 1, under mild  
 215 assumptions. We view this theorem as support for our (somewhat involved) formalism of Conjecture 1.  
 216 Indeed, without Theorem 1 below, it is unclear if our statement of Conjecture 1 can ever be satisfied by  
 217 any classifier, or if it is simply too strong to be true. This theorem applies generically to a wide class  
 218 of distributions; the only assumption is a weak regularity condition: sampling the nearest-neighbor  
 219 train point to a random test point should yield (close to) a uniformly random test point.

220 **Theorem 1.** *Let  $\mathcal{D}$  be a distribution over  $\mathcal{X} \times \mathcal{Y}$ , and let  $n \in \mathbb{N}$  be the number of train samples.*  
 221 *Assume the following regularity condition holds: Sampling the nearest-neighbor train point to a*  
 222 *random test point yields (close to) a uniformly random test point. That is, suppose that for some*  
 223 *small  $\delta \geq 0$ , the distributions:  $\{\text{NN}_S(x)\}_{S \sim \mathcal{D}^n} \approx_{\delta} \{x\}_{x \sim \mathcal{D}}$ . Then, Conjecture 1 holds. That is,*  
 224 *for all  $(\varepsilon, \text{NN}, \mathcal{D}, n)$ -distinguishable partitions  $L$ , the following distributions are statistically close:*

$$\{(y, L(x))\}_{x, y \sim \mathcal{D}} \approx_{\varepsilon + \delta} \{(\text{NN}_S^{(y)}(x), L(x))\}_{S \sim \mathcal{D}^n, x, y \sim \mathcal{D}} \quad (5)$$

225 The proof of Theorem 1 is straightforward, and provided in Appendix D – but this strong property of  
 226 nearest-neighbors was not known before, to our knowledge.

### 227 3.4 Limitations: Natural Distributions

228 Technically, Conjecture 1 is not fully specified, since it does not specify exactly which classifiers or  
 229 distributions obey the conjecture. We do not claim that *all* classifiers and distributions satisfy our  
 230 conjectures. Nevertheless, we claim our conjectures hold in all “natural” settings, which informally  
 231 means settings with real data and classifiers that are actually used in practice. The problem of  
 232 understanding what separates “natural distributions” from artificial ones is not unique to our work,  
 233 and lies at the heart of deep learning theory. Many theoretical works handle this by considering  
 234 simplified distributional assumptions (e.g. smoothness, well-separatedness, gaussianity), which are  
 235 mathematically tractable, but untested in practice [2, 4, 35]. In contrast, we do not make untestable  
 236 mathematical assumptions. This benefit of realism comes at the cost of mathematical formalism.  
 237 We hope that as the theory of deep learning evolves, we will better understand how to formalize the  
 238 notion of “natural” in our conjectures.

## 239 4 Experiments: Feature Calibration

240 We now give empirical evidence for our conjecture in a variety of settings in machine learning,  
 241 including neural networks, kernel machines, and decision trees. In each experiment, we consider  
 242 a feature that is (verifiably) distinguishable, and then test our Feature Calibration conjecture for  
 243 this feature. Each of the experimental settings below highlights a different aspect of interpolating  
 244 classifiers, which may be of independent interest. Selected experiments are summarized here, with  
 245 full details and further experiments in Appendix C.

246 **Constant Partition:** Consider the trivially-distinguishable *constant* feature:  $L(x) = 0$  everywhere.  
 247 For this feature, Conjecture 1 reduces to the statement that the marginal distribution of class labels for  
 248 any interpolating classifier is close to the true marginals  $p(y)$ . To test this, we construct a variant of  
 249 CIFAR-10 with class-imbalance and train classifiers with varying levels of test errors to interpolation  
 250 on it. As shown in Figure 2B, the marginals of the classifier outputs are close to the true marginals,  
 251 even for a classifier that only achieves 37% test error.

252 **Coarse Partition:** Consider AlexNet trained on ILSVRC-2012 ImageNet [56], a 1000-class image  
 253 classification problem with 116 varieties of dogs. The network achieves only 56.5% accuracy

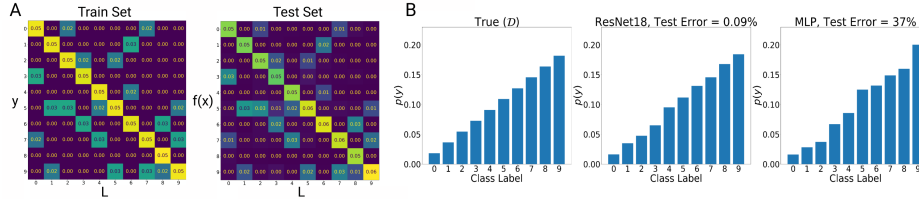


Figure 2: **Feature Calibration.** (A) Random confusion matrix on CIFAR-10, with a WideResNet28-10 trained to interpolation. Left: Joint density of labels  $y$  and original class  $L$  on the train set. Right: Joint density of classifier predictions  $f(x)$  and original class  $L$  on the test set. These two joint densities are close, as predicted by Conjecture 1. (B) Constant partition: The CIFAR-10 train set is class-rebalanced according to the left panel distribution. The center and right panels show that both ResNets and MLPs have the correct marginal distribution of outputs, even though the MLP has high test error.

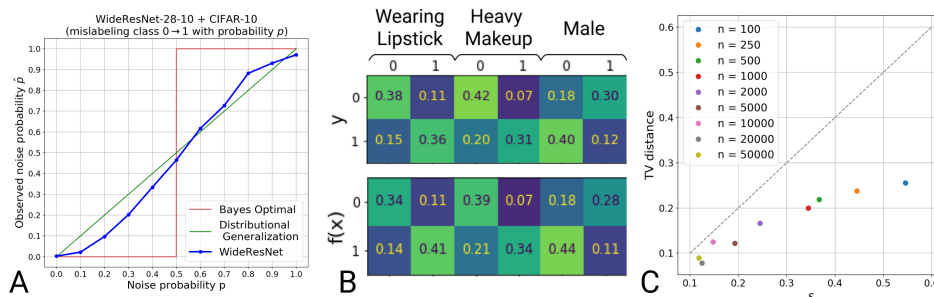


Figure 3: **Feature Calibration.** (A) CIFAR-10 with  $p$  fraction of class  $0 \rightarrow 1$  mislabeled on the train set. Plotting observed noise on classifier outputs vs. applied noise on the train set. (B) Multiple feature calibration on CelebA. (C) TV-distance between  $(L(x), f(x))$  and  $(L(x), y)$  for a variant of Experiment 1 with error on the distinguishable partitions ( $\epsilon$ ). The error was changed by changing the number of samples  $n$ .

254 on the test set. But it will at least classify most dogs as dogs (with 98.4% accuracy), making  
 255  $L(x) \in \{\text{dog}, \text{not-dog}\}$  a distinguishable feature. Moreover, as predicted by Conjecture 1,  
 256 the network is *calibrated* with respect to dogs: 22.4% of all dogs in ImageNet are Terriers, and indeed  
 257 the network classifies 20.9% of all dogs as Terriers (though it has 9% error on which specific dogs  
 258 it classifies as Terriers). See Appendix Table 2 for details, and related experiments on ResNets and  
 259 kernels in Appendix C.

260 **Class Partition:** We now consider settings where the class labels are themselves distinguishable  
 261 features (eg: CIFAR-10 classes are distinguishable by ResNets). Here our conjecture predicts the  
 262 behavior of interpolating classifiers under structured label noise. As an example, we generate a  
 263 random sparse confusion matrix and apply this to the labels of CIFAR-10 as shown in Figure 2A.  
 264 We find that a WideResNet trained to interpolation outputs the same confusion matrix on the test  
 265 set as well (Figure 2B). Now, to test that this phenomenon is indeed robust to the level of noise, we  
 266 mislabel class  $0 \rightarrow 1$  with probability  $p$  in the CIFAR-10 train set for varying levels of  $p$ . We then  
 267 observe  $\hat{p}$ , the fraction of samples mislabeled by this network from  $0 \rightarrow 1$  in the test set (Figure 3A  
 268 shows  $p$  versus  $\hat{p}$ ). The Bayes optimal classifier for this distribution behaves as a step function (in  
 269 red), and a classifier that obeys Conjecture 1 exactly would follow the diagonal (in green). The actual  
 270 experiment (in blue) is close to the behavior predicted by Conjecture 1. This experiment shows a  
 271 contrast with classical learning theory. While most existing theory focuses on whether classifiers  
 272 converge to the Bayes optimal solution, we show that interpolating classifiers behave “optimally” in a  
 273 different sense: they match the distribution of their train set. We discuss this further in Section 5. See  
 274 Appendix C.4 for more experiments, including other classifiers such as Decisions Trees.

275 **Multiple features:** Conjecture 1 states that the network should be automatically calibrated for  
 276 all distinguishable features, without any explicit labels for them. To do this, we use the CelebA  
 277 dataset [37], containing images with many binary attributes per image. (“male”, “blond hair”, etc).

278 We train a ResNet-50 to classify one of the hard attributes (accuracy 80%) and confirm that the  
 279 Feature Calibration holds for all the other attributes (Figure 3) that are themselves distinguishable.

280 **Quantitative predictions:** We now test the quantitative predictions made by Conjecture 1. This  
 281 conjecture states that the TV-distance between the joint distributions  $(L(x), f(x))$  and  $(L(x), y)$   
 282 is at most  $\varepsilon$ , where  $\varepsilon$  is the error of the training procedure in learning  $L$  (see Definition 1). To  
 283 test this, we consider binary task similar to Experiment 1 where (Ship, Plane) are labeled as  
 284 class 0 and (Cat, Dog) are labeled as class 1, with  $p = 0.3$  fraction of cats mislabeled to class 0.  
 285 Then, we train a convolutional network to interpolation on this task. To vary the error  $\varepsilon$  on these  
 286 distinguishable features systematically, we train networks with varying number of train samples.  
 287 Networks with fewer samples have larger  $\varepsilon$  since they are worse at classifying the distinguishable  
 288 features of (Ship, Plane, Cat, Dog). Then, we use the same setup to train networks on the binary  
 289 task and measure the TV-distance between  $(L(x), f(x))$  and  $(L(x), y)$  in this task. The results are  
 290 shown in Figure 3C. As predicted, the TV distance on the binary task is upper bounded by  $\varepsilon$  error on  
 291 the 4-way classification task.

## 292 5 Distributional Generalization

293 In order to relate our results to the classical theory of generalization, we now propose a formal  
 294 notion of “Distributional Generalization”, which subsumes both Feature Calibration and classical  
 295 generalization. In fact, we will also give preliminary evidence that this new notion can apply even for  
 296 non-interpolating methods, unlike Feature Calibration.

297 A trained model  $f$  obeys classical generalization (with respect to test error) if its error on the train set  
 298 is close to its error on the test distribution. We first rewrite this using our definitions below.

299 **Classical Generalization (informal):** Let  $f$  be a trained classifier. Then  $f$  generalizes if:

$$\mathbb{E}_{\substack{x \sim \text{TrainSet} \\ \hat{y} \leftarrow f(x)}} [\mathbb{1}\{\hat{y} \neq y(x)\}] \approx \mathbb{E}_{\substack{x \sim \text{TestSet} \\ \hat{y} \leftarrow f(x)}} [\mathbb{1}\{\hat{y} \neq y(x)\}] \quad (6)$$

300 Above,  $y(x)$  is the true class of  $x$  and  $\hat{y}$  is the predicted class. The LHS of Equation 6 is the train  
 301 error of  $f$ , and the RHS is the test error. Using our definitions of  $\mathcal{D}_{\text{tr}}, \mathcal{D}_{\text{te}}$  from Section 3.1 and  
 302 defining  $T_{\text{err}}(x, \hat{y}) := \mathbb{1}\{\hat{y} \neq y(x)\}$ , we can write Equation 6 equivalently:

$$\mathbb{E}_{x, \hat{y} \sim \mathcal{D}_{\text{tr}}} [T_{\text{err}}(x, \hat{y})] \approx \mathbb{E}_{x, \hat{y} \sim \mathcal{D}_{\text{te}}} [T_{\text{err}}(x, \hat{y})] \quad (7)$$

303 That is, classical generalization states that a certain function ( $T_{\text{err}}$ ) has similar expectations on both the  
 304 Train Distribution  $\mathcal{D}_{\text{tr}}$  and Test Distribution  $\mathcal{D}_{\text{te}}$ . We can now introduce Distributional Generalization,  
 305 which is a property of trained classifiers. It is parameterized by a set of bounded functions (“tests”):  
 306  $\mathcal{T} \subseteq \{T : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]\}$ .

307 **Distributional Generalization:** Let  $f$  be a trained classifier. Then  $f$  satisfies Distributional Gener-  
 308 alization with respect to tests  $\mathcal{T}$  if:

$$\forall T \in \mathcal{T} : \mathbb{E}_{x, \hat{y} \sim \mathcal{D}_{\text{tr}}} [T(x, \hat{y})] \approx \mathbb{E}_{x, \hat{y} \sim \mathcal{D}_{\text{te}}} [T(x, \hat{y})] \quad (8)$$

309 This states that the train and test distribution have similar expectations for *all* functions in the family  
 310  $\mathcal{T}$ , which we can write as:  $\mathcal{D}_{\text{tr}} \approx^{\mathcal{T}} \mathcal{D}_{\text{te}}$ . For the singleton set  $\mathcal{T} = \{T_{\text{err}}\}$ , this is equivalent to  
 311 classical generalization, but it may hold for much larger sets  $\mathcal{T}$ . This definition of Distributional  
 312 Generalization, like the definition of classical generalization, is just defining an object— not stating  
 313 when or how it is satisfied. Feature Calibration turns this into a concrete conjecture.  
 314

### 315 5.1 Feature Calibration as Distributional Generalization

316 We can write our Feature Calibration Conjecture as a special case of Distributional Generalization,  
 317 for a certain family of tests  $\mathcal{T}$ . Informally, for a given setting, the family  $\mathcal{T}$  is all tests which take  
 318 input  $(x, y)$ , but only depend on  $x$  via a *distinguishable feature* (Definition 1). For example, a test  
 319 of the form  $T(x, y) = g(L(x), y)$  where  $L$  is a distinguishable feature, and  $g$  is arbitrary. Formally,  
 320 for a given problem setting, suppose  $\mathcal{L}$  is the set of  $(\varepsilon, \mathcal{A}, \mathcal{D}, n)$ -distinguishable features. Then  
 321 Conjecture 1 states that  $\forall L \in \mathcal{L} : (L(x), f(x)) \approx_{\varepsilon} (L(x), y)$ . This is equivalent to the statement

$$\mathcal{D}_{\text{te}} \approx_{\varepsilon}^{\mathcal{T}} \mathcal{D} \quad (9)$$



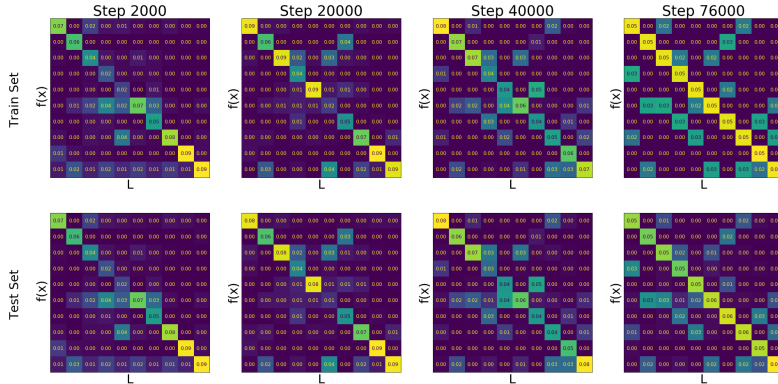


Figure 4: **Distributional Generalization for WideResNet on CIFAR-10.** The confusion matrices on the train set (top row) and test set (bottom row) remain close throughout training.

322 where  $\mathcal{T}$  is the set of functions  $\mathcal{T} := \{T : T(x, y) = g(L(x), y), L \in \mathcal{L}, g : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]\}$ .  
 323 For interpolating classifiers, we have  $\mathcal{D} \equiv \mathcal{D}_{\text{tr}}$ , and so Equation (9) is equivalent to  $\mathcal{D}_{\text{te}} \approx_{\varepsilon}^T \mathcal{D}_{\text{tr}}$ ,  
 324 which is a statement of Distributional Generalization. Since any classifier family will contain a large  
 325 number of distinguishable features, the set  $\mathcal{L}$  may be very large. Hence, the distributions  $\mathcal{D}_{\text{tr}}$  and  $\mathcal{D}_{\text{te}}$   
 326 can be thought of as being close *as distributions*.

## 327 5.2 Beyond Interpolating Methods

328 The previous sections have focused on *interpolating* classifiers, which fit their train sets exactly. Here  
 329 we informally discuss how to extend our results beyond interpolating methods. The discussion in this  
 330 section is not as precise as in previous sections, and is only meant to suggest that our abstraction of  
 331 Distributional Generalization can be useful in other settings.

332 For non-interpolating classifiers, we may still expect that they behave similarly on their test and  
 333 train sets – that is,  $\mathcal{D}_{\text{te}} \approx^T \mathcal{D}_{\text{tr}}$  for some family of tests  $\mathcal{T}$ . For example, the following is a possible  
 334 generalization of Feature Calibration to non-interpolating methods.

335 **Conjecture 2** (Generalized Feature Calibration, informal). *For trained classifiers  $f$ , the following*  
 336 *distributions are statistically close for many partitions  $L$  of the domain:*

$$\begin{aligned} (L(x), \hat{y})_{x, \hat{y} \sim \mathcal{D}_{\text{te}}} &\approx (L(x), \hat{y})_{x, \hat{y} \sim \mathcal{D}_{\text{tr}}} \end{aligned} \quad (10)$$

337 We leave unspecified the exact set of partitions  $L$  for which this holds, since we do not yet understand  
 338 the appropriate notion of “distinguishable feature” in this setting. However, we give experimental  
 339 evidence suggesting some refinement of Conjecture 2 is true. In Figure 4, we apply label noise from  
 340 a random sparse confusion to the CIFAR-10 train set. We then train a single WideResNet28-10, and  
 341 measure its predictions on the train and test sets over increasing train time (SGD steps). The top row  
 342 shows the confusion matrix of predictions  $f(x)$  vs true labels  $L(x)$  on the train set, and the bottom  
 343 row shows the corresponding confusion matrix on the test set. As the network is trained for longer, it  
 344 fits more of the noise on the train set, and this noise is mirrored almost identically on the test set. Full  
 345 experimental details, and an analogous experiment for kernels, are given in Appendix B.

## 346 6 Conclusion

347 This work initiates the study of a new kind of generalization— Distributional Generalization— which  
 348 considers the entire input-output behavior of classifiers, instead of just their test error. We presented  
 349 both new empirical behaviors, and new formal conjectures which characterize these behaviors.  
 350 Roughly, our conjecture states that the outputs of classifiers on the test set are “close in distribution”  
 351 to their outputs on the train set. These results build a deeper understanding of models used in practice,  
 352 and we hope our results inspire further work on distributional generalization in machine learning.

## References

- 353
- 354 [1] Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in  
355 neural networks. *arXiv preprint arXiv:1710.03667*, 2017.
- 356 [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overpa-  
357 rameterized neural networks, going beyond two layers. *arXiv preprint arXiv:1811.04918*,  
358 2018.
- 359 [3] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparam-  
360 eterized neural networks, going beyond two layers. In *Advances in neural information processing*  
361 *systems*, pages 6158–6169, 2019.
- 362 [4] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of op-  
363 timization and generalization for overparameterized two-layer neural networks. In *International*  
364 *Conference on Machine Learning*, pages 322–332, 2019.
- 365 [5] Susan Athey, Julie Tibshirani, Stefan Wager, et al. Generalized random forests. *The Annals of*  
366 *Statistics*, 47(2):1148–1178, 2019.
- 367 [6] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *The Journal*  
368 *of Machine Learning Research*, 18(1):629–681, 2017.
- 369 [7] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in  
370 linear regression. *Proceedings of the National Academy of Sciences*, 2020.
- 371 [8] Mikhail Belkin, Daniel J Hsu, and Partha Mitra. Overfitting or perfect fitting? risk bounds for  
372 classification and regression rules that interpolate. In *Advances in neural information processing*  
373 *systems*, pages 2300–2311, 2018.
- 374 [9] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to  
375 understand kernel learning. *arXiv preprint arXiv:1802.01396*, 2018.
- 376 [10] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-  
377 learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy*  
378 *of Sciences*, 116(32):15849–15854, 2019.
- 379 [11] Leo Breiman. Reflections after refereeing papers for nips. *The Mathematics of Generalization*,  
380 pages 11–15, 1995.
- 381 [12] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 382 [13] Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A Olshen. *Classification and*  
383 *regression trees*. CRC press, 1984.
- 384 [14] Niladri S Chatterji and Philip M Long. Finite-sample analysis of interpolating linear classifiers  
385 in the overparameterized regime. *arXiv preprint arXiv:2004.12019*, 2020.
- 386 [15] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural  
387 networks trained with the logistic loss. *arXiv preprint arXiv:2002.04486*, 2020.
- 388 [16] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL [http://archive](http://archive.ics.uci.edu/ml)  
389 [ics.uci.edu/ml](http://archive.ics.uci.edu/ml).
- 390 [17] Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds  
391 for deep (stochastic) neural networks with many more parameters than training data. *arXiv*  
392 *preprint arXiv:1703.11008*, 2017.
- 393 [18] Manuel Fernández-Delgado, Eva Cernadas, Senén Barro, and Dinani Amorim. Do we need  
394 hundreds of classifiers to solve real world classification problems? *The journal of machine*  
395 *learning research*, 15(1):3133–3181, 2014.
- 396 [19] Stanislav Fort, Paweł Krzysztof Nowak, Stanislaw Jastrzebski, and Sridhar Narayanan. Stiffness:  
397 A new perspective on generalization in neural networks. *arXiv preprint arXiv:1901.09491*,  
398 2019.

- 399 [20] Mario Geiger, Stefano Spigler, Stéphane d’Ascoli, Levent Sagun, Marco Baity-Jesi, Giulio  
400 Biroli, and Matthieu Wyart. Jamming transition as a paradigm to understand the loss landscape  
401 of deep neural networks. *Physical Review E*, 100(1):012115, 2019.
- 402 [21] Federica Gerace, Bruno Loureiro, Florent Krzakala, Marc Mézard, and Lenka Zdeborová.  
403 Generalisation error in learning with random features and the hidden manifold model. *arXiv*  
404 *preprint arXiv:2002.09339*, 2020.
- 405 [22] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized  
406 two-layers neural networks in high dimension. *arXiv preprint arXiv:1904.12191*, 2019.
- 407 [23] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation.  
408 *Journal of the American statistical Association*, 102(477):359–378, 2007.
- 409 [24] Sebastian Goldt, Madhu S Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborova.  
410 Generalisation dynamics of online learning in over-parameterised neural networks. *arXiv*  
411 *preprint arXiv:1901.09085*, 2019.
- 412 [25] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias  
413 in terms of optimization geometry. In *International Conference on Machine Learning*, pages  
414 1832–1841. PMLR, 2018.
- 415 [26] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural  
416 networks. *arXiv preprint arXiv:1706.04599*, 2017.
- 417 [27] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning:  
418 data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- 419 [28] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-  
420 dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
- 421 [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
422 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
423 pages 770–778, 2016.
- 424 [30] Úrsula Hébert-Johnson, Michael Kim, Omer Reingold, and Guy Rothblum. Multicalibration:  
425 Calibration for the (computationally-identifiable) masses. In *International Conference on  
426 Machine Learning*, pages 1939–1948, 2018.
- 427 [31] Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on  
428 document analysis and recognition*, volume 1, pages 278–282. IEEE, 1995.
- 429 [32] Rashidedin Jahandideh, Alireza Tavakoli Targhi, and Maryam Tahmasbi. Physical attribute  
430 prediction using deep residual neural networks. *arXiv preprint arXiv:1812.07857*, 2018.
- 431 [33] Alex Krizhevsky et al. Learning multiple layers of features from tiny images. 2009.
- 432 [34] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning  
433 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 434 [35] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial  
435 large learning rate in training neural networks. *arXiv preprint arXiv:1907.04595*, 2019.
- 436 [36] Tengyuan Liang and Alexander Rakhlin. Just interpolate: Kernel" ridgeless" regression can  
437 generalize. *arXiv preprint arXiv:1808.00387*, 2018.
- 438 [37] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the  
439 wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 440 [38] Song Mei and Andrea Montanari. The generalization error of random features regression:  
441 Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- 442 [39] Nicolai Meinshausen. Quantile regression forests. *Journal of Machine Learning Research*, 7  
443 (Jun):983–999, 2006.

- 444 [40] Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless  
445 interpolation of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*,  
446 2020.
- 447 [41] Elizbar A Nadaraya. On estimating regression. *Theory of Probability & Its Applications*, 9(1):  
448 141–142, 1964.
- 449 [42] Vaishnavh Nagarajan and J. Zico Kolter. Uniform convergence may be unable to explain  
450 generalization in deep learning, 2019.
- 451 [43] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever.  
452 Deep double descent: Where bigger models and more data hurt. In *International Conference on*  
453 *Learning Representations*, 2020.
- 454 [44] Hariharan Narayanan and Sanjoy Mitter. Sample complexity of testing the manifold hypothesis.  
455 In *Advances in neural information processing systems*, pages 1786–1794, 2010.
- 456 [45] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning  
457 with noisy labels. In *Advances in neural information processing systems*, pages 1196–1204,  
458 2013.
- 459 [46] Brady Neal, Sarthak Mittal, Aristide Baratin, Vinayak Tantia, Matthew Scicluna, Simon Lacoste-  
460 Julien, and Ioannis Mitliagkas. A modern take on the bias-variance tradeoff in neural networks.  
461 *arXiv preprint arXiv:1810.08591*, 2018.
- 462 [47] Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro.  
463 Towards understanding the role of over-parametrization in generalization of neural networks.  
464 *arXiv preprint arXiv:1805.12076*, 2018.
- 465 [48] Alexandru Niculescu-Mizil and Rich Caruana. Predicting good probabilities with supervised  
466 learning. In *Proceedings of the 22nd international conference on Machine learning*, pages  
467 625–632, 2005.
- 468 [49] Matthew A Olson and Abraham J Wyner. Making sense of random forest probabilities: a kernel  
469 perspective. *arXiv preprint arXiv:1812.05792*, 2018.
- 470 [50] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,  
471 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in  
472 pytorch. 2017.
- 473 [51] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,  
474 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,  
475 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*  
476 *Learning Research*, 12:2825–2830, 2011.
- 477 [52] Taylor Pospisil and Ann B Lee. Rfcde: Random forests for conditional density estimation.  
478 *arXiv preprint arXiv:1804.05753*, 2018.
- 479 [53] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances*  
480 *in neural information processing systems*, pages 1177–1184, 2008.
- 481 [54] David Rolnick, Andreas Veit, Serge Belongie, and Nir Shavit. Deep learning is robust to massive  
482 label noise. *arXiv preprint arXiv:1705.10694*, 2017.
- 483 [55] Jonas Rothfuss, Fabio Ferreira, Simon Walther, and Maxim Ulrich. Conditional density estima-  
484 tion with neural networks: Best practices and benchmarks. *arXiv preprint arXiv:1903.00954*,  
485 2019.
- 486 [56] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
487 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.  
488 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*  
489 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 490 [57] Robert E Schapire. Theoretical views of boosting. In *European conference on computational*  
491 *learning theory*, pages 1–10. Springer, 1999.

- 492 [58] Robert E Schapire, Yoav Freund, Peter Bartlett, Wee Sun Lee, et al. Boosting the margin: A new  
493 explanation for the effectiveness of voting methods. *The annals of statistics*, 26(5):1651–1686,  
494 1998.
- 495 [59] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to*  
496 *algorithms*. Cambridge university press, 2014.
- 497 [60] Vaishaal Shankar, Alex Fang, Wenshuo Guo, Sara Fridovich-Keil, Ludwig Schmidt, Jonathan  
498 Ragan-Kelley, and Benjamin Recht. Neural kernels without tangents. *arXiv preprint*  
499 *arXiv:2003.02237*, 2020.
- 500 [61] Utkarsh Sharma and Jared Kaplan. A neural scaling law from the dimension of the data manifold.  
501 *arXiv preprint arXiv:2004.10802*, 2020.
- 502 [62] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The  
503 implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*,  
504 19(1):2822–2878, 2018.
- 505 [63] Sunil Thulasidasan, Tanmoy Bhattacharya, Jeff Bilmes, Gopinath Chennupati, and Jamal  
506 Mohd-Yusof. Combating label noise in deep learning using abstention. *arXiv preprint*  
507 *arXiv:1905.10964*, 2019.
- 508 [64] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David  
509 Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J.  
510 van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew  
511 R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, CJ Carey, İlhan Polat, Yu Feng, Eric W.  
512 Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen,  
513 E. A. Quintero, Charles R Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa,  
514 Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for  
515 Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: <https://doi.org/10.1038/s41592-019-0686-2>.  
516
- 517 [65] Larry Wasserman. *All of statistics: a concise course in statistical inference*. Springer Science &  
518 Business Media, 2013.
- 519 [66] Geoffrey S Watson. Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics,*  
520 *Series A*, pages 359–372, 1964.
- 521 [67] Abraham J Wyner, Matthew Olson, Justin Bleich, and David Mease. Explaining the success  
522 of adaboost and random forests as interpolating classifiers. *The Journal of Machine Learning*  
523 *Research*, 18(1):1558–1590, 2017.
- 524 [68] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for  
525 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 526 [69] Mohammad Yaghini, Bogdan Kulynych, and Carmela Troncoso. Disparate vulnerability: On  
527 the unfairness of privacy attacks against machine learning. *arXiv preprint arXiv:1906.00389*,  
528 2019.
- 529 [70] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint*  
530 *arXiv:1605.07146*, 2016.
- 531 [71] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding  
532 deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- 533 [72] Liu Ziyin, Blair Chen, Ru Wang, Paul Pu Liang, Ruslan Salakhutdinov, Louis-Philippe Morency,  
534 and Masahito Ueda. Learning not to learn in the presence of noisy labels. *arXiv preprint*  
535 *arXiv:2002.06541*, 2020.

536 **Checklist**

- 537 1. For all authors...
- 538 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
539 contributions and scope? [Yes]
- 540 (b) Did you describe the limitations of your work? [Yes]
- 541 (c) Did you discuss any potential negative societal impacts of your work? [No] This paper  
542 does not introduce any new methods or applications, and we thus cannot predict any  
543 near-term societal impact (positive or negative).
- 544 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
545 them? [Yes]
- 546 2. If you are including theoretical results...
- 547 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 548 (b) Did you include complete proofs of all theoretical results? [Yes]
- 549 3. If you ran experiments...
- 550 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
551 perimental results (either in the supplemental material or as a URL)? [No] No new  
552 methods were introduced, so the code is standard. We fully specify all experimental  
553 hyperparameters for the sake of reproduction.
- 554 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
555 were chosen)? [Yes]
- 556 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
557 ments multiple times)? [No] The experiments we consider all exhibit concentration  
558 around their expected values, and this is well-known in the community. Notably, we  
559 only consider supervised learning.
- 560 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
561 of GPUs, internal cluster, or cloud provider)? [No]
- 562 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 563 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 564 (b) Did you mention the license of the assets? [No]
- 565 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 566 (d) Did you discuss whether and how consent was obtained from people whose data you're  
567 using/curating? [No]
- 568 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
569 information or offensive content? [No]
- 570 5. If you used crowdsourcing or conducted research with human subjects...
- 571 (a) Did you include the full text of instructions given to participants and screenshots, if  
572 applicable? [N/A]
- 573 (b) Did you describe any potential participant risks, with links to Institutional Review  
574 Board (IRB) approvals, if applicable? [N/A]
- 575 (c) Did you include the estimated hourly wage paid to participants and the total amount  
576 spent on participant compensation? [N/A]