# Datamodels: Predicting Predictions from Training Data

Andrew Ilyas*
ailyas@mit.edu
MIT

Sung Min Park*
sp765@mit.edu
MIT

Logan Engstrom*
engstrom@mit.edu
MIT

Guillaume Leclerc
leclerc@mit.edu
MIT

Aleksander Mądry
madry@mit.edu
MIT

## Abstract

We present a conceptual framework, *datamodeling*, for analyzing the behavior of a model class in terms of the training data. For any fixed "target" example $x$, training set $S$, and learning algorithm, a *datamodel* is a parameterized function $2^S \to \mathbb{R}$ that for any subset of $S' \subset S$—using only information about which examples of $S$ are contained in $S'$—predicts the outcome of training a model on $S'$ and evaluating on $x$. Despite the potential complexity of the underlying process being approximated (e.g., end-to-end training and evaluation of deep neural networks), we show that even simple *linear* datamodels can successfully predict model outputs. We then demonstrate that datamodels give rise to a variety of applications, such as: accurately predicting the effect of dataset counterfactuals; identifying brittle predictions; finding semantically similar examples; quantifying train-test leakage; and embedding data into a well-behaved and feature-rich *representation space*.[1]

## 1  Introduction

*What kinds of biases does my (machine learning) system exhibit? What correlations does it exploit? On what subpopulations does it perform well (or poorly)?*

A recent body of work in machine learning suggests that the answers to these questions lie within both the learning algorithm and the training data used [GDG17; CLK+19; IST+19; Hoo21; JTM21]. However, it is often difficult to understand *how* algorithms and data combine to yield model predictions. In this work, we present *datamodeling*—a framework for tackling this question by forming an explicit model for predictions in terms of the training data.

**Setting.**   Consider a typical machine learning setup, starting with a training set $S$ comprising $d$ input-label pairs. The focal point of this setup is a *learning algorithm* $\mathcal{A}$ that takes in such a training set of input-label pairs, and outputs a trained model. (Note that this learning algorithm does not have to be deterministic—for example, $\mathcal{A}$ might encode the process of training a deep neural network from random initialization using stochastic gradient descent.)

Now, consider a *fixed* input $x$ (e.g., a photo from the test set of a computer vision benchmark) and define

$$f_{\mathcal{A}}(x; S) := \text{the outcome of training a model on } S \text{ using } \mathcal{A} \text{ and evaluating it on the input } x, \qquad (1)$$

where we leave "outcome" intentionally broad to capture a variety of use cases. For example, $f_{\mathcal{A}}(x; S)$ may be the cross-entropy loss of a classifier on $x$, or the squared-error of a regression model on $x$. The potentially stochastic nature of $\mathcal{A}$ means that $f_{\mathcal{A}}(x; S)$ is a random variable.

---

*Equal contribution.

[1]Data for this paper (including pre-computed datamodels as well as raw predictions from four million trained deep networks) is available at https://github.com/MadryLab/datamodels-data.

**Goal.** Broadly, we aim to understand how the training examples in $S$ combine through the learning algorithm $\mathcal{A}$ to yield $f_{\mathcal{A}}(x; S)$ (again, for the *specifically* chosen input $x$). To this end, we leverage a classic technique for studying black-box functions: *surrogate modeling* [SWM+89]. In surrogate modeling, one replaces a complex black-box function with an inexact but significantly easier-to-analyze approximation, then uses the latter to shed light on the behavior of the original function.

In our context, the complex black-box function is $f_{\mathcal{A}}(x; \cdot)$. We thus aim to find a simple *surrogate* function $g(S')$ whose output roughly matches $f_{\mathcal{A}}(x; S')$ for a variety of training sets $S'$ (but again, for a *fixed* input $x$). Achieving this goal would reduce the challenge of scrutinizing $f_{\mathcal{A}}(x; \cdot)$—and more generally, the map from training data to predictions through learning algorithm $\mathcal{A}$—to the (hopefully easier) task of analyzing $g$.

**Datamodeling.** By parameterizing the surrogate function $g$ (e.g., as $g_\theta$, for a parameter vector $\theta$), we transform the challenge of constructing a surrogate function into a *supervised learning* problem. In this problem, the "training examples" are subsets $S' \subset S$ of the original task's training set $S$, and the corresponding "labels" are given by $f_{\mathcal{A}}(x; S')$ (which we can compute by simply training a new model on $S'$ with algorithm $\mathcal{A}$, and evaluating on $x$). Our goal is then to fit a parametric function $g_\theta$ mapping the former to the latter.

We now formalize this idea as *datamodeling*—a framework that forms the basis of our work. In this framework, we first fix a distribution over subsets that we will use to collect "training data" for $g_\theta$,

$$\mathcal{D}_S := \text{a fixed distribution over subsets of } S \text{ (i.e., support}(\mathcal{D}_S) \subseteq 2^S), \tag{2}$$

and then use $\mathcal{D}_S$ to collect a *datamodel training set*, or a collection of pairs

$$\{(S_1, f_{\mathcal{A}}(x; S_1)), \ldots, (S_m, f_{\mathcal{A}}(x; S_m))\},$$

where $S_i \sim \mathcal{D}_S$, and again $f_{\mathcal{A}}(x; S_i)$ is the result of training a model on $S_i$ and evaluating on $x$ (cf. (1)).

We next focus on how to parameterize our surrogate function $g_\theta$. In theory, $g_\theta$ can be any map that takes as input subsets of the training set, and returns estimates of $f_{\mathcal{A}}(x; \cdot)$. However, to simplify $g_\theta$ we ignore the actual *contents* of the subsets $S_i$, and instead focus solely on the *presence* of each training example of $S$ within $S_i$. In particular, we consider the *characteristic vector* corresponding to each $S_i$,

$$\mathbf{1}_{S_i} \in \{0, 1\}^d \qquad \text{such that} \qquad (\mathbf{1}_{S_i})_j = \begin{cases} 1 & \text{if } z_j \in S_i \\ 0 & \text{otherwise}, \end{cases} \tag{3}$$

a vector that indicates which elements of the original training set $S$ belong to a given subset $S_i$. We then define a <u>datamodel</u> for a given input $x$ as a function

$$g_\theta : \{0, 1\}^d \to \mathbb{R}, \qquad \text{where} \qquad \theta = \arg\min_w \frac{1}{m} \sum_{i=1}^m \mathcal{L}\left(g_w(\mathbf{1}_{S_i}), f_{\mathcal{A}}(x; S_i)\right), \tag{4}$$

and $\mathcal{L}(\cdot, \cdot)$ is a fixed loss function (e.g., squared-error). This setup (4) places datamodels squarely within the realm of supervised learning: e.g., we can easily test a given datamodel by sampling new subset-output pairs $\{(S_i, f_{\mathcal{A}}(x; S_i))\}$ and computing average loss. For completeness, we restate the entire datamodeling framework below:

---

**Definition 1** (Datamodeling). *Consider a fixed training set $S$, a learning algorithm $\mathcal{A}$, a target example $x$, and a distribution $\mathcal{D}_S$ over subsets of $S$. For any set $S' \subset S$, let $f_{\mathcal{A}}(x; S')$ be the (stochastic) output of training a model on $S'$ using $\mathcal{A}$, and evaluating on $x$. A <u>datamodel</u> for $x$ is a parametric function $g_\theta$ optimized to predict $f_{\mathcal{A}}(x; S_i)$ from training subsets $S_i \sim \mathcal{D}_S$, i.e.,*

$$g_\theta : \{0, 1\}^{|S|} \to \mathbb{R}, \qquad \text{where} \qquad \theta = \arg\min_w \widehat{\mathbb{E}}^{(m)}_{S_i \sim \mathcal{D}_S}\left[\mathcal{L}\left(g_w(\mathbf{1}_{S_i}), f_{\mathcal{A}}(x; S_i)\right)\right],$$

*$\mathbf{1}_{S_i} \in \{0, 1\}^{|S|}$ is the characteristic vector of $S_i$ in $S$ (see (3)), $\mathcal{L}(\cdot, \cdot)$ is a loss function, and $\widehat{\mathbb{E}}^{(m)}$ is an $m$-sample empirical estimate of the expectation.*

---

The pseudocode for computing datamodels is in Appendix A. Before proceeding further, we highlight two critical (yet somewhat subtle) properties of the datamodeling framework:

- **Datamodeling studies model *classes*, not specific models**: Datamodeling focuses on the entire distribution of models induced by the algorithm $\mathcal{A}$, rather than a specific model. Recent work suggests this distinction is particularly significant for modern learning algorithms (e.g., neural networks), where models can exhibit drastically different behavior depending on only the choice of random seed during training [DHM+20; NB20; JNB+21; ZGK+21]—we discuss this further in Section 6.

- **Datamodels are target example-specific**: A datamodel $g_\theta$ predicts model outputs on a specific but arbitrary target example $x$. This $x$ might be an example from the test set, a synthetically generated example, or even (as we will see in Section 3.1) an example from the training set $S$ itself. We will often work with *collections* of datamodels corresponding to a set of target examples (e.g., we might consider a test set $\{x_1, \ldots x_n\}$ with corresponding datamodels $\{g_{\theta_1}, \ldots g_{\theta_n}\}$). In Section 3 we show that as long as the learning algorithm $\mathcal{A}$ and the training set $S$ are fixed, computing a collection of datamodels simultaneously is not much harder than computing a single one.

## 1.1 Roadmap and contributions

The key contribution of our work is the *datamodeling framework* described above, which allows us to analyze the behavior of a machine learning algorithm $\mathcal{A}$ in terms of the training data. In the remainder of this work, we show how to instantiate, implement, and apply this framework.

We begin in Section 2 by considering a concrete instantiation of datamodeling in which the map $g_\theta$ is a *linear* function. Then, in Section 3 we develop the remaining machinery required to apply this instantiation to deep neural networks trained on standard image datasets. In the rest of the paper, we find that:

- **Datamodels successfully predict model outputs (§ 3.2, Figure 1)**: despite their simplicity, datamodels yield predictions that match expected model outputs on new sets $S$ drawn from the same distribution $\mathcal{D}_S$. (For example, the Pearson correlation between predicted and ground-truth outputs is $r > 0.99$.)

- **Datamodels successfully predict counterfactuals (§ 4.1, Figure 2)**: predictions correlate with model outputs even on out-of-distribution training subsets (Figures 8, J.3 and Appendix F.1) allowing us to estimate the *causal effect* of removing training images on a given test prediction. Leveraging this ability, we find that *for 50% of CIFAR-10 [Kri09] test images, models can be made incorrect by removing less than 200 target-specific training points (i.e., 0.4% of the total training set size)*. If one mislabels the training examples instead of only removing them, *35 label-specific points* suffice.
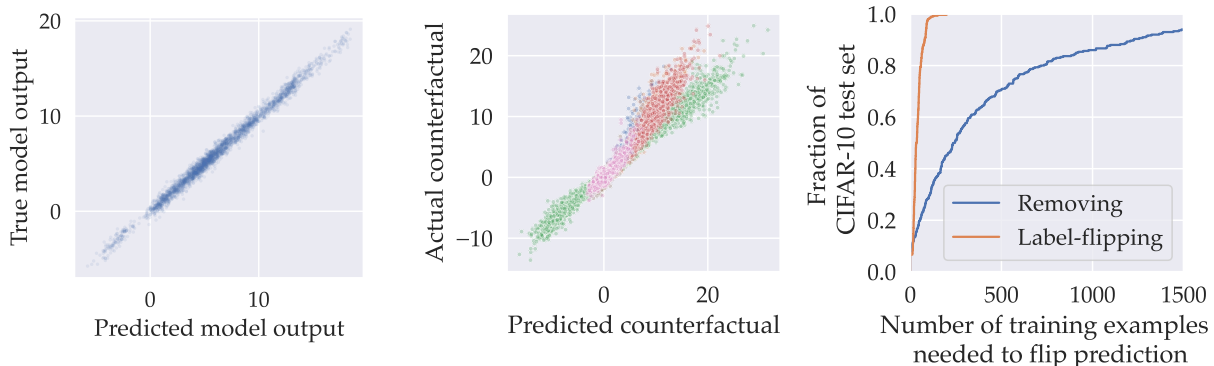


**Figure 1:** Datamodels predict $(g_\theta(S'),$ x-axis) the outcome of training models on subsets $S'$ of the training set $S$ sampled from $\mathcal{D}_S$ and evaluating on $x$ $(\mathbb{E}[f_\mathcal{A}(x; S')],$ y-axis)

**Figure 2:** Datamodels predict out-of-distribution dataset counterfactuals (left) and identify brittle predictions (right). As seen on the right, approximately 50% of predictions on the CIFAR-10 test set can be flipped by removing less than 200 (target-specific) training images. If we flip the labels of chosen training images instead of removing them, just 35 images suffices.

- **Datamodel weights encode similarity (§ 4.2, Figure 3)**: the most positive (resp., negative) datamodel weights tend to correspond to similar training images from the same (resp., different) class as the target example $x$. We use this property to identify significant train-test leakage across both datasets we study (CIFAR-10 and Functional Map of the World [CFW+18; KSM+20]).
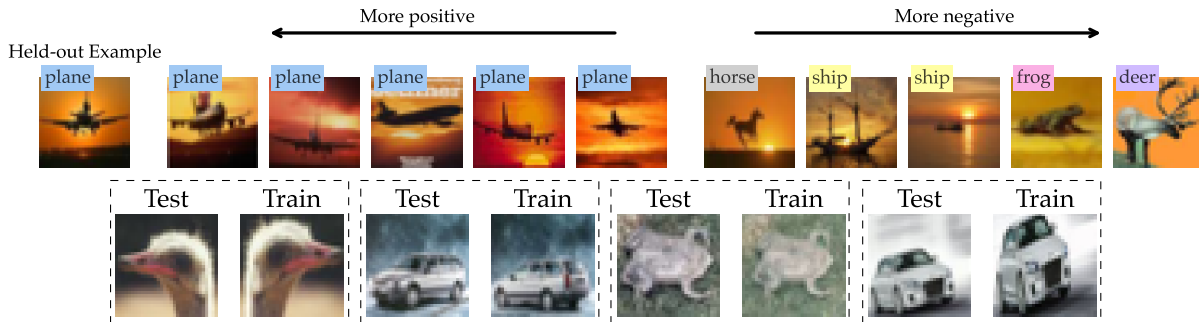


**Figure 3:** High-magnitude datamodel weights identify semantically similar training examples (top), which we can use to find train-test leakage in benchmark computer vision datasets (bottom).

- **Datamodels yield a well-behaved embedding (§ 4.3, Figure 4)**: viewing datamodel weights as a *feature embedding* of each image into $\mathbb{R}^d$ (where $d$ is the training set size), we discover a well-behaved representation space. In particular, we find that so-called *datamodel embeddings*:

  (a) enable (qualitatively) high-quality clustering;

  (b) allow us to identify model-relevant *subpopulations* that we can causally verify in a natural sense;

  (c) have a number of advantages over representations derived from, e.g., the penultimate layer of a fixed pre-trained network, such as higher effective dimensionality and (*a priori*) human-meaningful coordinates.
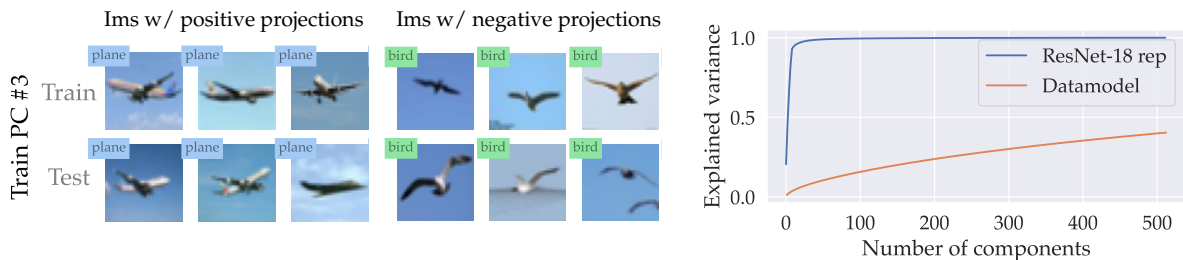


**Figure 4:** Datamodels yield a well-behaved (left), relatively dense (right) embedding of any given input into $\mathbb{R}^d$, where $d$ is the training set size. Applying natural manipulations to these embeddings enables a variety of applications which we discuss more thoroughly in § 4.3.

More broadly, datamodels turn out to be a versatile tool for understanding how learning algorithms leverage their training data. In Section 6, we contextualize datamodeling with respect to several ongoing lines of work in machine learning and statistics. We conclude, in Section 7, by outlining a variety of directions for future work on both improving and applying datamodels.

## 2 Constructing (linear) datamodels

As described in Section 1, building datamodels comprises the following steps:

  (a) pick a parameterized class of functions $g_\theta$;

(b) sample a collection of subsets $S_i \subset S$ from a fixed training set according to a distribution $\mathcal{D}_S$;

(c) for each subset $S_i$, train a model using algorithm $\mathcal{A}$, evaluate the model on target input $x$ using the relevant metric (e.g., loss); collect the resulting pair $(\mathbf{1}_{S_i}, f_\mathcal{A}(x; S_i))$;

(d) split the collected dataset of subset-output pairs into a datamodel training set of size $m$, a datamodel validation set of size $m_{val}$, and a datamodel test set of size $m_{test}$;

(e) estimate parameters $\theta$ by fitting $g_\theta$ on subset-output pairs, i.e., by minimizing

$$\frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\left(g_\theta(\mathbf{1}_{S_i}),\ f_\mathcal{A}(x; S_i)\right)$$

over the collected datamodel training set, and use the validation set to perform model selection.

We now explicitly instantiate this framework, with the goal of understanding the predictions of (deep) *classification* models. To this end, we revisit steps (a)-(e) above, and consider each relevant aspect—the sampling distribution $\mathcal{D}_S$, the output function $f_\mathcal{A}(x; S)$, the parameterized family $g_\theta$, and the loss function $\mathcal{L}(\cdot, \cdot)$—separately:

**(a) What surrogate function $g_\theta$ should we use?**   The first design choice to make is which family of parameterized surrogate functions $g_\theta$ to optimize over. At first, one might be inclined to use a complex family of functions in the hope of reducing potential misspecification error. After all, $g_\theta$ is meant to be a surrogate for the end-to-end training of a deep classifier. In this work, however, we will instantiate datamodeling by taking $g_\theta(\cdot)$ to be a simple *linear* mapping

$$g_\theta(\mathbf{1}_{S_i}) := \theta^\top \mathbf{1}_{S_i} + \theta_0, \tag{5}$$

where we recall that $\mathbf{1}_{S_i}$ is the size-*d characteristic vector* of $S_i$ within $S$ (cf. (3)).

**Remark 1.** *While we will allow $g_\theta(\cdot)$ to fit a <u>bias</u> term as above, for notational convenience we omit $\theta_0$ throughout this work and will simply write $\theta^\top \mathbf{1}_{S_i}$ to represent a datamodel prediction for the set $S_i$.*

**(b) What distribution $\mathcal{D}_S$ over training subsets do we use?**   In step (a) of the estimation process above, we collect a "datamodel training set" by sampling subsets $S_i \subset S$ from a distribution $\mathcal{D}_S$. A simple first choice for $\mathcal{D}_S$—and indeed, the one we consider for the remainder of this work—is the distribution of random $\alpha$-fraction subsets of the training set. Formally, we set

$$\mathcal{D}_S = \mathrm{Uniform}\left(\{S' \subset S : |S'| = \alpha d\}.\right) \tag{6}$$

This design choice reduces the choice of $\mathcal{D}_S$ to a choice of *subsampling fraction* $\alpha \in (0, 1)$, a decision whose impact we explore in Section 5. In practice, we estimate datamodels for *several* choices of $\alpha$, as it turns out that the value of $\alpha$ corresponding to the most useful datamodels can vary by setting.

**(c) What outputs $f_\mathcal{A}(x; S')$ should we track?**   Recall that for any subset $S' \subset S$ of the training set $S$, $f_\mathcal{A}(x; S')$ is intended to be a specific (potentially stochastic) function representing the output of a model trained on $S'$ and evaluated on a target example $x$. There are, however, several candidates for $f_\mathcal{A}(x; S')$ based on which model output we opt to track.

In the context of understanding *classifiers*, perhaps the simplest such candidate is the *correctness* function (i.e., a stochastic function that is 1 if the model trained on $S'$ is correct on $x$, and 0 otherwise). However, while the correctness function may be a natural choice for $f_\mathcal{A}(x; S')$, it turns out to be suboptimal in two ways. First, fitting to the correctness function ignores potentially valuable information about the model's confidence in a given decision. Second, recall that our procedure fits model outputs using a least-squares linear model, which is not designed to properly handle discrete (binary) dependent variables.

A natural way to improve over our initial candidate would thus be to use continuous output function, such as cross-entropy loss or correct-label confidence. But which exact function should we choose? In Appendix C, we describe a heuristic that we use to guide our choice of the *correct-class margin*:

$$f_\mathcal{A}(x; S') := (\text{logit for correct class}) - (\text{highest incorrect logit}). \tag{7}$$

**(e) What loss function $\mathcal{L}$ should we minimize?** In step (d) above, we are free to pick any estimation algorithm for $\theta$. This freedom of choice allows us to incorporate *priors* into the datamodeling process. In particular, one might expect that predictions on a given target example will not depend on every training example. We can incorporate a corresponding *sparsity prior* by adding $\ell_1$ regularization, i.e., setting

$$\theta = \min_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \left( w^\top \mathbf{1}_{S_i} - f_{\mathcal{A}}(x; S_i) \right)^2 + \lambda \|w\|_1, \tag{8}$$

where we recall that $d$ is the size of the original training set $S$. We can use cross-validation to select the regularization parameter $\lambda$ for each specific target example $x$.

# 3 Accurately predicting outputs with datamodels

We now demonstrate how datamodels can be applied in the context of deep neural networks—specifically, we consider deep image classifiers trained on two standard datasets: CIFAR-10 [Kri09] and Functional Map of the World (FMoW) [KSM+20] (see Appendix D.1 for more information on each dataset).

**Goal.** As discussed in Section 1, our goal is to construct a *collection* of datamodels for each dataset, with each datamodel predicting the model-training outcomes for a *specific* target example. Thus, for both CIFAR and FMoW, we fix a deep learning algorithm (architecture, random initialization, optimizer, etc.), and aim to estimate a datamodel for each test set example *and* training set example. As a result, we will obtain $n = 10,000$ "test set datamodels" and $n = 50,000$ "training set datamodels" for CIFAR (each being a linear model $g_\theta$ parameterized by a vector $\theta \in \mathbb{R}^d$, for $d = 50,000$); as well as $n = 3,138$ test set datamodels and $n = 21,404$ training set datamodels for FMoW (again, parameterized by $\theta \in \mathbb{R}^d$ where $d = 21,404$).

## 3.1 Implementation details

Before applying datamodels to our two tasks of interest, we address a few remaining technical aspects of datamodel estimation:

**Simultaneously estimating datamodels for a collection of target examples.** Rather than repeat the entire datamodel estimation process for each target example $x$ of interest separately, we can estimate datamodels for an entire *set* of target examples simultaneously through model reuse. Specifically, we train a large pool of models on subsets $S_i \subset S$ sampled from the distribution $\mathcal{D}_S$, and use the *same* models to compute outputs $f_{\mathcal{A}}(x; S_i)$ for each target example $x$.

**Collecting a (sufficiently large) datamodel training set.** The cost of obtaining a single subset-output pair can be non-trivial—in our case, it involves training a ResNet from scratch on CIFAR-10. It turns out, however, that recent advances in fast neural network training [Pag18; LIE+22] allow us to train a wealth of models on different $\alpha$-subsets of each dataset *very* efficiently. (For example, for $\alpha = 50\%$ we can use [LIE+22] to train 40,000 models/day on an $8 \times$ A100 GPU machine; see Appendix D.2 for details.) We train $m = 300,000$ CIFAR models and $m = 150,000$ FMoW models on $\alpha = 50\%$ subsets of each dataset. We also train $m$ models for each subsampling fraction $\alpha \in \{10\%, 20\%, 75\%\}$, using $\alpha$ to scale $m$. See Table 1 for a summary of the models trained.

**Computing datamodels when the target example is a training input.** Recall that the target example $x$ for which we estimate a datamodel can be arbitrary. In particular, $x$ could itself be a training example—indeed, as we mention above, our goal is to estimate a datamodel for every image in the FMoW and CIFAR-10 test *and* training sets. When $x$ is in the training set, however, we slightly alter the datamodel estimation objective (8) to exclude training sets $S_i$ containing the target example:

$$\theta = \min_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{x \notin S_i\} \cdot \left( w^\top \mathbf{1}_{S_i} - f_{\mathcal{A}}(x; S_i) \right)^2 + \lambda \|w\|_1. \tag{9}$$

| | Models trained | |
|---|---|---|
| Subset size ($\alpha$) | CIFAR-10 | FMoW |
| 0.1 | 1,500,000 | – |
| 0.2 | 750,000 | 375,000 |
| 0.5 | 300,000 | 150,000 |
| 0.75 | 600,000 | 300,000 |

Table 1: The number of models (ResNet-9 for CIFAR and ResNet-18 for FMoW) used to estimate datamodels for each dataset. All models are trained from scratch using optimized code [LIE+22] (e.g., each $\alpha = 0.5$ model on CIFAR-10 takes 17s to train (on a single A100 GPU) to 90% accuracy; see Appendix D.2 for details).

**Running LASSO regression at scale.** After training the models, we record the correct-class margin for all the (train and test) images as well as the training subsets. Our task now is to estimate, for each example in the train and test set, a datamodel $g_\theta$ mapping subsets $\mathbf{1}_{S_i}$ to observed margins. Recall that we compute datamodels via $\ell_1$-regularized least-squares regression (cf. (8)), where we set the regularization parameter $\lambda$ for each test image independently via a fixed validation set of trained models.

However, most readily available LASSO solvers require too much memory or are prohibitively slow for our values of $n$ (the number of datamodels to estimate), $m$ (the number of models trained and thus the size of the datamodel training set of subset-output pairs), and $d$ (the size of the original tasks training set and thus the input dimensionality of the regression problem in (8)). We therefore built a custom solver leveraging the works of [WSM21] and [LIE+22]—details of our implementation are in Appendix E.1.

## 3.2 Results: linear datamodels can predict deep network training

For both datasets considered (CIFAR-10 and FMoW), we minimize objectives (8) (respectively, (9)) yielding a datamodel $g_{\theta_i}$ for each example $x_i$ in the test set (respectively, training set). We now assess the quality of these datamodels in terms of how well they predict model outputs on *unseen* subsets (i.e., fresh samples from $\mathcal{D}_S$). We refer to this process as *on-distribution* evaluation because we are interested in subsets $S_i$, sampled from the same distribution $\mathcal{D}_S$ as the datamodel training set, but *not* the exact ones used for estimation. (In fact, recall that we explicitly held out $m_{test}$ subset-output pairs for evaluation in Section 2.)

We focus here on the collection of datamodels corresponding to the CIFAR-10 test set, i.e., a set of linear datamodel parameters $\{\theta_1, \ldots, \theta_n\}$ corresponding to examples $\{x_1, \ldots, x_n\}$ for $n = 10,000$ (analogous results for FMoW are in Appendix E.2). In Figure 5, aggregating over both datamodels $\{g_{\theta_j}\}_{j=1}^n$ and held-out subsets $\{S_i\}_{i=1}^m$, we compare datamodel predictions $\theta_j^\top \mathbf{1}_{S_i}$ to *expected* true model outputs $\mathbb{E}[f_\mathcal{A}(x_j; S_i)]$ (which we estimate by training 100 models on the same subset $S_i$ and averaging their output on $x_j$). The results show a near-perfect correspondence between datamodel predictions and ground truth. Thus, for a given target example $x$, we can accurately predict the outcome of "training a neural network on a random ($\alpha$-)training subset and computing correct-class margin on $x$" (a process that involves hundreds of SGD steps on a non-convex objective) as a simple *linear* function of the characteristic vector of the subset.

**Sample complexity.** We next study the dependence of datamodel estimation on the size of the datamodel training set $m$. Specifically, we can measure the *on-distribution* average mean-squared error (MSE) as

$$\text{MSE}(\{\theta_1, \ldots, \theta_n\}) = \frac{1}{2n} \sum_{j=1}^n \left( \mathbb{E}_{S_i \sim \mathcal{D}_S} \left[ \left( \theta_j^\top \mathbf{1}_{S_i} - f_\mathcal{A}(x_j; S_i) \right)^2 \right] \right). \tag{10}$$

To evaluate (10), we replace the inner expectation with an empirical average, again using a heldout set of samples that was not used for estimation.

In Figure 6, we plot average MSE as a function of the number of trained models $m$. To put the results into context, we introduce the *optimal mean-squared error loss* (OPT), which is the MSE (10) with datamodel
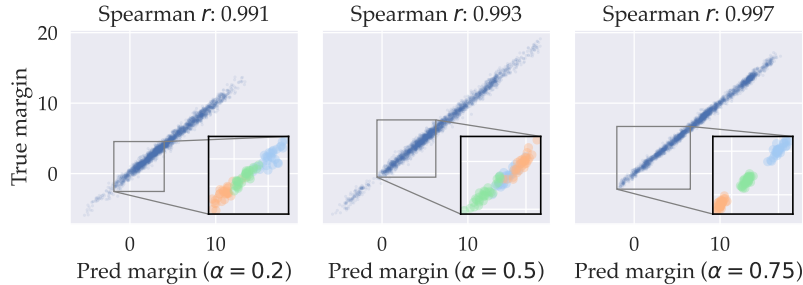
**Figure 5:** *Linear* **datamodels accurately predict margins.** Each point in the graphs above corresponds to a specific target example $x_j$ *and* a specific held-out training set $S_i$ from CIFAR-10. The $y$-coordinate represents the ground-truth margin $f_{\mathcal{A}}(x_j; S_i)$, averaged across $T{=}100$ models trained on $S_i$. The $x$-coordinate represents the *datamodel-predicted* value of the same quantity. We observe a strong linear correlation (as seen in the main blue line) that persists even at the level of individual examples (the bottom-right panel shows data for three random target examples $x_j$ color-coded by example). Corresponding plots for $\alpha = 10\%$ and FMoW are in Figures E.3 and E.4.

**Figure 6:** Average mean-squared error (Eqn. (10)) for CIFAR-10 test set datamodels ($\alpha = 0.5$) as a function of the size of the datamodel training set $m$. The red line denotes optimal error (Eqn. (11)) based on inherent noise in training.

predictors $\theta_j^{\top} \mathbf{1}_{S_i}$ replaced by the optimal deterministic predictors $\mathbb{E}[f_{\mathcal{A}}(x_j; S_i)]$:

$$\text{OPT} = \frac{1}{2n} \sum_{j=1}^{n} \left( \mathbb{E}\left[ \left( \mathbb{E}\left[ f_{\mathcal{A}}(x_j; S_i) \right] - f_{\mathcal{A}}(x_j; S_i) \right)^2 \right] \right). \tag{11}$$

Note that OPT is independent of the estimator $g_{\theta}$ and measures only the inherent variance in the prediction problem, i.e., loss that will necessarily be incurred due only to inherent noise in deep network training.

**The role of regularization.** Finally, in Appendix E.2, we study the effect of the regularization parameter $\lambda$ (cf. (8) and (9)) on datamodel performance. In particular, in Figure E.2 we plot the variation in average MSE, on both *on-distribution* subsets (i.e., the exact subsets that we used to optimize (8)) and unseen subsets, as we vary the regularization parameter $\lambda$ in (8). We find that—as predicted by classical learning theory—setting $\lambda = 0$ leads to *overfit* datamodels, i.e., estimators $g_{\theta}$ that perform well on the exact subsets that were used to estimate them, but are poor output predictors on *new* subsets $S_i$ sampled from $\mathcal{D}_S$. (In fact, using $m = 300,000$ trained models with $\lambda = 0$ results in higher MSE than using only $m = 10,000$ with optimal $\lambda$, i.e., the left-most datapoint in Figure 1).

## 4 Leveraging datamodels

Now that we have introduced (Section 1), instantiated (Section 2), and implemented (Section 3) the data-modeling framework, we turn to some of its applications. Specifically, we will now show how to apply datamodels within three different contexts:

**Counterfactual prediction.** We originally constructed datamodels to predict the outcome of training a model on *random* ($\alpha$-)subsets of the training set. However, it turns out that we can also use datamodels to predict model outputs on *arbitrary* subsets (i.e., subsets that are "off-distribution" from the perspective of the datamodel prediction task). To illustrate the utility of this capability, we will use datamodels to (a) identify predictions that are *brittle* to removal of relatively few training points, and (b) estimate *data counterfactuals*, i.e., the causal effects of removing groups of training examples.

8

**Train-test similarity.** We demonstrate that datamodels can identify, for any given target example $x$, a set of visually similar examples in the training data. Leveraging this ability, we will identify instances of *train-test leakage*, i.e., when test examples are duplicated (or nearly duplicated) within the training set.

**Data embeddings.** We show that for a given target example $x$ with corresponding datamodel $g_\theta$, we can use the parameters $\theta$ of the datamodel as an *embedding* of the target example $x$ into $\mathbb{R}^d$ (where $d$ is the training set size). By applying two standard data exploration techniques to such embeddings, we demonstrate that they capture *latent structure* within the data, and enable us to find model-relevant *data subpopulations*.

## 4.1 Counterfactual prediction

So far, we have computed and evaluated datamodels entirely within a supervised learning framework. In particular, we constructed datamodels with the goal of predicting the outcome of training on *random* subsets of the training set (sampled from a distribution $\mathcal{D}_S$ (6)) and evaluating on a fixed target example $x$. Accordingly, for each target example $x$, we evaluated its datamodel $g_\theta$ by (a) sampling new random subsets $S_i$ (from the same distribution); (b) training (a neural network) on each one of these subsets; (c) measuring correct-class margin on the target example $x$; and (d) comparing the results to the datamodel's predictions (namely, $g_\theta(S_i)$) in terms of *expected* mean-squared error (see (10)) over the distribution of subsets.

We will now go beyond this framework, and use datamodels to predict the outcome of training on *arbitrary* subsets of the training set. In particular, consider a fixed target example $x$ with corresponding datamodel $g_\theta$. For any subset $S'$ of the training set $S$, we will use the *datamodel-predicted* outcome of training on $S'$ and evaluating on $x$, i.e., $g_\theta(\mathbf{1}_{S'})$, in place of the *ground-truth* outcome $f_\mathcal{A}(x; S')$. Since $S'$ is an *arbitrary* subset of the training set, it is "out-of-distribution" with respect to the distribution of fixed-size subsets $\mathcal{D}_S$ that we designed the datamodel to operate on. As such, using datamodel predictions in place of end-to-end-model training in this manner is not a priori guaranteed to work. Nevertheless, we will demonstrate through two applications that datamodels *can* in fact be effective proxies for end-to-end model training, even for such out-of-distribution subsets.

> **Use Case 1** (Proxy for end-to-end training)**.** *We can use datamodel predictions as an efficient, closed-form proxy for end-to-end model training. That is, for a test example $x$ with datamodel $g_\theta$, and an <u>arbitrary</u> subset $S'$ of the training set $S$, we can leverage the approximation*
>
> $$f_\mathcal{A}(x; S') \approx g_\theta(\mathbf{1}_{S'})$$

### 4.1.1 Measuring brittleness of individual predictions to training data removal

We first illustrate the utility of datamodels as a proxy for model training by using them to answer the question: *how brittle are model predictions to removing training data?* While all useful learning algorithms are data-dependent, cases where model behavior is sensitive to just a few data points are often of particular interest or concern [DKM+06; BGM21]. To quantify such sensitivity, we define the *data support* $\text{SUPPORT}(x)$ of a target example $x$ as

$$\text{SUPPORT}(x) = \text{the smallest training subset } R \subset S \text{ such that classifiers trained on } S \setminus R \quad (12)$$
$$\text{misclassify } x \text{ on average.}^2$$

Intuitively, examples with a small data support are the examples for which removing a small subset of the training data significantly changes model behavior, i.e., they are "brittle" examples by our criterion of interest. By computing $\text{SUPPORT}(x)$ for every image in the test set, we can thus get an idea of how brittle model predictions are to removing training data.

---

[2]We define misclassification here as having expected margin (Eq. (7)) less than 0.

**Computing data support.** One way to compute SUPPORT($x$) for a given target example $x$ would be to train several models on every possible subset of the training set $S$, then report the largest subset for which the example was misclassified on average—the complement of this set would be *exactly* SUPPORT($x$). However, exhaustively computing data support in this manner is simply intractable.

Using datamodels as a proxy for end-to-end model training provides an (efficient) alternative approach. Specifically, rather than training models on every possible subset of the training set, we can use datamodel-predicted outputs $g_\theta(S')$ to perform a *guided search*, and only train on subsets for which predicted margin on the target example is small. This strategy (described in detail in Algorithm 1 and in Appendix F.3) allows us to compute estimates of the data support while training only a handful of models per target example.

**Results.** We apply our algorithm to estimate SUPPORT($x$) for 300 random target examples in the CIFAR-10 test set. For over 90% of these 300 examples, we are able to *certify* that our estimated data support is *strictly larger than* the true data support SUPPORT($x$) (i.e., that we are not over-estimating brittleness) by training several models after excluding the estimated data support and checking that the target example is indeed misclassified on average.

We plot the distribution of estimated data support sizes in Figure 7. Around *half* of the CIFAR-10 test images have a datamodel-estimated data support comprising 250 images or less, meaning that removing a specific 0.4% of the CIFAR-10 training set induces misclassification. Similarly, 20% of the images had an estimated data support of less than *40 training images* (which corresponds to *0.08% of the training set*).
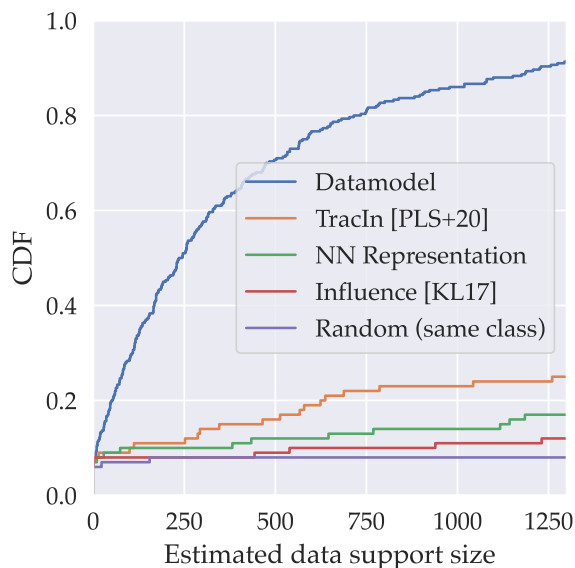


**Figure 7: Characterizing brittleness.** We use data-models to estimate *data support* (i.e., the minimal set of training examples whose removal causes misclassification) for 300 random CIFAR-10 test examples, and plot the cumulative distribution of estimated sizes. Over 25% of examples can be misclassified by removing *less than 100* (example-specific) training images. Also, datamodels yield substantially better upper bounds on support size than baselines.

**Algorithm 1:** The (intractable) exhaustive search algorithm (top) and (efficient) datamodel-guided algorithm (bottom) for estimating data support. Note that in our linear datamodel setting, $G$ (the quantity on line 4 of the second algorithm) actually has a closed-form solution: it is the subset of the training set corresponding to the largest $k$ indices of the datamodel parameter $\theta$. See Appendix F.3 for implementation and setup details.

---

1: **procedure** EXHAUSTIVE(target ex. $x$, trainset $S$)
2:   **for** $k$ in $\{1, \ldots, d\}$ **do**
   *Try every subset of size $k$:*
3:    **for** $G$ in $2^S$ with $|G| = k$ **do**
4:     **if** $\mathbb{E}[f_\mathcal{A}(x; S \setminus G)] < 0$ **then**
5:      **return** $G$

1: **procedure** GUIDED($x$, $S$, datamodel $g_\theta$)
2:   $A \leftarrow []$
3:   **for** $k \in \{10, 20, 40, \ldots\}$ **do**
   *Find subset $G_k$ of lowest predicted margin:*
4:    $G_k \leftarrow \arg\min_{|G|=k} g_\theta(S \setminus G)$
5:    Estimate $\mathbb{E}[f_\mathcal{A}(x; S \setminus G_k)]$ by re-training
6:    Append $(k, \mathbb{E}[f_\mathcal{A}(x; S \setminus G_k)])$ to $A$
   *Piecewise-linear interp. mapping $k$ to margin:*
7:   $h(\cdot) \leftarrow$ INTERPOLATE($A$)
8:   $\widehat{k} \leftarrow k$ for which $h(k) = 0$
   *Conservative estimate of data support:*
9:   **return** TOP-K($\theta, \widehat{k} \times 1.2$)

---

**Baselines.** To contextualize these findings, we compare our estimates of data support to a few natural baselines. We provide the exact comparison setup in Appendix F.3.1: in summary, each baseline technique

can be cast as a swap-in alternative to datamodels for guiding the data support search described above.

It turns out that every baseline we tested provides much looser estimates of data support (Figure 7). For example, even the best-performing baseline predicts that one would need to remove over 600 training images per test image to force misclassification on 20% of the test set[3]. In contrast, our datamodel-guided estimates indicate that removing 40 train examples is sufficient for misclassifying 20% of test examples.

**Removing versus mislabeling.** Note that the brittleness we consider in this section (i.e., brittleness to *removing* training examples) is substantively different than brittleness to *mislabeling* examples (as in *label-flipping attacks* [XXE12; KL17; RWR+20]). In particular, brittleness to removal indicates that there exists a small set of training images whose presence is *necessary* for correct classification of the target example (thus motivating the term "data support"). Meanwhile, label-flipping attacks can succeed even when the target example has a large data support, as (consistently) mislabeling a set of training examples provide a much stronger signal than simply removing them. Nevertheless, we can easily adapt the above experiment to test brittleness to mislabeling—we do so in Appendix F.4. As one might expect, test predictions are even *more* brittle to data mislabeling than removal—for 50% of the CIFAR-10 test set, mislabeling *35 target-specific training examples* suffices to flip the corresponding prediction (see Figure F.3 for a CDF).

### 4.1.2 Predicting data counterfactuals

As we have already seen, a simple application of datamodels as a proxy for model training (on *arbitrary* subsets of the training set) enabled us to identify brittle predictions. We now demonstrate another, more intricate application of datamodels as a proxy for end-to-end training: predicting *data counterfactuals*.

For a fixed target example $x$, and a specific subset of the training set $R(x) \subset S$, a data counterfactual is the causal effect of removing the set of examples $R(x)$ on model outputs for $x$. In terms of our notation, this effect is precisely

$$\mathbb{E}\left[f_{\mathcal{A}}(x; S) - f_{\mathcal{A}}(x; S \setminus R(x))\right].$$

Such data counterfactuals can be helpful tools for finding brittle predictions (as in the previous subsection), estimating *group influence* (as done by [KAT+19] for linear models), and more broadly for understanding how training examples combine (through the lens of the model class) to produce test-time predictions.

**Estimating data counterfactuals.** Just as in the last section, we again use datamodels beyond the supervised learning regime in which they were developed. In particular, we predict the outcome of a data counterfactual as

$$g_\theta(\mathbf{1}_S) - g_\theta(\mathbf{1}_{S \setminus R(x)}),$$

where again $g_\theta$ is the datamodel for a given target example of interest. Since $g_\theta$ is a linear function in our case, the above *predicted data counterfactual* actually simplifies to

$$\theta^\top \mathbf{1}_S - \theta^\top \mathbf{1}_{S \setminus R(x)} = \theta^\top \mathbf{1}_{R(x)}.$$

Our goal now is to demonstrate that datamodels are useful predictors of data counterfactuals across a variety of removed sets $R(x)$. To accomplish this, we use a large set of target examples. Specifically, for each such target example, we consider different subset sizes $k$; for each such $k$, we use a variety of heuristics to select a set $R(x)$ comprising $k$ "examples of interest." These heuristics are:

(a) setting $R(x)$ to be the nearest $k$ training examples to the target example $x$ in terms of *influence score* [KL17], *TracIn score* [PLS+20], or *distance in pre-trained representation space* [BCV13][4];

(b) setting $R(x)$ to be the *maximizer* of the datamodel-predicted counterfactual, i.e.,

$$R(x) = \arg\max_{|R|=k} g_\theta(S) - g_\theta(S \setminus R) = \arg\max_{|R|=k} \theta^\top \mathbf{1}_R.$$

---

[3]Moreover, the data support estimates derived from the baselines are only "certifiable" in the above-described sense (see the beginning of the "Results" paragraph) for 60% of the 300 test examples we study (as opposed to 90% for datamodel-derived estimates).
[4]Note that these methods are precisely the ones used as baselines in the previous section.

(Note that since our datamodels are linear, this simplifies to excluding the training examples corresponding to the top $k$ coordinates of the datamodel parameter $\theta$.)

(c) setting $R(x)$ to be the training images corresponding to the *bottom* (i.e., most negative) $k$ coordinates of the datamodel weight $\theta$.

We consider six values of $k$ (the size of the removed subset) ranging from 10 to 1280 examples (i.e., $0.02\% - 2.6\%$ of the training set). Thus, the outcome of our procedure is, for each target example, both *true* and *datamodel-predicted* data counterfactuals for *30 different training subsets $R(x)$* (six values of $k$ and five different heuristics).

**Results.** In Figure 8, we plot datamodel-predicted data counterfactuals against true data counterfactuals, aggregating across all target examples $x$, values of $k$, and selection heuristics for $R(x)$. We find a strong correlation between these two quantities. In particular, across all factors of variation, predicted and true data counterfactuals have Spearman correlation $\rho = 0.98$ and $\rho = 0.94$ for CIFAR-10 and FMoW respectively. In fact, the two quantities are correlated roughly *linearly*: we obtain (Pearson) correlations of $r = 0.96$ (CIFAR-10) and $r = 0.90$ (FMoW) between counterfactuals and their estimates on aggregate. Correlations are even more pronounced when restricting to any single class of removed sets (i.e., any single hue in Figure 8).



**Figure 8: Datamodels predict data counterfactuals.** Each point in the graphs above corresponds to a test example and a subset $R(x)$ of the original training set $S$, identified by a (color-coded) heuristic. The $y$-coordinate of each point represents the *true* difference, in terms of model output on $x$, between training on $S$, and training on $S \setminus R(x)$. The $x$-coordinate of each point represents the *datamodel-predicted* value of this quantity. We plot results for **(right)** CIFAR-10 and **(left)** FMoW. Datamodel predictions are predictive of the underlying counterfactuals, with Pearson coefficients $r$ being 0.96/0.90 for CIFAR/FMoW respectively. Predictions are computed with datamodels estimated with $\alpha = 0.5$ for CIFAR-10 and $\alpha = 0.75$ for FMoW (cf. Appendix F.9 other values of $\alpha$). See Appendix F for more experimental details and results.

**Limits of datamodel predictions.** We have seen that datamodels accurately predict the outcome of many natural data counterfactuals, despite only being constructed to predict outcomes for random subsets of a fixed size ($\alpha \cdot d$ for $\alpha \in (0, 1)$ and $d$ the training set size). Of course, due to both estimation error (i.e., we might not have trained enough models to identify *optimal* linear datamodels) and misspecification error (i.e., the optimal datamodel might not be linear), we don't expect a perfect correspondence between datamodel-predicted outputs $g_\theta(\mathbf{1}_{S'})$ and true outputs $f(x; S)$ for *all* $2^d$ possible subsets of the training set. Indeed, this is part of the reason why we estimated datamodels for several values of $\alpha$, only one of which is shown in Figure 8. As shown in Appendix F.9, datamodels estimated for other values of $\alpha$ still display strong correlation between true and predicted model outputs, but behave qualitatively differently than the ones shown above (i.e., each value of $\alpha$ is better or worse at predicting the outcomes of certain types of counterfactuals).

12

## 4.2 Using datamodels to find similar training examples

We now turn to another application of datamodels: identifying training examples that are similar to a given test example. One can use this primitive to identify issues in datasets such as duplicated training examples [LIN+21] or train-test leakage [BD20] (test examples that have near-duplicates in the training set).

Recall that in our instantiation of the framework, datamodels predict model output (for a fixed target example) as a *linear* function of the presence of each training example in the training set. That is, we predict the output of training on a subset $S'$ of the training set $S$ as

$$g_\theta(\mathbf{1}_{S'}) = \theta^\top \mathbf{1}_{S'}.$$

A benefit of parameterizing datamodels as simple linear functions is that we can use the magnitude of the coordinates of $\theta$ to ascertain *feature importance* [GE03]. In particular, since in our case each feature coordinate (i.e., each coordinate of $\mathbf{1}_{S'}$) actually represents the presence of a particular training example, we can interpret the highest-magnitude coordinates of $\theta$ as the indices of the training examples whose presence (or absence) is most predictive of model behavior (again, on the fixed target example in context).

We now show that these high-magnitude training examples (a) they visually resemble the target image, yielding a method for finding similar training examples to a given target; and (b) as a result, datamodels can automatically detect train-test leakage.

> **Use Case 2** (Train-test similarity). *For a test example x with a linear datamodel $g_\theta$, we can interpret the training examples corresponding to the highest-magnitude coordinates of $\theta$ as the "nearest neighbors" of x.*

### 4.2.1 Finding similar training examples

Motivated by the feature importance perspective described above, we visualize (in Figures 9 and G.1) a random set of target examples from the CIFAR-10 test set together with the CIFAR-10 training images that correspond to the highest-magnitude datamodel coordinates for each test image.

**Results.** We find that for a given target example, the highest-magnitude datamodel coordinates—both positive and negative—consistently correspond to visually similar training examples.
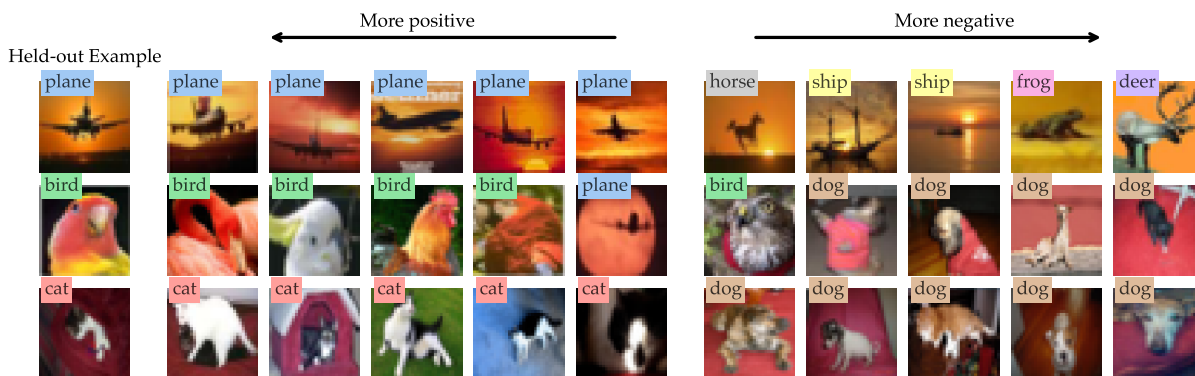


**Figure 9: Large datamodel weights correspond to similar images.** Randomly choosing test examples and visualizing their most negative- and positive-weight examples for $\alpha = 50\%$, we find that large magnitude train examples share similarities with their test examples. Top negative weights generally correspond to visually similar images from other classes. See Appendix G for more examples.

Furthermore, the exact training images that are surfaced by looking at high-magnitude weights differ depending on the subsampling parameter $\alpha$ that we use while constructing the datamodels. (Recall from Section 2 that $\alpha$ controls the size of the random subsets used to collect the datamodel training set—a datamodel estimated with parameter $\alpha$ is constructed to predict outcomes of training on random training

subsets of size $\alpha \cdot d$, where $d$ is the training set size.) In Figure 10 (and G.2), we consider a pair of target examples from the CIFAR-10 test set, and, for each target example, compare the top training images from two different datamodels: one estimated using $\alpha = 10\%$, and the other using $\alpha = 50\%$. We find that in some cases (e.g., Figure 10 left), the $\alpha = 10\%$ datamodel identifies training images that are highly similar to the target example but do not correspond to the highest-magnitude coordinates for the $\alpha = 50\%$ datamodel (in other cases, the reverse is true). Our hypothesis here—which we expand upon in Section 5—is that datamodels estimated with lower $\alpha$ (i.e., based on smaller random training subsets) find train-test relationships driven by larger groups of examples (and vice-versa).



**Figure 10: Datamodels corresponding to different $\alpha$ surface qualitatively different images**. For each target example (taken from the CIFAR-10 test set), we consider two different datamodels: one estimated with $\alpha = 10\%$ (i.e., constructed to predict model outputs on the target example after training on random 10% subsets of the CIFAR-10 training set), and the other estimated with $\alpha = 50\%$. For each datamodel, we visualize the training examples corresponding to the largest coordinates of the parameter vector $\theta$. On the left we see an example where the datamodel estimated with $\alpha = 10\%$ (top row) detects a set of near-duplicates of the target example that the $\alpha = 50\%$ datamodel (bottom row) does not identify. See Appendix G for more examples.

**Influence functions.** Another method for finding similar training images is *influence functions*, which aim to estimate the effect of removing a single training image on the loss (or correctness) for a given test image. A standard technique from robust statistics [HRR+11] (applied to deep networks by Koh and Liang [KL17]) uses first-order approximation to estimate influence of each training example. We find (cf. Appendix Figure G.3), that the high-influence and low-influence examples yielded by this approximation (and similar methods) often fail to find similar training examples for a given test example (also see [BPF21; HYH+21]).

Another approach based on *empirical* influence approximation was used by Feldman and Zhang [FZ20], who (successfully) use their estimates to identify similar train-test pairs in image datasets as we do above. We discuss empirical influence approximation and its connection with datamodeling in Section 6.1.

### 4.2.2 Identifying train-test leakage

We now leverage datamodels' ability to surface training examples similar to a given target in order to identify *same-scene* train-test leakage: cases where test examples are near-duplicates of, or clearly come from the same scene as, training examples. Below, we use datamodels to uncover evidence of train-test leakage on both CIFAR and FMoW, and show that datamodels outperform a natural baseline for this task.

**Train-test leakage in CIFAR.** To find train-test leakage in CIFAR-10, we collect ten candidate training examples for each image in the CIFAR-10 test set—the ones corresponding to the ten largest coordinates of the test example's datamodel parameter. We then show crowd annotators (using Amazon Mechanical Turk) tasks that consist of a random CIFAR-10 test example accompanied by its candidate training examples. We ask the annotators to label any of the candidate training images that constitute instances of same-scene leakage (as defined above). We show each task (i.e., each test example) to multiple annotators, and compute the "annotation score" for each of the test example's candidate training examples as the fraction of annotators who marked it as an instance of leakage. Finally, we compute the "leakage score" for each test example as

the highest annotation score (over all of its candidate train images). We use the leakage score as a proxy for whether or not the given image constitutes train-test leakage.

In Figure 11, we plot the distribution of leakage scores over the CIFAR-10 test set, along with random train-test pairs stratified by their annotation score. As the annotation score increases, pairs (qualitatively) appear more likely to correspond to leakage (see Appendix H for more pairs). Furthermore, *roughly 10% of test set images were labeled as train-test leakage by over half of the annotators that reviewed them*.



**Figure 11: Finding CIFAR train-test leakage candidates with datamodels.** Nine MTurk annotators view each test image alongside the train images with largest datamodel weight. The annotators then select the train images judged as belonging to the same scene as the test image. We measure the *annotation score* for a given (train, test) pair as the frequency with which annotators selected the pair as clearly coming from the same scene. The *leakage score* for a test image is defined as the maximum annotation score over all of its candidate train images. (See Appendix Section H for a more detailed setup.) **(Left)** Histogram of the leakage score for each image of the CIFAR test set. **(Right)** Train-test pairs stratified by their leakage score. A majority of annotators (annotation score of more than $\frac{1}{2}$) consider 10% of the test set as train-test leakage. Many of these pairs are near-duplicate; see more examples in Appendix H.
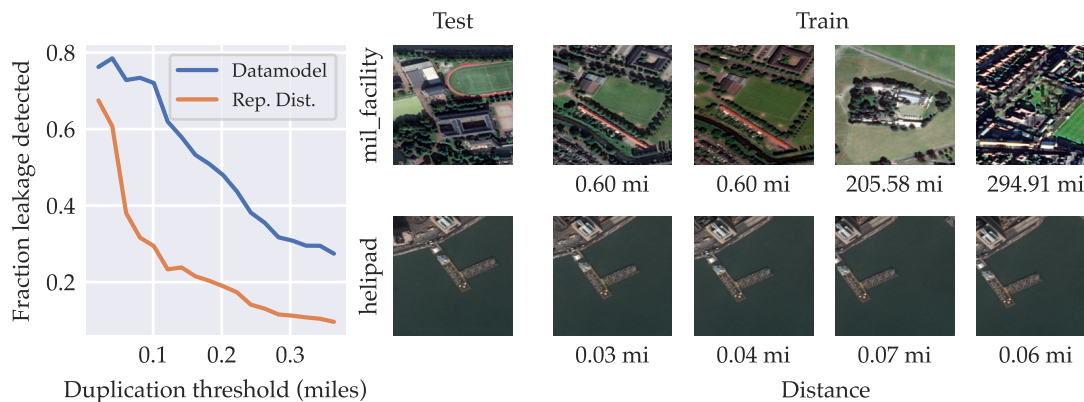


**Figure 12: Datamodels detect same-scene train-test leakage on FMoW.** FMoW images are annotated with geographic coordinates. For any distance $d$, we call a test image $x$ "leaked" if it is within $d$ miles of *any* training image $x_s$. A leaked test image $x$ is considered "detected" if the corresponding training image $x_s$ has one of the 10 largest datamodel weights for $x$. **(Left)** With $d$ on the x-axis, we plot the fraction of leaked test images that are also detected. As a baseline, we replace datamodel weights with (negative) distances in neural network representation space. **(Right)** for two test examples (top: random; bottom: selected), we show the most similar train examples (by datamodel weight), labeled by their distance to the test example.

15

**Train-test leakage in FMoW.** To identify train-test leakage on FMoW, we begin with the same candidate-finding process that we used for CIFAR-10. However, FMoW differs from CIFAR in that the examples (satellite images labeled by category, e.g., "port" or "arena") are annotated with *geographic coordinates*. These coordinates allow us to avoid crowdsourcing—instead, we compute the geodesic distance between the test image and each of the candidates, and use a simple threshold $d$ (in miles) to decide whether a given test example constitutes train-test leakage.

Furthermore, we can calculate a "ground-truth" number of train-test leakage instances by counting the test examples whose *geodesic* nearest-neighbor in the training set is within the specified threshold $d$.[5] Comparing this ground truth to the number of instances of leakage found within the candidate examples yields a qualitative measure of the efficacy of our method (i.e., the quality of candidates we generate).

In Figure 12, we plot this measure of efficacy (# instances found / # ground truth) as a function of the threshold $d$, and also visualize examples images from the FMoW test set together with their corresponding datamodel-identified training set candidates. To put our quantitative results into context, we compare the efficacy of candidates derived from top datamodel coordinates (i.e., the ones we use here and for CIFAR-10) to that of candidates derived from *nearest neighbors* in the representation space of a pretrained neural network [BCV13; ZIE+18] (examining such nearest neighbors is a standard way of finding train-test leakage, e.g., used by [BD20] to study CIFAR-10 and CIFAR-100). Datamodels consistently outperform this baseline.

## 4.3 Using datamodels as a feature embedding

Sections 4.1 and 4.2 illustrate the utility of datamodels on a *per-example* level, i.e., for predicting the outcome of training on arbitrary training subsets and evaluating on a specific target example, or for finding similar training images (again, to a specific target). We'll conclude this section by demonstrating that datamodels can also help uncover *global structure* in datasets of interest.

Key to this capability is the following shift in perspective. Consider a target example $x$ with corresponding *linear* datamodel $g_\theta$, and recall that the datamodel is parameterized by a vector $\theta \in \mathbb{R}^d$, where $d$ is the training set size. We optimized $\theta$ to minimize the squared-error of $g_\theta$ (8) when predicting the outcome of training on random training subsets and evaluating on the target $x$. Now, however, instead of viewing the vector $\theta$ as just a parameter of the predictor $g_\theta$, we cast it as a *feature representation* for the target example itself, i.e., a *datamodel embedding* of $x$ into $\mathbb{R}^d$. Since the datamodel $g_\theta$ is a linear function of the presence of each point in the training set, each coordinate of this "datamodel embedding" corresponds to a weight for a specific training example. One can thus think of a datamodel embedding as a feature vector that represents a target example $x$ in terms of how predictive each training example is of model behavior on $x$.

Critically, the coordinates of a datamodel embedding have a *consistent* interpretation across datamodel embeddings, even for different target examples. That is, we expect similar target examples to be acted upon similarly by the training set, and thus have similar datamodel embeddings. In the same way, if model performance on two unrelated target examples is driven by two disjoint sets of training examples, their datamodel embeddings will be orthogonal. This intuition suggests that by embedding an entire dataset of examples $\{x_i\}$ as a set of feature vectors $\{\theta_i \in \mathbb{R}^d\}$, we may be able to uncover structure in the set of examples by looking for structure in their datamodel embeddings, i.e., in the (Euclidean) space $\mathbb{R}^d$.

In this section we demonstrate, through two applications, the potential for such datamodel embeddings to discover dataset structure in this way. In Section 4.3.2, we use datamodel embeddings to partition datasets into disjoint clusters, and in Section 4.3.1 we use principal component analysis to get more fine-grained insights into dataset structure. To emphasize our shift in perspective (i.e., from $\theta$ being just a parameter of a datamodel $g_\theta$, to $\theta$ being an embedding for the target example $x$), we introduce an *embedding function* $\varphi(x) \mapsto \theta$ which maps a particular target example to the weights of its corresponding datamodel.

---

> **Use Case 3** (Datamodel embeddings). *We can use datamodels as a way to embed any given target example into the same (Euclidean) space $\mathbb{R}^d$, where d is the training set size. Specifically, we can define the* datamodel embedding $\varphi(x)$ *for a target example x as the weight vector $\theta \in \mathbb{R}$ of the datamodel corresponding to x.*

---

[5]It turns out that despite having already been de-duplicated, about 20% and 80% of FMoW test images are within 0.25 and 2.6 miles of a training image, respectively—see Appendix Figure H.3.

### 4.3.1 Spectral clustering with datamodel embeddings

We begin with a simple application of datamodel embeddings, and show that they enable high-quality clustering. Specifically, given two examples $x_1$ and $x_2$, datamodel embeddings induce a natural *similarity measure* between them:

$$d(x_1, x_2) := K(\varphi(x_1), \varphi(x_2)), \tag{13}$$

where we recall that $\varphi(\cdot)$ is the *datamodel embedding function* mapping target examples to the weights of their corresponding datamodels, and $K(\cdot, \cdot)$ is any kernel function (below, we use the RBF kernel)[6]. Taking this even further, for a set of $k$ target examples $\{x_1, \ldots, x_k\}$, we can compute a full *similarity matrix* $A \in \mathbb{R}^{k \times k}$, whose entries are

$$A_{ij} = d(x_i, x_j). \tag{14}$$

Finally, we can view this similarity matrix as an *adjacency matrix* for a (dense) graph connecting all the examples $\{x_1, \ldots x_k\}$: the edge between two examples will be $d(x_i, x_j)$, which is in turn the kernelized inner product between their two datamodel weights. We expect similar examples to have high-weight edges between them, and unrelated examples to have (nearly) zero-weight edges between them.

Such a graph unlocks a myriad of graph-theoretic tools for exploring datasets through the lens of datamodels (e.g., cliques in this graph should be examples for which model behavior is driven by the same subset of training examples). However, a complete exploration of these tools is beyond the scope of our work: instead, we focus on just one such tool: spectral clustering.

At a high level, spectral clustering is an algorithm that takes as input any similarity graph $G$ as well as the number of clusters $C$, and outputs a partitioning of the vertices of $G$ into $C$ disjoint subsets, in a way that (roughly) minimizes the total weight of inter-cluster edges. We run an off-the-shelf spectral clustering algorithm on the graph induced by the similarity matrix $A$ above for the images in the CIFAR-10 test set. The result (Figure 13 and Appendix I) demonstrates a simple unsupervised method for uncovering subpopulations in datasets.
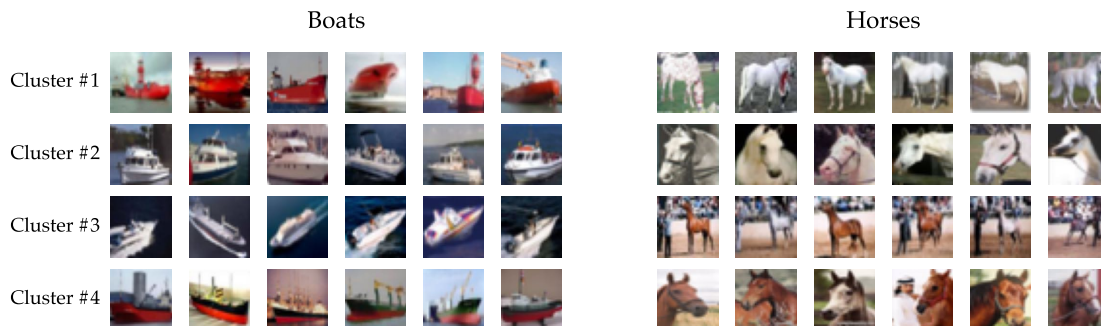


**Figure 13: Spectral clustering on datamodel embeddings finds subpopulations.** For each CIFAR-10 class, we first compute a similarity score between all datamodel embeddings (we use $\alpha = 20\%$ datamodels), then run spectral clustering on the resulting matrix. We show the top clusters with the lowest average distance to the cluster center (in the embedding space); each row shows six random images from the given cluster. Each cluster seems to correspond to a specific subpopulation with shared, distinctive visual features. See Figures I.1 and I.2 for more examples from other classes and comparison across $\alpha$.

### 4.3.2 Analyzing datamodel embeddings with PCA

We observed above that datamodel embeddings encode enough information about their corresponding examples to cluster them into (at least qualitatively) coherent groups. We now attempt to gain even further

---

[6]A kernel function $K(\cdot, \cdot)$ is a similarity measure that computes the inner product between its two arguments in a transformed inner product space (see [SC04] for an introduction). The RBF kernel is $K(v_1, v_2) = \exp\{-\|v_1 - v_2\|^2/2\sigma^2\}$

insight into the structure of these datamodel embeddings, in the hopes of shedding light on the structure of the underlying dataset itself.

Datamodel embeddings are both high-dimensional and sparse, making analyzing them directly (e.g., by looking at the variation of each coordinate) a daunting task. Instead, we leverage a canonical tool for finding structure in high-dimensional data: principal component analysis (PCA).

PCA is a dimensionality reduction technique which—given a set of embeddings $\{\varphi(x_i) \in \mathbb{R}^d\}$ and any $k \ll d$—returns a *transformation function* that maps any embedding $\varphi(x) \in \mathbb{R}^d$ to a new embedding $\widetilde{\varphi}(x) \in \mathbb{R}^k$, such that:

(a) each of the $k$ coordinates of the transformed embeddings is a (fixed) linear combination of the coordinates of the initial datamodel embeddings, i.e., $\widetilde{\varphi}(x) = M \cdot \varphi(x)$ for a fixed $k \times d$ matrix $M$;

(b) transformed embeddings preserve as much information as possible about the original ones. More formally, we find the matrix $M$ that allows us to *reconstruct* the given set of embeddings $\{\varphi(x_i) \in \mathbb{R}^d\}$ from their transformed counterparts with minimal error.

Note that in (a), the $i$-th coordinate of a transformed embedding is always the *same* linear combination of the corresponding original embedding (and thus, each coordinate of the transformed embedding has a concrete interpretation as a weighted combination of datamodel coefficients). The exact coefficients of this combination (i.e., the rows of the matrix $M$ above) are called the first $k$ *principal components* of the dataset.

We apply PCA to the collection of datamodel embeddings $\{\varphi(x_i) \in \mathbb{R}^d\}_{i=1}^d$ for the CIFAR-10 training set, and use the result to to compute new $k$-dimensional embeddings for each target example in both the training set and the test set (i.e., by computing each target example's datamodel embedding then transforming it to an embedding in $\mathbb{R}^k$). We can then look at each coordinate in the new, much more manageable ($k$-dimensional) embeddings.[7]

**Coordinates identify subpopulations.** Our point of start in analyzing these transformed embeddings is to examine each transformed coordinate separately. In particular, in Figure 14 we visualize, for a few sample coordinate indices $i \in [k]$, the target examples whose transformed embeddings have particularly high or low values of the given coordinate (equivalently, these are the target examples whose datamodel embeddings have the highest or lowest projections onto the $i$-th principal component). We find that:

(a) the examples whose transformed embeddings have a large $i$-th coordinate all (visually) share a common feature: e.g., the first-row images in Figure 14 share similar pose and color composition;

(b) this (visual) feature is consistent across both train and test set examples[8]; and

(c) for a given coordinate, the most positive images and most negative images (i.e., the left and right side of each row of Figure 14, respectively) either (a) have a differing label but share the same common feature or (b) have the same label but differ along the relevant feature.

**Principal components are *model-faithful*.** In Appendix J, we verify that not only are the groups of images found by PCA visually coherent, they are in fact rooted in how the model class makes predictions. In particular, we show that one can find, for any coordinate $i \in [k]$ of the transformed embedding, the training examples that are most important to that coordinate. Furthermore, retraining without these examples significantly decreases (increases) accuracy on the target examples with the most positive (negative) coordinate $i$, suggesting that the identified principal components actually reflect model class behavior.

---

[7]We first *normalize* each datamodel embedding before transforming them (i.e., we transform $\varphi(x)/\|\varphi(x)\|$).

[8]Recall that we computed the PCA transformation to preserve the information in only the *training set* datamodel embeddings. Thus, this result suggests that the transformed embeddings computed by PCA are not "overfit" to the specific examples that we used to compute it.

**Figure 14: PCA on datamodel embeddings.** We visualize the top three principal components (PCs) and a randomly selected PC from the top 100. In the $i$-th row, the left-most (right-most) images are those whose datamodel embeddings have the highest (lowest) normalized projections onto the $i$-th principal component $v_i$. Highest magnitude images along each direction share qualitative features; moreover, images at opposite ends suggest a *feature tradeoff*—a combination of images in the training set that helps accuracy on one subgroup but hurts accuracy on the other. See Figure J.5 for more datamodel PCA components.

### 4.3.3   Advantages over penultimate-layer embeddings

In the context of deep neural networks, the word "embedding" typically refers to features extracted from the penultimate layer of a fixed pre-trained model (see [BCV13] for an overview). These "deep representations" can serve as an effective proxy for visual similarity [ZIE+18; BD20], and also enable a suite of applications such as clustering [GGT+17] and feature visualization [BBC+07; ARS+15; OMS17; EIS+19].

Here, we briefly discuss a few advantages of datamodel-based embeddings over their standard penultimate layer-based counterparts.

- **Axis-alignment**: datamodel embeddings are *axis-aligned*—each embedding component directly corresponds to index into the training set, as opposed to a more abstract or qualitative concept. As a corollary, aggregating or comparing different datamodel embeddings for a given dataset is straightforward, and does not require any alignment tools or additional heuristics. This is not the case for network-based representations, for which the right way to combine representations—even for two models of the same architecture—is still disagreed upon [KNL+19; BNB21]. In particular, we can straightforwardly compare datamodel embeddings across different target examples, model architectures, training paradigms, or even datamodel estimation techniques—as long as the set of training examples being stays the same, any resulting datamodel has a uniform interpretation.

- **Richer representation**: the space of datamodel embeddings seems significantly richer than that of standard representation space. In particular, Appendix Figure J.1 shows that for standard representation space, *ten linear directions* suffice to capture *90% of the variation* in training set representations.

The "effective dimension" of datamodel representations is much higher, with the top *500 principal components* explaining only *50% of the variation* in training set datamodel embeddings. This difference manifests qualitatively when we redo our PCA study on standard representations (Appendix Figure J.4): principal components beyond the 10th lack both the perceptual quality and train-test consistency exhibited by those of datamodel embeddings (e.g., for datamodels even the 76th principal component, shown in Figure 14, exhibits these qualities).

- **Ingrained causality**: datamodel embeddings inherently encode information about how the model class generalizes. Indeed, in Section 4.3.2 we verified via counterfactuals that insights extracted from the principal components of $\Theta$ actually reflect underlying model class behavior.

# 5   Discussion: The role of the subsampling fraction $\alpha$



**Figure 15: Datamodels capture data relationships at varying levels of granularity.** We illustrate the role of *subsampling fraction $\alpha$* of datamodels by considering a nearest-neighbor classifier in two dimensions. In the datamodel for the target example ($\star$, yellow), the red (blue) examples have positive (negative) weights, with the shade indicating the magnitude. At large values of $\alpha$ (right), the model identifies only local relationships. Meanwhile, at small values of $\alpha$ (left), we can identify more global relationships, but at the cost of granularity. Intermediate values of $\alpha$ (middle) provide a smooth tradeoff between these two regimes.

We have used datamodels estimated using several choices of the subsampling fraction $\alpha$, and saw that the value of $\alpha$ corresponding to the most useful datamodels can vary by setting. In particular, the visualizations in Figure 10 suggest that datamodels estimated with lower $\alpha$ (i.e., based on smaller random training subsets) find train-test relationships driven by larger groups of examples (and vice-versa). Here, we explore this intuition further using thought experiment, toy example, and numerical simulation. Our goal is to intuit how different choices of $\alpha$ can lead to substantively different datamodels.

First, consider the task of estimating a datamodel for a prototypical image $x$—for example, a plane on a blue sky background. As $\alpha \to 1$, the sets $S_i$ sampled from $\mathcal{D}_S$ are relatively large—if these sets have enough other images of planes on blue skies, we will observe little to no variation in $f_{\mathcal{A}}(x; S_i)$, since any predictor trained on $S_i$ will perform very well on $x$. As a result, a datamodel for $x$ estimated with $\alpha \to 1$ may assign very little weight to any particular image, even if in reality their *total* effect is actually significant.

Decreasing $\alpha$, then, offers a solution to this problem. In particular, we allow the datamodel to observe cases where entire *groups* of training examples are not present, and *re-distribute* the corresponding effect back to the constituents of the group (i.e., assigning them all a share of the weight).

Now, consider a highly atypical yet correctly classified example, whose correctness relies on just the presence of just a few images from the training set. In this setting, datamodels estimated with a small value of $\alpha$ may be unable to isolate these training points, since they will constantly distribute variation in $f_{\mathcal{A}}(x; S_i)$ among a large group of non-present images. Meanwhile, using a large value of $\alpha$ allows the estimated datamodel to place weight on the correct training images (since $x$ will be classified correctly until some of the important training images are not present in $S_i$).

In line with this intuition, decreasing $\alpha$ in Figure 15 (i.e., moving from right to left) leads to datamodels that assign weight to increasingly large neighborhoods of points around the target input. This example and the above reasoning lead us to hypothesize that larger (respectively, smaller) $\alpha$ are better-suited to

20

cases where model predictions are driven by smaller (respectively, larger) groups of training examples. In Appendix B, we perform a more quantitative analysis of the role of $\alpha$, this time by studying an underdetermined linear regression model on data that is organized into overlapping *subpopulations*. Our findings in this setting (see Figure B.1) mirror our intuition thus far—in particular, smaller values of $\alpha$ result in datamodels that were more predictive on *larger* subpopulations in the training set, whereas higher values of $\alpha$ tended to work better *smaller* subpopulations.

# 6 Related work

Datamodels build on a rich and growing body of literature in machine learning, statistics, and interpretability. In this section, we illustrate some of the connections to these fields, highlight a few of the most closely related works to ours.

## 6.1 Connecting datamodeling to empirical influence estimation

We start by discussing the particularly important connection between datamodels and another well-studied concept that has recently been applied to the machine learning setting: influence estimators. In particular, a recent line of work aims to compute the *empirical influence* [HRR+11] of training points $x_i$ on predictions $f(x_j)$, i.e.,

$$\text{Infl}[x_i \to x_j] := \mathbb{P}\left(\text{model trained on } S \text{ is correct on } x_j\right) - \mathbb{P}\left(\text{model trained on } S \setminus \{x_i\} \text{ is correct on } x_j\right),$$

where randomness is taken over the training algorithm. Evaluating these influence functions naively requires training $C \cdot d$ models where $d$ is again the size of the train set and $C$ is the number of samples necessary for an accurate empirical estimate of the probabilities above. To circumvent this prohibitive sample complexity, a recent line of work has proposed approximation schemes for $\text{Infl}[x_i \to x_j]$. We discuss these approximations (and their connection to our work) more generally in Section 6.2, but here we focus on a specific approximation used by Feldman and Zhang [FZ20] (and in a similar form, by [GZ19] and [JDW+19])[9]:

$$\widehat{\text{Infl}}[x_i \to x_j] = \mathbb{P}_{S \sim \mathcal{D}_S}\left(\text{model trained on } S \text{ is correct on } x_j | x_i \in S\right)$$
$$- \mathbb{P}_{S \sim \mathcal{D}_S}\left(\text{model trained on } S \text{ is correct on } x_j | x_i \notin S\right). \tag{15}$$

This estimator improves sample efficiency by reusing the same set of models to compute influences between different input pairs. More precisely, Feldman and Zhang [FZ20] show that the size of the random subsets trades off sample efficiency (model reuse is maximized when the subsets are exactly half the size of the training set, since this maximizes the number of samples available to estimate each term in (15)) and accuracy with respect to the true empirical influence (which is maximized as the subsets $S_i$ get larger). Despite its different goal, formulation, and estimation procedure, it turns out that we can cast the difference-of-probabilities estimator (15) above as a rescaled datamodel (in the infinite-sample limit). In particular, in Appendix K.1 we show:

**Lemma 1.** *Fix a training set $S$ of size $n$, and a test example $x$. For $i \in [m]$, let $S_i$ be a random variable denoting a random 50%-subset of the the training set $S$. Let $\boldsymbol{w}_{infl} \in \mathbb{R}^n$ be the estimated empirical influences (15) onto $x$ estimated using the sets $S_i$. Let $\boldsymbol{w}_{OLS}$ be the least-squares estimator of whether a particular model will get image $x$ correct, i.e.,*

$$\boldsymbol{w}_{OLS} := \arg\min_w \frac{1}{m} \sum_{i=1}^{m} \left(w^\top \boldsymbol{z_i} - \mathbf{1}\{\text{model trained on } S_i \text{ correct on } x\}\right)^2, \qquad \text{where } \boldsymbol{z}_i = 2 \cdot \mathbf{1}_{S_i} - \mathbf{1}_n.$$

*Then, as $m \to \infty$,*

$$\left\|(1 + 2/n)\boldsymbol{w}_{OLS} - \frac{1}{2}\boldsymbol{w}_{infl}\right\|_2 \to 0.$$

---

[9]In fact, (15) is ubiquitous—e.g., in causal inference, it is called the *average treatment effect* of training on $x_i$ on the correctness of $x_j$.

We illustrate this result quantitatively in Appendix K and perform an in-depth study of influence estimators as datamodels. As one might expect given their different goal, influence estimates significantly underperform explicit datamodels in terms of predicting model outputs with respect to every metric we studied (Table K.1, Figure K.1). We then attempt to explain this performance gap and reconcile it with Lemma 1 in terms of the estimation algorithm (OLS vs. LASSO), scale (number of models trained), and output function (0/1 loss vs. margins).

In addition to forging a connection between datamodels and influence estimates, this result also provides an alternate perspective on the parameter $\alpha$. Specifically, in light of our discussion in Section 5, it suggests that $\alpha$ may control the *kinds* of correlations that are surfaced by empirical influence estimates.

## 6.2  Other connections

**Influence functions and instance-based explanations.**   Above, we contrasted datamodels with *empirical influence functions*, which measure the counterfactual effect of removing individual training points on a given model output. Specifically, in that section and the corresponding Appendix K, we discussed the sub-sampled influence estimator of Feldman and Zhang [FZ20], who use influences to study the memorization behavior of standard vision models. We now provide a brief overview of a variety of other methods for influence estimation developed in prior works.

First-order influence functions are a canonical tool in robust statistics that allows one to approximate the impact of removing a data point on a given parameter without re-estimating the parameter itself [HRR+11]. Koh and Liang [KL17] apply influence functions to both a variety of classical machine learning models and to penultimate-layer embeddings from neural network architectures, to trace model's predictions back to individual training examples. In classical settings (namely, for a logistic regression model), Koh et al. [KAT+19] find that influence functions are also useful for estimating the impact of *groups* of examples. On the other hand, Basu et al. [BPF21] finds that approximate influence functions scale poorly to deep neural network architectures; and Feldman and Zhang [FZ20] argue that understanding the dynamics of the penultimate layer is insufficient for understanding deep models' decision mechanisms. Other methods for influence approximation (or more generally, instance-level attribution) include gradient-based methods [PLS+20] and metrics based on representation similarity [YKY+18; CGF+19]—see [HYH+21] for a more detailed overview. Finally, another related line of work [GZ19; JDW+19; WZJ+21] uses *Shapley values* [Sha51] to assign a value to datapoints based on their contribution to some *aggregate* metric (e.g., test accuracy).

As discussed in Section 6.1, datamodels serve a different purpose to influence functions—the former constructs an explicit statistical model, whereas the latter measures the counterfactual value of each training point. Nevertheless, we find that wherever efficient influence approximations and datamodels are quantitatively comparable (e.g., see Section 4.1 or Appendix K) datamodels predict model behavior better.

**Pixel-space surrogate models for interpretability.**   Datamodels are essentially surrogate models for the function mapping training data to predictions. Surrogate models from *pixel-space* to predictions are popular tools in machine learning interpretability [RSG16; LL17; SHS+19]. For example, LIME [RSG16] constructs a local linear model mapping test images to model predictions. Such surrogate models try to understand, for a *fixed* model, how the features of a given test example change the prediction. In contrast, datamodels hold the test example fixed and instead study how the images present in the training set change the prediction.

In addition to the advantages of our data-based view stated in Section 1, datamodels have two further advantages over pixel-level surrogate models: (a) a clear notion of *missingness* (i.e., it is easy to remove a training example but usually hard to "remove" pixels [SLL20; JSW+22]); and (b) *globality* of predictions— pixel-level surrogate models are typically accurate within a small neighborhood of a given input in pixel space, whereas datamodels model entire distribution over subsets of the training set, and remain useful both on- and off-distribution.

In other contexts, surrogate models are also used to evaluate data points for active learning and coreset selection [LC94; CYM+20]. Coleman et al. [CYM+20] find that shallow neural networks trained with fewer epochs can be a good proxy for a larger model when evaluating data for these applications.

**Model understanding beyond fixed weights.**   Recall (from Section 1) that datamodels are, in part, inspired by the fact that re-training deep neural networks using the same data and model class leads to

models with similar accuracies but vastly different individual predictions. This phenomenon has been observed more broadly. For example, Sellam et al. [SYW+21] make this point explicitly in the context of BERT [DCL+19] pre-trained language models. Similarly, Nakkiran and Bansal [NB20] make note of this non-determinism for networks trained on the same training *distribution* (but not the same data), while Jiang et al. [JNB+21] find that the same is true for networks trained on the same exact data. D'Amour et al. [DHM+20] find that on out-of-distribution data even overall accuracy is highly random. More closely to the spirit to our work, Zhong et al. [ZGK+21] find that non-determinism of individual predictions poses a challenge for comparing different model architectures. (They also propose a set of statistical techniques for overcoming this challenge.) More traditionally, the non-determinism is leveraged by Bayesian [Nea96] and ensemble methods [LPB17], which use a distribution over model weights to improve aspects of inference such as calibration of uncertainty.

**Learning and memorization.** Recent work (see [ZBH+16; Cha18; Fel19; BN20] and references therein) brings to light the interplay between learning and memorization, particularly in the context of deep neural networks. While memorization and generalization may seem to be at odds, the picture is more subtle. Indeed, Chatterjee [Cha18] builds a network of small lookup tables on small vision datasets to show that purely memorization-based systems can still generalize-well. Feldman [Fel19] suggests that memorization of atypical examples may be *necessary* to generalize well due to a long tail of subpopulations that arises in standard datasets. Feldman and Zhang [FZ20] find some empirical support for this hypothesis by identifying memorized images on CIFAR-100 and ImageNet and showing that removing them hurts overall generalization. Relatedly, Brown et al. [BBF+21] proves that for certain natural distributions, memorization of a large fraction of data, even data irrelevant to the task at hand, is necessary for close to optimal generalization. For state of the art models, recent works (e.g., [CLK+19; CTW+21]) show that one can indeed extract sensitive training data, indicating models' tendency to memorize.

Conversely, it has been observed that differentially private (DP) machine learning models—whose aim is precisely to avoid memorizing the training data—tend to exhibit poorer generalization than their memorizing counterparts [ACG+16]. Moreover, the impact on generalization from DP is disparate across subgroups [BPS19]. A similar effect has been noted in the context of neural network pruning [HCD+19]. Data-modeling may be a useful tool for studying these phenomena and, more broadly, the mechanisms mapping data to predictions for modern learning algorithms.

**Brittleness of conclusions.** A long line of work in statistics focuses on testing the *robustness* of statistical conclusions to the omission of datapoints. Broderick et al. [BGM21] study the robustness of econometric analyses to removing a (small) fraction of data. Their method uses a Taylor-approximation based metric to estimate the most influential subset of examples on some target quantity, similar in spirit to our use of datamodels to estimate data support for a target example (as in Figure 7). Datamodels may be a useful tool for extending such robustness analyses to the context of state-of-the-art machine learning models.

# 7 Future work

Our instantiation of the datamodeling framework yields both good predictors of model behavior and a variety of direct applications. However, this instantiation is fairly basic and thus leaves significant room for improvement along several axes. More broadly, datamodeling provides a lens under which we can study a variety of questions not addressed in this work. In this section, we identify (a subset of) these questions and provide connections to existing lines of work on them across machine learning and statistics.

## 7.1 Improving datamodel estimation

In Section 2, we outlined our basic procedure for fitting datamodels: we first sample subsets uniformly at random, then fit a sparse linear model from (the characteristic vectors of) training subsets to model outputs (margins) via $\ell_1$ regularization. We first discuss various ways in which this paradigm might be improved to yield even better predictions.

- **Correlation-aware estimation.** One key feature of our estimation methodology is that the same set of models is used to estimate datamodel parameters for an entire test set of images at once. This significantly reduces the sample complexity of estimating datamodels but also introduces a correlation between the errors in the estimated parameters. This correlation is driven by the fact that model outputs are not i.i.d. across inputs—for example, if on a picture of a dog $x$ a given model has very large output (compared to the "average" model, i.e., if $f_{\mathcal{A}}(x; S_i) - \mathbb{E}[f_{\mathcal{A}}(x; S_i)]$ is large), the model is also more likely to have large output on another picture of a dog (as opposed to, e.g., a picture of a cat).

  Parameter estimation in the presence of such correlated outputs is an active area of research in statistics (see [DDP19; LLZ19] and references therein). Applying the corresponding techniques (or modifications thereof) to datamodels may help calibrate predictions and improve sample-efficiency.

- **Confidence intervals for datamodels.** In this work we have focused on attaining point estimates for datamodel parameters via simple linear regression. A natural extension to these results would be to obtain *confidence intervals* around the datamodel weights. These could, for example, (a) provide interval estimates for model outputs rather than simple point estimates; and (b) decide if a training input is indeed a "significant" predictor for a given test input.

- **Post-selection inference.** Relatedly, the high input-dimensionality of our estimation problem and the sparse nature of the solutions suggests that a *two-stage* procedure might improve sample efficiency. In such procedures, one first selects (often automatically, e.g., via LASSO) a subset of the coefficients deemed to be "significant" for a given test example, then re-fits a linear model for *only* these coefficients. This two-stage approach is particularly attractive in settings where the number of subset-output pairs $(S_i, f_{\mathcal{A}}(x; S_i))$ is less than the size of the training set $|S|$ being subsampled.

  Unfortunately, using the data itself to perform model selection in this manner—a paradigm known as *post-selection inference*—violates the assumptions of classical statistical inference (in particular, that the model class is chosen independently of the data) and can result in significantly miscalibrated confidence intervals. Applying *valid* two-stage estimation to datamodeling would be an area for further improvement upon the protocol presented in our work.

- **Improving subset sampling.** Recall (cf. Section 2) that our framework uses a distribution over subsets $\mathcal{D}_S$ to generate the "datamodel training set." In this paper, we fixed $\mathcal{D}_S$ to be random $\alpha$-subsets of the training set, and used a nearest-neighbors example (see Figure 15) to provide intuition around the role of $\alpha$. While this design choice did yield useful datamodels, it is unclear whether this class of distributions is optimal. In particular, a long line of literature in causal inference focuses on *intervention design* [ES07]; drawing upon this line of work may lead to a better choice of subsampling distribution. Furthermore, one might even go beyond a fixed distribution $\mathcal{D}_S$ and instead choose subsets $S_i$ *adaptively* (i.e., based on the datamodels estimated with the previously sampled subsets) in order to reduce sample complexity.

- **Devising better priors.** Finally, in this paper we employed simple least-squares regression with $\ell_1$ regularization (tuned through a held-out validation set). While the advantage of this rather simple prior—namely, that datamodels are *sparse*—is that the resulting estimation methodology is largely data-driven, one may consider incorporating domain-specific knowledge to design better priors. For instance, one can use structured-sparsity [HZM11] to take advantage of any additional structure.

## 7.2 Studying generalization

Datamodels also present an opportunity to study generalization more broadly:

- **Understanding linearity.** The key simplifying assumption behind our instantiation of the datamodeling framework is that we can approximate the final output of training a model on a subset of the trainset as a *linear* function of the presence of each training point. While this assumption certainly leads to a simple estimation procedure, we have very little justification for why such a linear model should be able to capture the complexities of end-to-end model training on data subsets. However, we find that datamodels *can* accurately predict ground-truth model outputs (cf. Sections 2). In fact,

we find a tight *linear* correlation between datamodel predictions and model outputs even on out-of-distribution (i.e., not in the support of $\mathcal{D}_S$) counterfactual datasets. Understanding *why* a simple linearity assumption leads to effective datamodels for deep neural networks is an interesting open question. Tackling this question may necessitate a better understanding of the training dynamics and implicit biases behind overparameterized training [SRK+20; BMR21].

- **Using sparsity to study generalization.** A recent line of work in machine learning studies the interplay between learning, overparameterization, and memorization [ZBH+16; Cha18; Fel19; BN20; ZBH+20]. Datamodeling may be a helpful tool in this pursuit, as it connects predictions of machine learning models directly to the data used to train them. For example, the *data support* introduced in Section 4.1.1 provides a quantitative measure of "how memorized" a given test input is.

- **Theoretical characterization of the role of $\alpha$.** In line with our intuitions in Section 5, we have observed both qualitatively (e.g., Figure 10) and quantitatively (e.g., Appendix B) that estimating datamodels using different values of $\alpha$ identifies correlations at varying granularities. However, despite empirical results around the clear role of $\alpha$—Appendix B even isolates its effect on datamodels for simple underdetermined linear regression—we lack a crisp *theoretical* understanding of how $\alpha$ affects our estimated datamodels. A better theoretical understanding of the role of $\alpha$, even for simple models trained on structured distributions, can provide us with more rigorous intuition for the phenomena observed here, and can in turn guide the development of better choices of sampling distribution for datamodeling.

## 7.3 Applying datamodels

Finally, each of the presented perspectives in Section 4 can be taken further to enable even better data and model understanding. For example:

- **Interpreting predictions.** For a given test example, the training images corresponding to the largest-magnitude datamodel weights both (a) share features in common with the test example; and (b) seem to be causally linked to the test example (in the sense that removing the training images flips the test prediction). This immediately suggests the potential utility of datamodels as a tool for *interpreting* test-time predictions in a counterfactual-centric manner. Establishing them as such requires further evaluation through, for example, human-in-the-loop studies.

- **Building data exploration tools.** In a similar vein, another opportunity for future work is in building user-friendly *data exploration* tools that leverage datamodel embeddings. In this paper we present the simplest such example in the form of PCA, but leave the vast field of data bias and feature discovery methods (cf. [CAS+19] and Leclerc et al. [LSI+21] for a survey) unexplored.

# 8 Conclusion

We present datamodeling, a framework for viewing the output of model training as a simple function of the presence of each training data point. We show that a simple linear instantiation of datamodeling enables us to predict model outputs accurately, and facilitates a variety of applications.

# Acknowledgements

# References

[ACG+16]    Martín Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. "Deep Learning with Differential Privacy". In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. Vienna, Austria: ACM, 2016, pp. 308–318.

[ARS+15]    Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson. "Factors of transferability for a generic convnet representation". In: *IEEE transactions on pattern analysis and machine intelligence* (2015).

[BPS19]     Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. "Differential privacy has disparate impact on model accuracy". In: *Neural Information Processing Systems (NeurIPS)*. 2019.

[BNB21]     Yamini Bansal, Preetum Nakkiran, and Boaz Barak. "Revisiting Model Stitching to Compare Neural Representations". In: *Neural Information Processing Systems (NeurIPS)*. 2021.

[BMR21]     Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. "Deep learning: a statistical viewpoint". In: *arXiv preprint arXiv:2103.09177*. 2021.

[BD20]      Björn Barz and Joachim Denzler. "Do we train on test data? purging cifar of near-duplicates". In: *Journal of Imaging*. 2020.

[BPF21]     Samyadeep Basu, Phillip Pope, and Soheil Feizi. "Influence Functions in Deep Learning Are Fragile". In: *International Conference on Learning Representations (ICLR)*. 2021.

[BBC+07]    Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. "Analysis of representations for domain adaptation". In: *Neural Information Processing Systems (NeurIPS)*. 2007.

[BCV13]     Y. Bengio, A. Courville, and P. Vincent. "Representation Learning: A Review and New Perspectives". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2013.

[BN20]      Guy Bresler and Dheeraj Nagaraj. "A corrective view of neural networks: Representation, memorization and learning". In: *Conference on Learning Theory (COLT)*. 2020.

[BGM21]     Tamara Broderick, Ryan Giordano, and Rachael Meager. "An Automatic Finite-Sample Robustness Metric: Can Dropping a Little Data Change Conclusions?" In: *Arxiv preprint arXiv:2011.14999*. 2021.

[BBF+21]    Gavin Brown, Mark Bun, Vitaly Feldman, Adam Smith, and Kunal Talwar. "When is memorization of irrelevant training data necessary for high-accuracy learning?" In: *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*. 2021.

[CLK+19]    Nicholas Carlini, Chang Liu, Jernej Kos, Úlfar Erlingsson, and Dawn Song. "The Secret Sharer: Measuring Unintended Neural Network Memorization & Extracting Secrets". In: *USENIX Security Symposium*. 2019.

[CTW+21]    Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. "Extracting training data from large language models". In: *30th USENIX Security Symposium (USENIX Security 21)*. 2021.

[CAS+19]    Shan Carter, Zan Armstrong, Ludwig Schubert, Ian Johnson, and Chris Olah. "Activation atlas". In: *Distill* (2019).

[CGF+19]    Guillaume Charpiat, Nicolas Girard, Loris Felardos, and Yuliya Tarabalka. "Input similarity from the neural network perspective". In: *Neural Information Processing Systems (NeurIPS)*. 2019.

[Cha18]     Satrajit Chatterjee. "Learning and Memorization". In: *Proceedings of the 35th International Conference on Machine Learning*. 2018.

[CFW+18]    Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. "Functional Map of the World". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2018.

[CYM+20]   Cody Coleman, Christopher Yeh, Stephen Mussmann, Baharan Mirzasoleiman, Peter Bailis, Percy Liang, Jure Leskovec, and Matei Zaharia. "Selection via proxy: Efficient data selection for deep learning". In: *International Conference on Learning Representations (ICLR)*. 2020.

[DHM+20]   Alexander D'Amour, Katherine A. Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D. Hoffman, Farhad Hormozdiari, Neil Houlsby, Shaobo Hou, Ghassen Jerfel, Alan Karthikesalingam, Mario Lucic, Yi-An Ma, Cory Y. McLean, Diana Mincu, Akinori Mitani, Andrea Montanari, Zachary Nado, Vivek Natarajan, Christopher Nielson, Thomas F. Osborne, Rajiv Raman, Kim Ramasamy, Rory Sayres, Jessica Schrouff, Martin Seneviratne, Shannon Sequeira, Harini Suresh, Victor Veitch, Max Vladymyrov, Xuezhi Wang, Kellie Webster, Steve Yadlowsky, Taedong Yun, Xiaohua Zhai, and D. Sculley. "Underspecification Presents Challenges for Credibility in Modern Machine Learning". In: *Arxiv preprint arXiv:2011.03395*. 2020.

[DDP19]   Constantinos Daskalakis, Nishanth Dikkala, and Ioannis Panageas. "Regression from dependent observations". In: *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*. 2019, pp. 881–889.

[DCL+19]   Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. "Bert: Pre-training of deep bidirectional transformers for language understanding". In: (2019).

[DT17]   Terrance DeVries and Graham W Taylor. "Improved Regularization of Convolutional Neural Networks with Cutout". In: *arXiv preprint arXiv:1708.04552*. 2017.

[DKM+06]   Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. "Our data, ourselves: Privacy via distributed noise generation". In: *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. 2006.

[ES07]   Frederick Eberhardt and Richard Scheines. "Interventions and Causal Inference". In: *Philosophy of Science*. 2007.

[EIS+19]   Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Brandon Tran, and Aleksander Madry. "Adversarial Robustness as a Prior for Learned Representations". In: *ArXiv preprint arXiv:1906.00945*. 2019.

[Fel19]   Vitaly Feldman. "Does Learning Require Memorization? A Short Tale about a Long Tail". In: *Symposium on Theory of Computing (STOC)*. 2019.

[FZ20]   Vitaly Feldman and Chiyuan Zhang. "What Neural Networks Memorize and Why: Discovering the Long Tail via Influence Estimation". In: *Advances in Neural Information Processing Systems (NeurIPS)*. Vol. 33. 2020, pp. 2881–2891.

[FHT10]   Jerome Friedman, Trevor Hastie, and Rob Tibshirani. "Regularization paths for generalized linear models via coordinate descent". In: *Journal of statistical software* (2010).

[GGS19]   Nidham Gazagnadou, Robert M Gower, and Joseph Salmon. "Optimal mini-batch and step sizes for SAGA". In: *International Conference on Machine Learning (ICML)*. 2019.

[GZ19]   Amirata Ghorbani and James Zou. "Data shapley: Equitable valuation of data for machine learning". In: *International Conference on Machine Learning (ICML)*. 2019.

[GDG17]   Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. "Badnets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain". In: *arXiv preprint arXiv:1708.06733* (2017).

[GGT+17]   Joris Guérin, Olivier Gibaru, Stéphane Thiery, and Eric Nyiri. "CNN features are also great at unsupervised classification". In: *Arxiv preprint arXiv:1707.01700*. 2017.

[GE03]   Isabelle Guyon and André Elisseeff. "An introduction to variable and feature selection". In: *Journal of Machine Learning Research (JMLR)*. 2003.

[HRR+11]   Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust statistics: the approach based on influence functions*. Vol. 196. John Wiley & Sons, 2011.

[HYH+21]   Kazuaki Hanawa, Sho Yokoi, Satoshi Hara, and Kentaro Inui. "Evaluation of similarity-based explanations". In: *International Conference on Learning Representations (ICLR)*. 2021.

[HZR+16]  Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep Residual Learning for Image Recognition". In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016.

[Hoo21]  Sara Hooker. "Moving beyond "algorithmic bias is a data problem"". In: *Patterns*. 2021.

[HCD+19]  Sara Hooker, Aaron Courville, Yann Dauphin, and Andrea Frome. "Selective Brain Damage: Measuring the Disparate Impact of Model Pruning". In: *arXiv preprint arXiv:1911.05248*. 2019.

[HZM11]  Junzhou Huang, Tong Zhang, and Dimitris Metaxas. "Learning with Structured Sparsity." In: *Journal of Machine Learning Research (JMLR)*. 2011.

[IST+19]  Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. "Adversarial Examples Are Not Bugs, They Are Features". In: *Neural Information Processing Systems (NeurIPS)*. 2019.

[JSW+22]  Saachi Jain, Hadi Salman, Eric Wong, Pengchuan Zhang, Vibhav Vineet, Sai Vemprala, and Aleksander Madry. "Missingness Bias in Model Debugging". In: *International Conference on Learning Representations*. 2022.

[JTM21]  Saachi Jain, Dimitris Tsipras, and Aleksander Madry. "Co-Priors: Combining Biases on Learned Features". In: *Preprint*. 2021.

[JDW+19]  Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel, Bo Li, Ce Zhang, Dawn Song, and Costas J. Spanos. "Towards Efficient Data Valuation Based on the Shapley Value". In: *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*. 2019.

[JNB+21]  Yiding Jiang, Vaishnavh Nagarajan, Christina Baek, and J. Zico Kolter. "Assessing Generalization of SGD via Disagreement". In: *Arxiv preprint arXiv:2106.13799*. 2021.

[KAT+19]  Pang Wei Koh, Kai-Siang Ang, Hubert HK Teo, and Percy Liang. "On the accuracy of influence functions for measuring group effects". In: *Neural Information Processing Systems (NeurIPS)*. 2019.

[KL17]  Pang Wei Koh and Percy Liang. "Understanding Black-box Predictions via Influence Functions". In: *International Conference on Machine Learning*. 2017.

[KSM+20]  Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Sara Beery, et al. "WILDS: A Benchmark of in-the-Wild Distribution Shifts". In: *arXiv preprint arXiv:2012.07421* (2020).

[KNL+19]  Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. "Similarity of Neural Network Representations Revisited". In: *Proceedings of the 36th International Conference on Machine Learning (ICML)*. 2019.

[Kri09]  Alex Krizhevsky. "Learning Multiple Layers of Features from Tiny Images". In: *Technical report*. 2009.

[LPB17]  Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. "Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles". In: *Neural Information Processing Systems (NeurIPS)*. 2017.

[LIE+22]  Guillaume Leclerc, Andrew Ilyas, Logan Engstrom, Sung Min Park, Hadi Salman, and Aleksander Madry. *ffcv*. https://github.com/libffcv/ffcv/. 2022.

[LSI+21]  Guillaume Leclerc, Hadi Salman, Andrew Ilyas, Sai Vemprala, Logan Engstrom, Vibhav Vineet, Kai Xiao, Pengchuan Zhang, Shibani Santurkar, Greg Yang, et al. "3DB: A Framework for Debugging Computer Vision Models". In: *arXiv preprint arXiv:2106.03805*. 2021.

[LIN+21]  Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. "Deduplicating Training Data Makes Language Models Better". In: *Arxiv preprint arXiv:2107.06499*. 2021.

[LC94]  David D Lewis and Jason Catlett. "Heterogeneous uncertainty sampling for supervised learning". In: *Machine learning proceedings 1994*. 1994, pp. 148–156.

[LLZ19]    Tianxi Li, Elizaveta Levina, and Ji Zhu. "Prediction models for network-linked data". In: *The Annals of Applied Statistics*. 2019.

[LL17]     Scott Lundberg and Su-In Lee. "A unified approach to interpreting model predictions". In: *Neural Information Processing Systems (NeurIPS)*. 2017.

[MMS+19]   Horia Mania, John Miller, Ludwig Schmidt, Moritz Hardt, and Benjamin Recht. "Model similarity mitigates test set overuse". In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2019, pp. 9993–10002.

[MT20]     PG Martinsson and JA Tropp. "Randomized numerical linear algebra: foundations & algorithms". In: *arXiv preprint arXiv:2002.01387*. 2020.

[MGS18]    Mathurin Massias, Alexandre Gramfort, and Joseph Salmon. "Celer: a Fast Solver for the Lasso with Dual Extrapolation". In: *Proceedings of the 35th International Conference on Machine Learning (ICML)*. 2018.

[NB20]     Preetum Nakkiran and Yamini Bansal. "Distributional generalization: A new kind of generalization". In: *Arxiv preprint arXiv:2009.08092*. 2020.

[Nea96]    Radford Neal. *Bayesian Learning for Neural Networks*. Springer, 1996.

[OMS17]    Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. "Feature Visualization". In: *Distill*. 2017.

[Owe72]    Guillermo Owen. "Multilinear Extensions of Games". In: *Management Science*. 1972.

[Pag18]    David Page. *CIFAR-10 Fast*. GitHub Repository. Oct. 2018. URL: https://github.com/davidcpage/cifar10-fast.

[PVG+11]   F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. "Scikit-learn: Machine Learning in Python". In: *Journal of Machine Learning Research*. Vol. 12. 2011, pp. 2825–2830.

[PJW+21]   Pouya Pezeshkpour, Sarthak Jain, Byron C Wallace, and Sameer Singh. "An Empirical Comparison of Instance Attribution Methods for NLP". In: *North American Chapter of the Association for Computational Linguistics (NAACL)*. 2021.

[PLS+20]   Garima Pruthi, Frederick Liu, Mukund Sundararajan, and Satyen Kale. "Estimating Training Data Influence by Tracing Gradient Descent". In: *Neural Information Processing Systems (NeurIPS)*. 2020.

[RSG16]    Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. ""Why Should I Trust You?": Explaining the Predictions of Any Classifier". In: *International Conference on Knowledge Discovery and Data Mining (KDD)*. 2016.

[RWD88]    Tim Robertson, F.T. Wright, and R. L. Dykstra. *Order Restricted Statistical Inference*. Wiley Series in Probability and Statistics, 1988.

[RWR+20]   Elan Rosenfeld, Ezra Winston, Pradeep Ravikumar, and Zico Kolter. "Certified robustness to label-flipping attacks via randomized smoothing". In: *International Conference on Machine Learning (ICML)*. 2020.

[SWM+89]   Jerome Sacks, William J. Welch, Toby J. Mitchell, and Henry P. Wynn. "Design and Analysis of Computer Experiments". In: *Statistical Science*. Vol. 4. 4. Institute of Mathematical Statistics, 1989, pp. 409–423. URL: http://www.jstor.org/stable/2245858.

[SRK+20]   Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. "An investigation of why overparameterization exacerbates spurious correlations". In: *International Conference on Machine Learning*. PMLR. 2020, pp. 8346–8356.

[SYW+21]   Thibault Sellam, Steve Yadlowsky, Jason Wei, Naomi Saphra, Alexander D'Amour, Tal Linzen, Jasmijn Bastings, Iulia Turc, Jacob Eisenstein, Dipanjan Das, Ian Tenney, and Ellie Pavlick. "The MultiBERTs: BERT Reproductions for Robustness Analysis". In: *Arxiv preprint arXiv:2106.16163*. 2021.

[Sha51]    LS Shapley. "Notes on the n-Person Game—II: The Value of an n-Person Game, The RAND Corporation, The RAND Corporation". In: *Research Memorandum*. 1951.

[SC04]     John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.

[SHS+19]   Kacper Sokol, Alexander Hepburn, Raul Santos-Rodriguez, and Peter Flach. "bLIMEy: Surrogate Prediction Explanations Beyond LIME". In: *Arxiv preprint arXiv:1910.13016*. 2019.

[Spe04]    Charles Spearman. "The Proof and Measurement of Association between Two Things". In: *The American Journal of Psychology*. 1904.

[SLL20]    Pascal Sturmfels, Scott Lundberg, and Su-In Lee. "Visualizing the Impact of Feature Attribution Baselines". In: *Distill* (2020). https://distill.pub/2020/attribution-baselines. DOI: 10.23915/distill.00022.

[TSC+19]   Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. "An Empirical Study of Example Forgetting during Deep Neural Network Learning". In: *ICLR*. 2019.

[WZJ+21]   Tianhao Wang, Yi Zeng, Ming Jin, and Ruoxi Jia. "A Unified Framework for Task-Driven Data Quality Management". In: *ArXiv preprint arXiv:2106.05484*. 2021.

[WSM21]    Eric Wong, Shibani Santurkar, and Aleksander Madry. "Leveraging Sparse Linear Layers for Debuggable Deep Networks". In: *International Conference on Machine Learning (ICML)*. 2021.

[XXE12]    Han Xiao, Huang Xiao, and Claudia Eckert. "Adversarial Label Flips Attack on Support Vector Machines." In: *European Conference on Artificial Intelligence (ECAI)*. 2012.

[YKY+18]   Chih-Kuan Yeh, Joon Sik Kim, Ian E. H. Yen, and Pradeep Ravikumar. "Representer Point Selection for Explaining Deep Neural Networks". In: *Neural Information Processing Systems (NeurIPS)*. 2018.

[ZBH+20]   Chiyuan Zhang, Samy Bengio, Moritz Hardt, Michael C Mozer, and Yoram Singer. "Identity crisis: Memorization and generalization under extreme overparameterization". In: *International Conference on Learning Representations (ICLR)*. 2020.

[ZBH+16]   Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. "Understanding deep learning requires rethinking generalization". In: *International Conference on Learning Representations (ICLR)*. 2016.

[ZIE+18]   Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. "The unreasonable effectiveness of deep features as a perceptual metric". In: *Computer Vision and Pattern Recognition (CVPR)*. 2018.

[ZGK+21]   Ruiqi Zhong, Dhruba Ghosh, Dan Klein, and Jacob Steinhardt. "Are Larger Pretrained Language Models Uniformly Better? Comparing Performance at the Instance Level". In: *Findings of the Association for Computational Linguistics (Findings of ACL)*. 2021.

# Appendices

# A   Pseudocode for Estimating Datamodels

---

**Algorithm A.1** An outline of the datamodeling framework: we use a simple parametric model as a proxy for the entire end-to-end training process.

---

1: **procedure** ESTIMATEDATAMODEL(target example $x$, trainset $S$ of size $d$, subsampling frac. $\alpha \in (0,1)$)
2: $\quad$ $T \leftarrow []$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Initialize *datamodel training set*
3: $\quad$ **for** $i \in \{1, \ldots, m\}$ **do**
4: $\qquad$ Sample a subset $S_i \subset S$ from $\mathcal{D}_S$ where $|S_i| = \alpha \cdot d$
5: $\qquad$ $y_i \leftarrow f_{\mathcal{A}}(x; S_i)$ $\qquad\qquad\qquad\qquad$ ▷ Train a model on $S_i$ using $\mathcal{A}$, evaluate on $x$
6: $\qquad$ Define $\mathbf{1}_{S_i} \in \{0,1\}^d$ as $(\mathbf{1}_{S_i})_j = 1$ if $x_j \in S_i$ else 0
7: $\qquad$ $T \leftarrow T + [(\mathbf{1}_{S_i}, y_i)]$ $\qquad\qquad\qquad\qquad\qquad$ ▷ Update datamodel training set
8: $\quad$ $\theta \leftarrow$ RUNREGRESSION(T) $\qquad\qquad\qquad$ ▷ Predict the $y_i$ from the $\mathbf{1}_{S_i}$ vectors
9: $\quad$ **return** $\theta$ $\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Result: a weight vector $\theta \in \mathbb{R}^d$

---

**Algorithm A.2** Assessing datamodels' ability to predict dataset counterfactuals.

---

1: **procedure** COUNTERFACTUALEVAL(target example $x$, datamodel $\theta$, trainset of size $d$)
2: $\quad$ $M_0 \leftarrow$ AVERAGE$[f_{\mathcal{A}}(x; [d])$ **for** $i \leftarrow 1 \ldots 100]$ $\qquad$ ▷ Average margin on $x$ with full training set
3: $\quad$ **for** $k \in \{20, 40, 80, 160, 320\}$ **do**
4: $\qquad$ $G \leftarrow$ TOP-K$(\theta, k)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Get the top-$k$ indices of $\theta$
5: $\qquad$ $M \leftarrow$ AVERAGE$[f_{\mathcal{A}}(x; [d] \setminus G)$ **for** $i \leftarrow 1 \ldots 20]$ $\qquad$ ▷ Average margin without $G$ over 20 trials
6: $\qquad$ $\Delta_{avg} \leftarrow M_0 - M$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Actual counterfactual
7: $\qquad$ $\widehat{\Delta} \leftarrow \theta^\top \mathbf{1}_G \ (= \sum_{i \in G} \theta_i)$ $\qquad\qquad\qquad$ ▷ Datamodel-predicted counterfactual
8: $\quad$ Compare all values of $\widehat{\Delta}$ and $\Delta_{avg}$

---

# B Understanding the Role of $\alpha$ through Simulation

At a high level, our intuition for the subsampling fraction[10] $\alpha$ is that datamodels estimated with higher $\alpha$ tend to detect more *local* effects (i.e., those driven by smaller groups of examples, such as near-duplicates or small subpopulations), while those estimated with lower $\alpha$ detect more *global* effects (i.e., those driven by larger groups of images, such as large subpopulations or subclass biases). To solidify and corroborate this intuition about $\alpha$, we analyze a basic simulated setting.

**Setup.** We consider an underdetermined linear regression model operating on $n$ data points with $d$ binary features, i.e., $x_i \in \{0,1\}^d, y_i \in \mathbb{R}$. Let $X \in \mathbb{R}^{n \times d}$ and $y \in \mathbb{R}^n$ denote their matrix and vector counterparts. $S$ is training set consisting of these $n$ samples, and we use an equally sized held-out set $S_V$ for evaluation.

The feature coordinates are distributed as Bernoulli variables of varying frequency:

$$x_{ik} \sim \text{Bernoulli}(p_k) \text{ for } 1 \le i \le n, \tag{16}$$

$$\text{where } p_k \in \left\{ \frac{1}{10}, \frac{2}{10}, \frac{3}{10}, \frac{4}{10}, \frac{5}{10} \right\} \text{ for } 1 \le k \le d.$$

Each feature $k \in [d]$ naturally defines a *subpopulation $S_k$*, the group of training examples with feature $k$ active, i.e., $S_k := \{x_i \in S : x_{ik} = 1\}$. Features with *lower* (resp. *higher*) frequency $p_k$ are intended to capture more *local* (resp. more *global*) effects.

The observed labels are generated according to a linear model $y := Xw + \mathcal{N}(0, \epsilon)$, where $w$ is the true parameter vector and $\epsilon > 0$ is a constant. We generate samples with $d = 150, n = 125$ and use linear regression[11] to estimate $w$.

Now, to use datamodels to analyze the above "training process" of fitting a linear regression model, we will model the output function $f_\mathcal{A}(;)$ given by the prediction of the linear model at point $x_j$ when $w$ is estimated with samples $S \subset S$, e.g.

$$f_\mathcal{A}(x_j; S) = (X_S^\top (X_S X_S^\top)^{-1} y_S) \cdot x_j \tag{17}$$

We generate $m = 1,000,000$ subsampled training subsets[12] along with their evaluations, and use ordinary least squares (OLS) to fit the datamodels. (Note that the use of OLS here is separate from the use of linear regression above as the original model class.)

**Analysis.** Our hypothesis is that datamodels estimated with lower (resp. higher) $\alpha$ are better at detecting the effect of features of higher (resp. lower) frequency. To test this, we estimate datamodels for the entire test set (stacking them into a matrix $\Theta \in \mathbb{R}^{n \times n}$, where $\Theta_{.j}$ is the datamodel for $x_j$) . We do this for varying values of $\alpha \in (0,1)$, and evaluate how well each set of datamodels predicts the effects of features across different frequencies. First, to evaluate a datamodel on some feature $k$, we can compare the following two quantities for different test examples $x_j \in S_V$:

(a) The *actual* effect of removing the subpopulation $S_k$ on $x_j$, i.e., $f_\mathcal{A}(x_j; S) - f_\mathcal{A}(x_j; S \setminus S_k)$,

(b) The *datamodel-predicted* effect of removing $S_k$, i.e., $\sum_{x_i \in S} \Theta_{ij} \cdot \mathbf{1}\{x_i \in S_k\}$.

To quantify the predictiveness of the datamodel at frequency $p$, we compute the *Pearson correlation* between the above two quantities over all features $k$ with frequency $p$ and all test examples; see Algorithm B.1 for a pseudocode. We repeat this evaluation varying $p$ and the datamodel (varying $\alpha$). According to our intuition, for features $k$ with lower (resp. higher) frequency $p_k$, this correlation should be maximized at higher (resp. lower) values of $\alpha$, where the datamodels capture more local (resp. global) effects. Figure B.1 accurately reflects this intuition: more local (i.e., less frequent) features are best detected at higher $\alpha$.

---

[10]See Section 2 for definition.

[11]As the system is underdetermined, we use the pseudoinverse of $X$ to find the solution with the smallest norm.

[12]Large sample size make sampling error negligible.

**Algorithm B.1** Evaluating datamodel's counterfactual predictiveness for features at a particular frequency.

1: **function** FEATURECORRELATION(datamodel matrix $\Theta$, feature frequency $p$)
2:     **for** $j \leftarrow 1 \ldots d$ if $p_j = p$ **do**                      $\triangleright$ For all features with frequency $p$
3:         **for** $x_j \in S_V$ **do**                                 $\triangleright$ For all test examples
4:             $\Delta \leftarrow f_{\mathcal{A}}(x_j; S) - f_{\mathcal{A}}(x_j; S \setminus S_k)$     $\triangleright$ Actual counterfactual where $f_{\mathcal{A}}(;)$ is given by (17)
5:             $\widehat{\Delta} \leftarrow \sum_{x_i \in S} \Theta_{ij} \cdot \mathbf{1}\{x_i \in S_k\}$               $\triangleright$ Datamodel-predicted counterfactual
6: **return** Pearson correlation between $\widehat{\Delta}$ and $\Delta_{avg}$            $\triangleright$ Across all features and test examples



**Figure B.1: The role of the subsampling fraction $\alpha$ in a simulated linear model.** The data consists of $d$-dimensional binary vectors $x_i$, which comprise (overlapping) *subpopulations* $S_k$ defined by a shared feature $k$, and their corresponding labels are generated according to a linear model, i.e. $y := Xw + \mathcal{N}(0, \epsilon)$. We estimate datamodels using various $\alpha$, and measure their ability to detect features at different frequencies $p$. To quantify latter, we compute the *Pearson correlation* between i) the *actual effect* of removing the subpopulation $S_k$ on a test example and ii) the *datamodel-predicted* effect, across all features with frequency $p$. Each line in the above plot corresponds to features of a particular frequency $p$, and shows the correlation ($y$-axis) while varying the datamodel ($\alpha$, $x$-axis). Consistent with our intuition, we observe that higher (resp. lower) values of $\alpha$ are better at detecting less (resp. more) frequent features, i.e. more local (resp. global) effects.

34

# C  Selecting Output Function to Model

In this section, we outline a heuristic method for selecting the output function $f_{\mathcal{A}}(x; S)$ to model. The heuristic is neither sufficient *nor* necessary for least-squares regression to work, but may provide some signal as to which output may yield better datamodels.

The first problem we would like to avoid is "output saturation," i.e., being unable to learn a good datamodel due to insufficient variation in the output. This effect is most pronounced when we measure model correctness: indeed, over 30% of the CIFAR-10 test set is either always correct or always incorrect over all models trained, making datamodel estimation impossible. However, this issue is not unique to correctness. We propose a very simple test inspired by the idealized ordinary least squares model to measure how normally distributed a given type of model output is.

**Normally distributed residuals.**  In the idealized ordinary least squares model, the observed outputs $f_{\mathcal{A}}(x; S)$ would follow a normal distribution with fixed mean $(\theta^\star)^\top \mathbf{1}_S$ and unknown variance, where $\theta^\star$ is the true parameter vector. Although we cannot guarantee this condition, we can measure the "normality" of the outputs (again, for a single fixed subset), with the intuition that the more normal the observed outputs are, the better a least-squares regression will work. Hence, compare different output functions by estimating the noise distribution of datamodels given each choice of output function. We leverage our ability—in contrast to typical settings for regression— to sample multiple response variables $f_{\mathcal{A}}(x; S)$ for a fixed $S$ (by retraining several models on the same data and recording the output on a fixed test example).

In Figure C.1, we show the results of normality test for residuals arising from different choices of $f_{\mathcal{A}}(\cdot; S)$: correctness function, confidence on the correct class, cross-entropy loss, and finally correct-class margin[13]. Correct-class margins is the only choice of $f_{\mathcal{A}}(\cdot; S)$ where the $p$-values are distributed nearly uniformly, which is consistent with the outputs being normally distributed. Hence, we choose to use the correct-class margins as the dependent variable for fitting our datamodels.



**Figure C.1: Correct-class margins are close to normally distributed.** For each choice of output function $f_{\mathcal{A}}$: we (a) fix a random subset $S' \sim \mathcal{D}_S$, where $S$ is the CIFAR-10 train set, (b) train 200 models on $S'$ and evaluate them on the entire CIFAR-10 test set, and (c) for each image $x_i$ in the test set, calculate a $p$-value for rejecting the normality of $f_{\mathcal{A}}(x_i; S')$. We plot a histogram of these $p$-values above. For every output function other than correct-class margin, almost every test is rejected, whereas for margins the distribution of $p$-values is uniform across $[0, 1]$, which is consistent with the null hypothesis (normality).

---

[13]Correct-class margin is the difference between the correct-class logit and the highest incorrect-class logit; it is unbounded by definition, and its sign indicates the correctness of the classification.

# D   Experimental Setup

## D.1   Datasets

**CIFAR-10.**   We use the standard CIFAR-10 dataset [Kri09].

**FMoW.**   FMoW [CFW+18] is a land use classification dataset based on satellite imagery. WILDS [KSM+20] uses a subset of FMoW and repurposes it as a benchmark for out-of-distribution (OOD) generalization; we use same the variant (presized to 224x224, single RGB image per example rather than a time sequence). We perform our analysis only on the in-distribution train/test splits (e.g. overlapping years) as our focus is not on OOD settings. Also, we limit our data to the year 2012. (These restrictions are only for convenience, and our framework can easily extend and scale to more general settings.)

Properties of both datasets are summarized in Table D.1.

Table D.1: Properties of datasets used.

| Dataset | Classes | Size (Train/Test) | Input Dimensions |
|---|---|---|---|
| CIFAR-10 | 10 | 50,000/10,000 | $3 \times 32 \times 32$ |
| FMoW | 62 | 21,404/3,138 | $3 \times 224 \times 224$ |

## D.2   Models and hyperparameters

**CIFAR-10.**   We use a ResNet-9 variant from Kakao Brain[14] optimized for fast training. The hyperparameters (Table D.2) were chosen using a grid search. We use the standard batch SGD. For data augmentation, we use random 4px random crop with reflection padding, random horizontal flip, and $8 \times 8$ CutOut [DT17].

For counterfactual experiments with ResNet-18 (Figure F.6), we use the standard variant [HZR+16].

**FMoW.**   We use the standard ResNet-18 architecture [HZR+16]. The hyperparameters (Table D.2) were chosen using a grid search, including over different optimizers (SGD, Adam) and learning rate schedules (step decay, cyclic, reduce on plateau). As in Koh et al. [KSM+20], we do not use any data augmentation. Unlike prior work, we do not initialize from a pre-trained ImageNet model; while this results in lower accuracy, this allows us to focus on the role of the FMoW dataset in isolation.

Table D.2: Hyperparameters for used model class.

| Dataset | Initial LR | Batch Size | Epochs | Cyclic LR Peak Epoch | Momentum | Weight Decay |
|---|---|---|---|---|---|---|
| CIFAR-10 | 0.5 | 512 | 24 | 5 | 0.9 | 5e-4 |
| FMoW | 0.4 | 512 | 15 | 6 | 0.9 | 1e-3 |

**Performance.**   In Table D.3, we show for each dataset the accuracies of the chosen model class (with its specific hyperparameters), across different values of $\alpha$.

## D.3   Training infrastructure

**Computing resources.**   We train our models on a cluster of machines, each with 9 NVIDIA A100 GPUs and 96 CPU cores. We also use half-precision to increase training speed.

**Data loading.**   We use FFCV [LIE+22], which removes the data loading bottleneck for smaller models and allows us achieve a throughput of over 5,000 CIFAR-10 models a day *per* GPU.

---

[14]https://github.com/wbaek/torchskeleton/blob/master/bin/dawnbench/cifar10.py

Table D.3: Accuracies for our chosen model classes on CIFAR-10 and FMoW across varying $\alpha$.

| | Accuracy (%) | |
| Subset size ($\alpha$) | CIFAR-10 | FMoW |
|---|---|---|
| 1.0 | 93.00 | 33.76 |
| 0.75 | 91.77 | 31.16 |
| 0.5 | 89.61 | 25.97 |
| 0.2 | 81.62 | 14.70 |
| 0.1 | 71.60 | N/A |

**Data processing.** Our datamodel estimation uses (the characteristic vectors) of training subsets and model outputs (margins) on train and test sets. Hence, we do not need to store any model checkpoints, as it suffices to store the training subset and the model outputs after evaluating at the end of training. In particular, training subsets and model outputs can be stored as $m \times n$ or $m \times d$ matrices, with one row for each model instance and one column for each train or test example. All subsequent computations only require the above matrices.

# E  Regression

## E.1  Solver details

As mentioned in Section 2, we construct datamodels by running $\ell_1$-regularized linear regression, predicting correct-class margins from characteristic vectors, or *masks*, $\mathbf{1}_{S_i}$. The resulting optimization problem is rather large: for example, estimating datamodels for $\alpha = 50\%$ requires running LASSO with a covariate matrix $X$ of size $50,000 \times 300,000$, which corresponds to about 60GB of data; for $\alpha = 10\%$, datamodels this increases five folds as there are 1.5 million models. Moreover, we need to solve up to $60,000$ regression problems (one datamodel each train / test example). The large-scale nature of our estimation problem rules out off-the-shelf solutions such as scikit-learn [PVG+11], GLMNet [FHT10], or Celer [MGS18], all of which either runs out of memory or does not terminate within reasonable time.

Note that solving large linear systems efficiently is an area of active research ([MT20]), and as a result we anticipate that datamodel estimation could be significantly improved by applying techniques from numerical optimization. In this paper, however, we take a rather simple approach based on the SAGA algorithm of [GGS19]. Our starting point is the GPU-enabled implementation of Wong et al. [WSM21]—while this implementation terminated (unlike the CPU-based off-the-shelf solutions), the regressions are still prohibitively slow (i.e., on the order of several GPU-hours per single datamodel estimation). To address this, we make the following changes:

**Fast dataloading.**  The first performance bottleneck turns out to be in dataloading. More specifically, SAGA is a minibatch-based algorithm: at each iteration, we have to read $B$ masks (50,000-dimensional binary vectors) and $B$ outputs (scalars) and move them onto the GPU for processing. If the masks are read from disk, I/O speed becomes a major bottleneck—on the other hand, if we pre-load the entire set of masks into memory, then we are not able to run multiple regressions on the same machine, since each regression will use essentially the entire RAM disk. To resolve this issue, we use the FFCV library [LIE+22] for dataloading—FFCV is based on memory mapping, and thus allows for multiple processes to read from the same memory (combining the benefits of the two aforementioned approaches). FFCV also supports batch pre-loading and parallelization of the data processing pipeline out-of-the-box—adapting the SAGA solver to use FFCV cut the runtime significantly.

**Simultaneous outputs.**  Next, we leverage the fact that the SAGA algorithm is trivially parallelizable across different instances (sharing the same input matrix), allowing us to estimate *multiple* datamodels at the same time. In particular, we estimate datamodels for the entire test set in one pass, effectively cutting the runtime of the algorithm by the test set size (e.g., 10,000 for CIFAR-10).

**Optimizations.**  In order to parallelize across test examples, we need to significantly reduce the GPU memory footprint of the SAGA solver. We accomplish this through a combination of simple code optimization (e.g., using in-place operations rather than copies) as well as writing a few custom CUDA kernels to speed up and reduce the memory consumption of algorithms such as soft thresholding or gradient updating.

**Experimental details.**  For each dataset considered, we chose a maximum $\lambda$: 0.01 for CIFAR-10 test, 0.1 for CIFAR-10 trainset, and 0.05 for FMoW datamodels. Next, we chose $k = 100$ logarithmically spaced intermediate values between $(\lambda/100, \lambda)$ as the regularization path. We ran one regression per intermediate $\lambda$, using $m - 50,000$ samples (where $m$ is as in the table in Figure 1 (right)) to estimate the parameters of the model and the remaining $50,000$ samples as a validation set. For each image in the test set, we select the $\lambda$ corresponding to the best-performing predictor (on the heldout set) along the regularization path. We then *re-run* the regression once more using these optimal $\lambda$ values and the full set of $m$ samples.
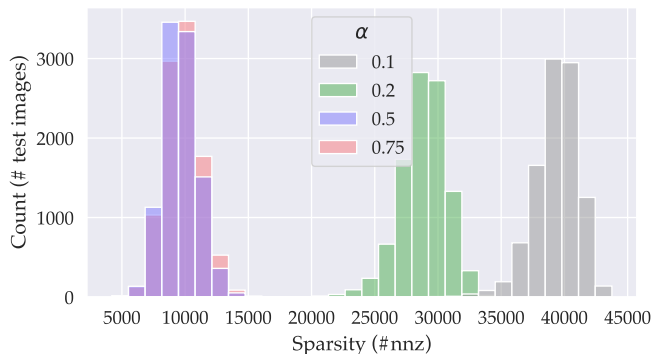
## E.2 Omitted results



**Figure E.1: Sparsity distribution of different datamodels.** Above shows the distribution of datamodel sparsity over test examples on CIFAR-10, compared across different $\alpha$. Sparsity decreases with higher $\alpha$, which is consistent with our intuitions (Section 5) that higher $\alpha$ captures relationships driven by smaller groups of images.
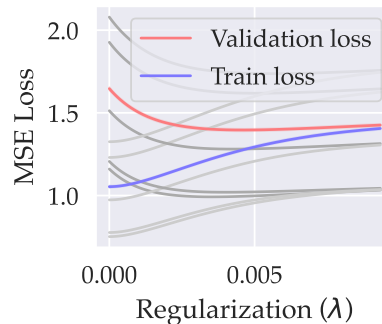
**Figure E.2: The role of regularization.** Average in-sample and out-of-sample MSE (i.e., (10)) for datamodels on CIFAR-10 estimated by optimizing the regularized least-squares objective (8) for varying $\lambda$.



**Figure E.3:** Identical results to Figure 5 for FMoW.



**Figure E.4:** Identical results to Figure 5 for CIFAR-10 for all values of $\alpha$.

# F  Counterfactual Prediction with Datamodels

## F.1  General setup

**Sample selection.**   For all of our counterfactual experiments, we use a random sample of the respective test datasets. We select at random 300 test images for CIFAR-10 (class-balanced; 30 per class) and 100 test images for FMoW. For the CIFAR-10 baselines, we consider counterfactuals for a 100 image subset of the 300.

**Size of counterfactuals.**   For CIFAR-10, we remove top $k = \{10, 20, 40, 80, 160, 320, 640, 1280\}$ images and bottom $k = \{20, 40, 80, 160, 320\}$ where applicable. For FMoW, we remove top and bottom $k = \{10, 20, 40, 80, 160, 320, 640\}$.

**Reducing noise by averaging.**   Each counterfactual (i.e., training models on a given training set $S'$) is evaluated over $T$ trials to reduce the variance that arises purely from non-determinism in model training. We use $T = 20$ for CIFAR-10 and FMoW, and $T = 10$ for CIFAR-10 baselines. In Appendix F.6, we show that using sufficiently high $T$ is important for reducing noise.

**Control values.**   To calculate the actual effects in all of our counterfactual evaluation, we need control values $\mathbb{E}[f_{\mathcal{A}}(x; S)]$ for the "null," i.e, margins when trained on 100% of the data. We estimate this by averaging over models on trained on the full training set (10,000 for CIFAR-10 and 500 for FMoW).

## F.2  Baselines

We describe the baseline methods used to generate data support estimates and counterfactuals. Each of the methods gives a way to select training examples that are most similar or influential to a target example. As in prior work [HYH+21; PJW+21], we consider a representative set of baselines spanning both methods based on representation similarity and gradient-based methods, such as influence functions.

**Representation distance.**   We use $\ell_2$ distances in the penultimate layer's representation to rank the training examples in order of similarity to the target test example. We also evaluated dot product, cosine, and mahalanobis distances, but they did not show much variation in their counterfactual effects.[15]

In order to more fairly compare with datamodels—so that we can disentangle the variance reduction from using many models and the additional signal captured by datamodels—we also averaged up to 1000 models' representation distances[16], but this had no discernible difference on the size of the counterfactual effects.

**Influence functions.**   We apply the influence function approximation introduced in [KL17]. In particular, we use the following first-order approximation for the influence of $z$ on the loss $L$ evaluated at $z_{\text{test}}$:

$$\mathcal{I}_{\text{up,loss}}(z, z_{\text{test}}) = -\nabla_\theta \ell(z_{\text{test}}, \widehat{\theta})^\top H_{\widehat{\theta}}^{-1} \nabla_\theta \ell(z, \widehat{\theta})$$

where $\widehat{\theta}$ is the empirical risk minimizer on the training set and $H$ is the Hessian of the loss. The influence here is just the dot product of gradients, weighted by the Hessian. We approximate these influence values by using the methods in [KL17] and as implemented (independently) in `pytorch-influence-functions`.[17] As in [KL17], we take a pretrained representation (of a ResNet-9 model, same as that modeled by our datamodels), and compute approximate influence functions with respect to only the parameters in the last linear layer.

---

[15]With the exception of dot product, which performs poorly due to lack of normalization; this is consistent with the findings in Hanawa et al. [HYH+21].

[16]We simply average the ranks from each model, but there are potentially better ways to aggregate them.

[17]https://github.com/nimarb/pytorch_influence_functions

**TracIn.** Pruthi et al. [PLS+20] define an alternative notion of influence: the influence of a training example $z$ on a test example $z'$ is the total change in loss on $z'$ contributed by updates from mini-batches containing $z$—intuitively, this measures whether gradient updates from $z$ are helpful to learning example $z'$. They approximate this in practice with TracInCP, which considers checkpoints $\theta_{t_1}, ..., \theta_{t_k}$ across training, and sums the dot product of the gradients at $z$ and $z'$ at each checkpoint:

$$\texttt{TracInCP}(z, z') = \sum_{i=1}^{k} \eta_i \nabla_\theta \ell(z, \theta_{t_i}) \cdot \nabla_\theta \ell(z', \theta_{t_i})$$

One can view TracInCP as a variant of the gradient dot product, but averaged over models at different epochs) and weighted by the learning rate $\eta_i$.

**Random baseline.** We also consider a random baseline of removing examples from the same class.

## F.3 Data support estimation

**Setup.** We use datamodels together with counterfactual evaluations in a guided search to efficiently estimate upper bounds on the size of data supports. For a given target example $x$ with corresponding datamodel $g_\theta$, we want to find candidate training subsets of small size $k$ whose removal most reduces the classification margin on $x$:

$$G_k := \arg \min_{|G|=k} g_\theta(S \setminus G_k). \tag{18}$$

Because $g_\theta$ is a linear model in our case, the solution to the above minimization problem is simply the set corresponding to the largest $k$ coordinates of the datamodel parameter $\theta$:

$$G_k = \arg \max_{G \subset S; |G|=k} \theta^\top \mathbf{1}_G = \text{top-}k \text{ indices of } \theta. \tag{19}$$

Our goal is to the find the smallest of these subsets $\{G_k\}_k$ so that $f_\mathcal{A}(x; S \setminus G_k) < 0$, i.e., the example is misclassified on average as per our definition.[18] Thus, for each target $x$, we try several values of $k \in \{10, 20, 40, 80, 160, 320, 640, 1280\}$, training models on the set $S \setminus G_k$ and evaluating the resulting models on $x$.[19] We train $T = 20$ models on each counterfactual $S \setminus G_k$ to reduce variance.

   (Given that we are using datamodels as surrogates after all, one might wonder if the above counterfactual evaluations are actually necessary—one could instead consider estimating the optimal $k$ directly from $\theta$. We revisit a heuristic estimation procedure based on this idea at the end of this subsection.)

**Estimation methodology.** We assume that the expected margin $h(k) := f_\mathcal{A}(x; S \setminus G_k)$ after removing $k$ examples decreases *monotonically* in $k$; this is expected from the linearity of our datamodels and is further supported empirically (see Figure F.1). Then, our goal is to estimate the unique zero[20] $\hat{k}$ of the above function $h(k)$ based on (noisy) samples of $h(k)$ at our chosen values of $k$. Note that by definition, $\hat{k}$ is an upper bound on SUPPORT$(x)$. Now, because of our monotonicity assumption, we can cast estimating $\hat{k}$ as instance of an isotonic regression problem [RWD88]); this effectively performs piecewise linear interpolation, while ensuring that monotonicity constraint is not violated. We use sklearn's IsotonicRegression to fit an estimate $h(k)$, and use this to estimate $\hat{k}$.

**Verifying support estimates.** Due to stochasticity in evaluating counterfactuals, the estimate $\hat{k}$ is noisy. Thus, it is possible that $\hat{k}$ is not a valid upper bound on SUPPORT$(x)$, e.g. removing top $\hat{k}$ examples do not misclassify $x$. In fact, removing $G_{\hat{k}}$ and re-training shows that only 67% of the images are actually misclassified. To establish an upper bound on SUPPORT$(x)$ that has sufficient coverage, we evaluate the

---

[18]Note that $\mathbb{E} f_\mathcal{A}(x; \cdot) < 0$ does not imply that the probability of misclassification is greater than 50%. Nonetheless, it is a natural threshold.

[19]While a binary search over $k$ for each $x$ would be more sample efficient, we collect the entire grid of samples for simplicity.

[20]More precisely, the upper ceiling as data support is defined as an integer quantity.
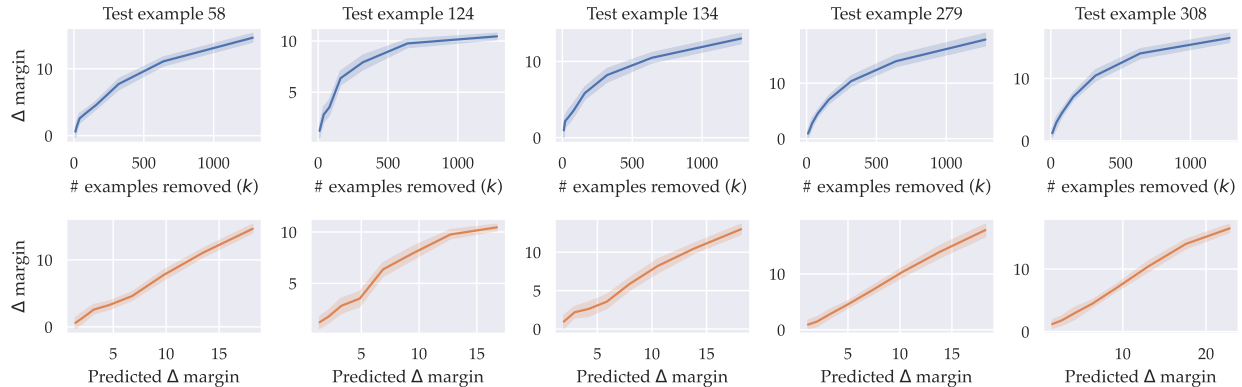
**Figure F.1: Counterfactuals for individual examples.** We plot the results of counterfactual evaluations (using $\alpha = 0.5$ datamodels) for five individual examples, shown in separate columns. **(Top)** The actual $\Delta$ margin changes monotonically with number of examples removed ($k$), corroborating the monotonicity assumption used in estimating data supports. **(Bottom)** On $x$, we instead plot the *predicted* $\Delta$ margin using datamodels. This shows that the linearity seen in Figure 8 manifests even at a local level.

counterfactuals after removing an additional 20% of highest datamodel weights, e.g. removing top $\widehat{k} \times 1.2$ examples for each test example. When an additional 20% of training examples are removed, 92% of test examples are misclassified. Hence, we use $\widehat{k} \times 1.2$ for our final estimates of SUPPORT($x$).

### F.3.1 Estimation using baselines

As baselines, we use the same guided search algorithm described above, but instead of using datamodel-predicted values to guide the search, we select the candidate subset using each of the baselines methods described in Appendix F.2. In particular, we choose the candidate subset $R_k$ for a given $k$ as follows:

1. *Representation distance*: top $k$ closest training examples to $x$ as measured by $\ell_2$ distance in the representation space of a pre-trained ResNet-9.

2. *Influence estimates (influence functions and TracIn)*: top $k$ training examples with highest (most positive) estimated influence on the target example $x$:

3. *Random*: first $k$ examples from a random ordering[21] of training examples from the same class as $x$.

### F.3.2 Heuristic estimates for data support

While we constructively estimate the data supports by training models on counterfactuals and using the above estimation procedure, we can also consider a simpler and cheaper heuristic to estimate SUPPORT($x$) assuming the fidelity of the linear datamodels: compute the smallest $k$ s.t. the sum of the $k$ highest datamodel weights for $x$ exceeds the average margin of $x$. In Figure F.2, we compare the predicted data supports based on this heuristic to the estimated ones from earlier, and find that they are highly correlated. In practice, this can be a more efficient alternative to quantify brittleness without additional model training (beyond the initial ones to estimate the datamodels).

## F.4 Brittleness to mislabeling

To study the brittleness of model predictions to mislabeling training images, we take the same 300 random CIFAR-10 test examples and analyze them as follows: First, we find for each example the *incorrect* class with the highest average logit (across $\sim$10,000 models trained on the full training set). Then, we construct counterfactual datasets similarly as in Appendix F.3 where we take the top $k = \{2, 4, ..., 256\}$ training examples

---

[21] The random ordering is fixed across different choices of $k$, but not across different targets.

**Figure F.2: Heuristic predictions for data supports.** For each of the 300 test examples shown, the *x*-coordinate represents the previous estimates based on counterfactuals, and the *y*-coordinate represents the heuristic estimate.



**Figure F.3: Brittleness to mislabeling.** We estimate an upper bound on the the smallest number of training images that can be mislabeled to flip a given target image. Much fewer images are required compared simply removing them (blue), as mislabeling provides additional signal to the model.

with the highest datamodel weights, but this time mislabel them with the incorrect class identified earlier. After training $T = 20$ models on each counterfactual, for each target example we estimate the number of mislabeled examples at which the expected margin becomes zero, using the same estimation procedure described in Appendix F.3. The resulting mislabeling brittleness estimates are shown in Figure F.3.

## F.5 Comparing raw effect sizes

Instead of comparing the data support estimates (which are derived quantities), here we directly compare the average counterfactual effect (i.e. delta margins) of groups selected using different methods. Figure F.4 shows again that datamodels identify much larger effects. Among baselines, we see that TracIn performs best, followed by representation distance. We also see that the representation baseline does not gain any additional signal from simple averaging over models.

## F.6 Effect of training stochasticity

As described in Appendix F.1, we re-train up to $T = 20$ models for *each* counterfactual to reduce noise that arises solely from stochasticity of model training. These additional samples significantly reduces unexplained variance: Figure F.5 shows the reduction in variance ("thickness" in the *y*-direction) and the resulting increase in correlation as the number of re-training runs is increased from $T = 1$ to $T = 20$.

## F.7 Transfer to different architecture

While the main premise of datamodeling is understanding how data is used by a given *fixed* learning algorithm, it is natural to ask how well datamodels can predict across different learning algorithms. We expect some degradation in predictiveness, as datamodels are fit to a particular learning algorithm; at the same time, we also expect some transfer of predictive power as modern deep neural networks are known to make similar predictions and errors [MMS+19; TSC+19].

Here, we study one of the factors in a learning algorithm, the choice of architecture. We take the same counterfactuals and evaluate them on ResNet-18 models, using the same training hyperparameters. As expected, the original datamodels continue to predict accurate counterfactuals for the new model class but with some degradation (Figure F.6).

**Figure F.4: Comparing effect sizes with baselines**. This shows the raw evaluations of counterfactuals generated using different methods, which were also used for estimating data supports. The *y*-axis shows the effect on margin averaged across all target examples when top *k* examples are removed for each target using each of the methods. Datamodels identify much larger effects compared to baselines. Among baselines, TracIn[PLS+20] performs the best. For representation distance, there is no noticeable gain from reducing stochasticity by averaging over more models (1000 vs 1).



**Figure F.5: Effect of model averaging.** *T* is the number of models trained per counterfactual.

**Figure F.6: Transfer across model classes.**

## F.8 Stress testing

Section 4.1 showed that datamodels excel at predicting counterfactuals across a variety of removal mechanisms. In an effort to find cases where datamodel predictions are not predictive of data counterfactuals, we evaluate the following additional counterfactuals:

- **Larger groups of examples (up to 20% of the dataset)**: we remove $k = 2560, 5120, 10240$ top weights using different datamodels $\alpha = 0.1, 0.2, 0.5, 0.75$. The changes in margin have more unexplained variance when larger number of images are removed; nonetheless, the overall correlation remains high ((Figure F.7)).

- **Groups of training examples whose predicted effects are *zero***: we remove $k = 20, 40, 80, 160, 320, 640, 1280, 2560$ examples with zero weight ($\alpha = 0.5$), chosen randomly among all such examples. All of tested counterfactuals had negligible impact on the actual margin, consistent with the prediction of datamodels (Figure F.8a).

- **Groups of examples whose predicted effect is *negative* according to baselines**: we test TracIn and influence functions. (We do not consider the representation distance baseline here is there is no obvious way of extracting this information from it.) Correlation degrades but remains high (Figure F.8b).

44

Note that the relative scale of the effects is much smaller compared to counterfactuals generated using datamodels (Figure F.8a).

In general, although there is some reduction in datamodels' predictiveness, we nevertheless find that datamodels continue to be accurate predictors of data counterfactuals.



**Figure F.7: Stress testing datamodel counterfactuals by removing a large number of images.** Plot shows datamodel counterfactuals from before ($k = 10, ..., 1280$) along with additional samples $k = 2560, 5120, 110240$ (shown with darker hue).



**(a)** Comparing counterfactuals with highest vs *zero* predicted effect using datamodels ($\alpha = 0.5$)

**(b)** Removing top $k$ positive and negative influence training examples according to baseline methods.

**Figure F.8: Stress testing counterfactual prediction.**

**Counterfactuals relative to a random control.** All of the counterfactuals studied so far are relative a fixed control (the entire training set). Here, we consider counterfactuals relative to a *random* control $S_0 \sim \mathcal{D}_S$ at $\alpha = 0.5$ (i.e. $|S_0| = \alpha|S|$). The motivation for considering the shifted control is two folds: first, the counterfactuals generated relative to such $S'$ are closer in distribution to the original distribution to which datamodels were fit to, so it is natural to study datamodels in this regime; second, this tests whether the counterfactual predictability is robust to the exact choice of the trainset. Latter is desirable, as ultimately we would like to understand how models behave on training sets similar in distribution to $S$, not the exact train set.

To implement above, after removing a target group $G$ from the full train set $S$, we subsample the remainder $S/G$ with probability $\alpha$. We adjust the control values accordingly to $\mathbb{E}_{S_0 \sim S}[f_{\mathcal{A}}(x; S_0)]$, where $\mathcal{D}_S$

is the $\alpha = 0.5$ subsampling distribution. The results show that one can indeed also predict counterfactuals relative a random control (Figure F.9).



**(a)** Random control ($\alpha = 0.5$ regime)     **(b)** Fixed control ($\alpha = 1.0$)

**Figure F.9: Datamodels can predict counterfactuals relative to random controls.** As in Figure 8, each point in the graphs above corresponds to a test example and a counterfactual trainset $S'$ (a subset of the full training set, $S$). The counterfactual is relative a *random* control $S_0 \sim \mathcal{D}_S$ with $\alpha = 0.5$, e.g. a set randomly subsampled at 50%. The $y$-coordinate of each point represents the expected *ground-truth* difference, in terms of model output on $x$, between training on a random $S_0$, and training on $S'$. The $x$-coordinate of each point represents the *datamodel-predicted* value of this quantity. **(a)** We use the $\alpha = 0.5$ datamodels to predict counterfactuals generated by removing, for each test example, the training inputs corresponding to the top-$k$ and bottom-$k$ (for several $k$) datamodel weights. **(b)** Same, but relative to a fixed control $S_0 = S$, e.g. the full train set.

## F.9   Additional plots for different $\alpha$ values

**Figure F.10: Varying $\alpha$ for counterfactual prediction (CIFAR-10).** Same plot as in Figure 8, except varying the datamodels used for prediction; each plot uses datamodels with the given $\alpha$. As before, each point in the graphs above corresponds to a test example and a subset $R(x)$ of the original training set $S$, identified by a (color-coded) heuristic. The $y$-coordinate of each point represents the *ground-truth* difference, in terms of model output on $x$, between training on $S$, and training on $S \setminus R(x)$. The $x$-coordinate of each point represents the *datamodel-predicted* value of this quantity.



**Figure F.11: Varying $\alpha$ for counterfactual prediction (FMoW).** Same plot as in Figure 8, except varying the datamodels used for prediction; each plot uses datamodels with the given $\alpha$. As before, each point in the graphs above corresponds to a test example and a subset $R(x)$ of the original training set $S$, identified by a (color-coded) heuristic. The $y$-coordinate of each point represents the *ground-truth* difference, in terms of model output on $x$, between training on $S$, and training on $S \setminus R(x)$. The $x$-coordinate of each point represents the *datamodel-predicted* value of this quantity.

47

# G  Nearest Neighbors

In this section we show additional examples of held-out images and their corresponding train image, data-model weight pairs.

## G.1  CIFAR

In Figure G.1 we show more *randomly* selected test images along with their positive and negative weight training examples. In Figure G.2 we show more examples of test i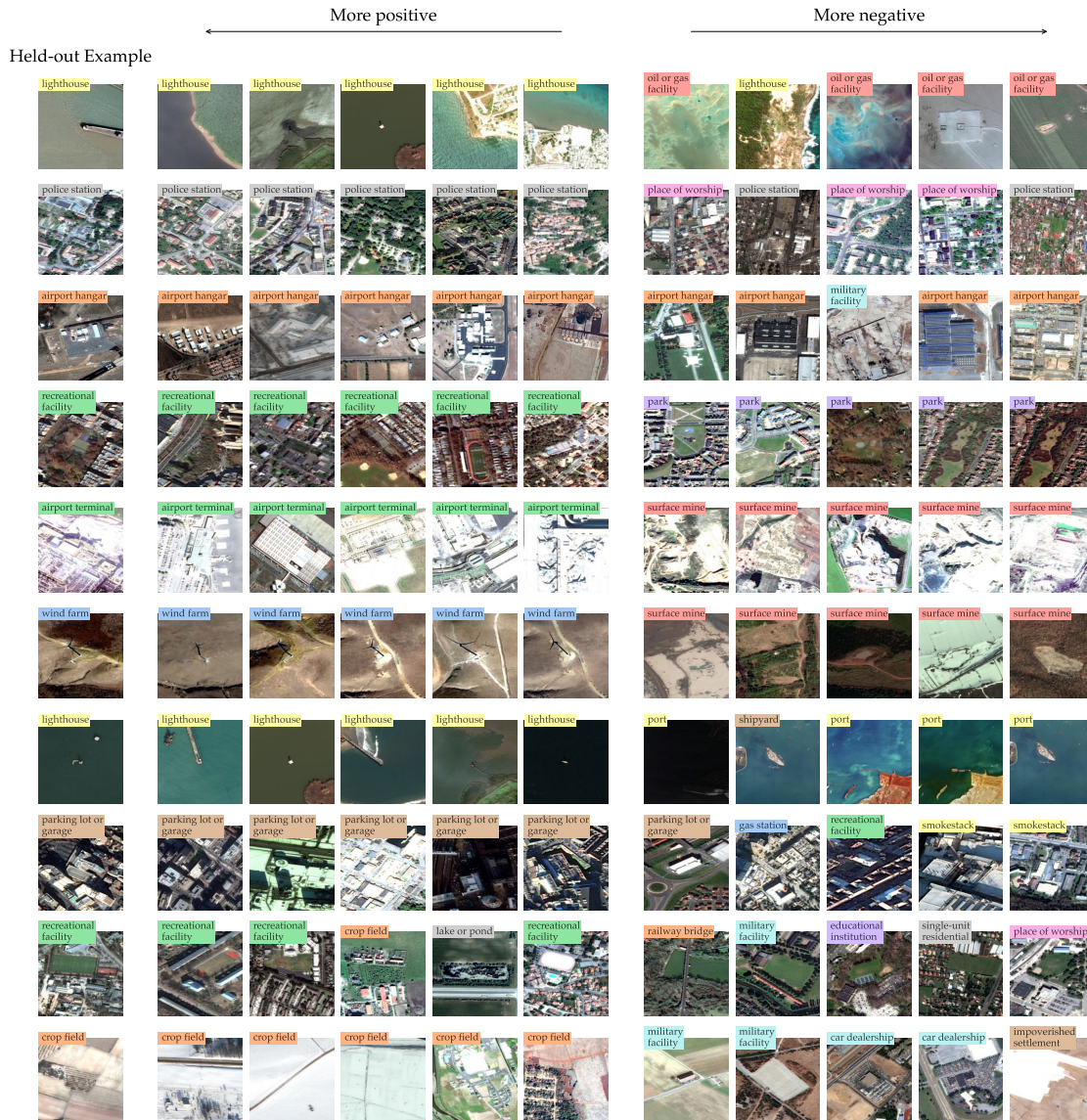mages and their corresponding top train images as we vary α. In Figure G.3 we compare most similar images to given test images identified using various baselines (see Appendix F.2 for their description).



**Figure G.1:** Additional examples of held-out images and their corresponding highest and lowest datamodel weight training images.

Held-out Example          Train Examples by Weight          Held-out Example          Train Examples by Weight

**Figure G.2:** Additional examples of held-out images and corresponding most relevant training examples, while varying $\alpha$.

**Figure G.3:** Comparisons of nearest neighbors found using different methods.

## G.2 FMoW

In Figure G.4 we show *randomly selected* target images along with their top-weight train images, using datamodels of different $\alpha$. In Figure G.5 we show more examples of test images and their corresponding top train images as we vary $\alpha$.



**Figure G.4:** FMoW examples of held-out images and their corresponding highest and lowest datamodel weight training images.

**Figure G.5:** FMoW examples of held-out images and corresponding most relevant training examples, while varying $\alpha$.

# H Train-Test Leakage

## H.1 CIFAR

**Task setup.** The general setup is as described Section 4.2.2. In the Amazon Mechanical Turk interface (Figure H.1), for each test image we displayed the top 5 and bottom 5 train examples by datamodel weight; the vast majority of potential leakage found corresponded to the top 5 examples. Nine different workers filled out each task. We paid 12 cents per task completed and used these qualifications: locale in US/CA/GB and percentage of hits approved $> 95\%$.



**Figure H.1:** The MTurk Interface, complete with instructions, shown to crowdsourced annotators. Note that there are 5 rows of images in the actual interface, some of which may require scrolling to get to.

**Examples.** Figure H.2 shows more examples of (train, test) pairs stratified by annotation score. While there is no ground truth due to lack of metadata, we see that the crowdsourced annotation combined with high quality candidates (as identified by datamodels) can effectively surface leaked examples.
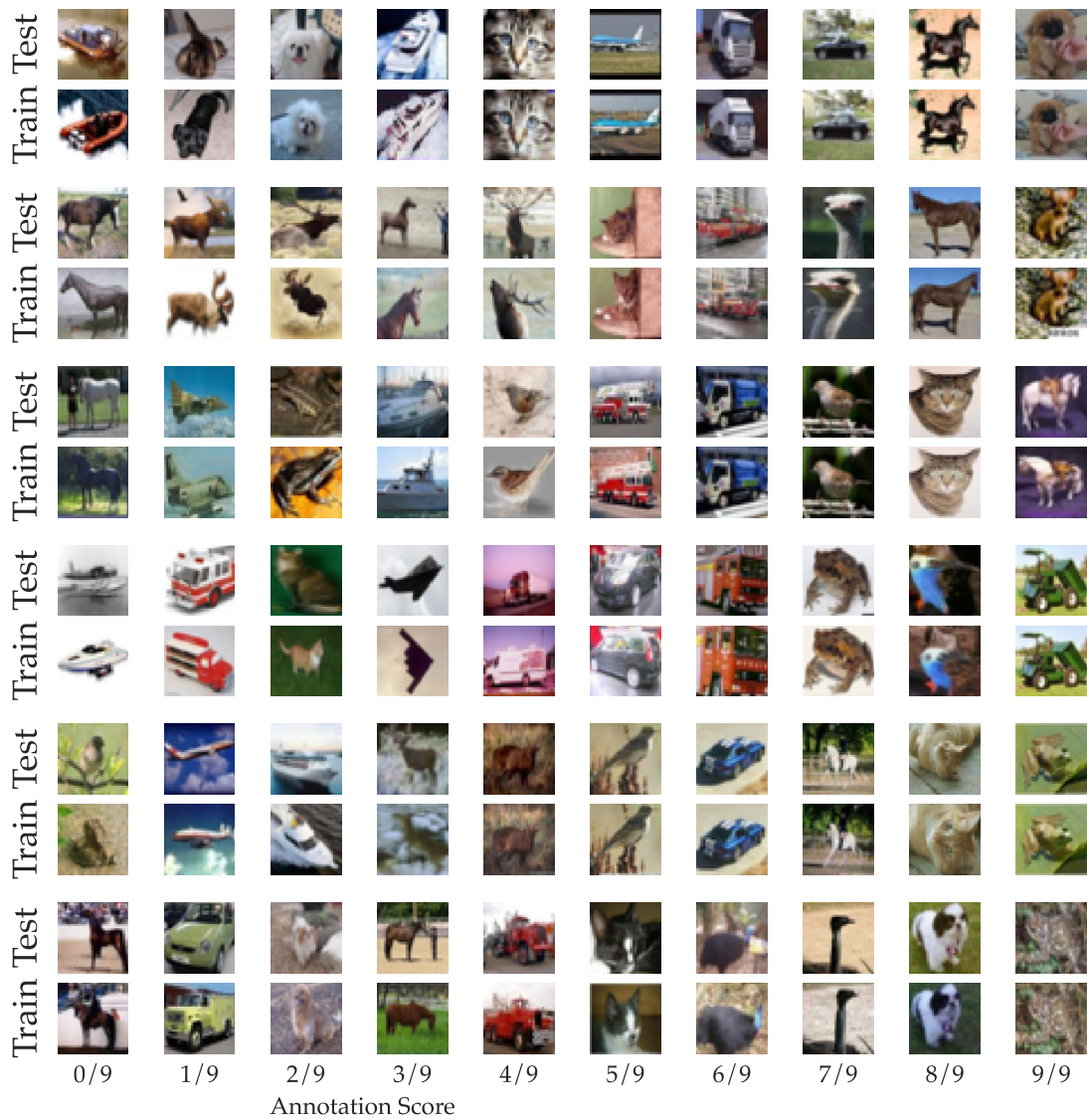
**Figure H.2:** More annotation scores paired with (train, test) leakage pairs. See Figure 11 and Section 4.2.2 for more information.

**Comparison with CIFAIR.** Barz and Denzler [BD20] present CIFAIR, a version of CIFAR with fewer duplicates. The authors define duplicates slightly differently than our definition of same scene train-test leakage (cf. Section 3.2 of their work and our interface shown in Figure H.1). They identify train-test leakage by using a deep neural network to measure representation space distances between images across training partitions and manually inspecting the lowest distances.

## H.2 FMoW

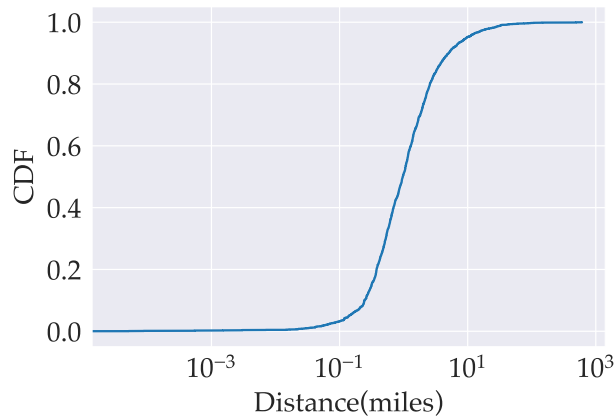Appendix Figure H.3 shows the CDF of each test image's minimum distance to a train image.



**Figure H.3:** CDF of distance in miles between each FMoW test set image and the nearest train set image.

# I Spectral Clustering

## I.1 Method

We use `sklearn`'s `cluster.SpectralClustering`. Internally, this computes similarity scores using the radial basis function (RBF) kernel on the datamodel embeddings. Then, it runs spectral clustering on the graph defined by the similarity matrix $A$: it computes a Laplacian $L$, represents each node using the first $k$ eigenvectors of $L$, and runs $k$-means clustering on the resulting feature representations. We use $k = 100$.

## I.2 Omitted results

Figure I.1 compares top clusters for the horse class across different $\alpha$. Figure I.2 shows additional clusters for eight other classes, apart from the ones shown in Figure 13.
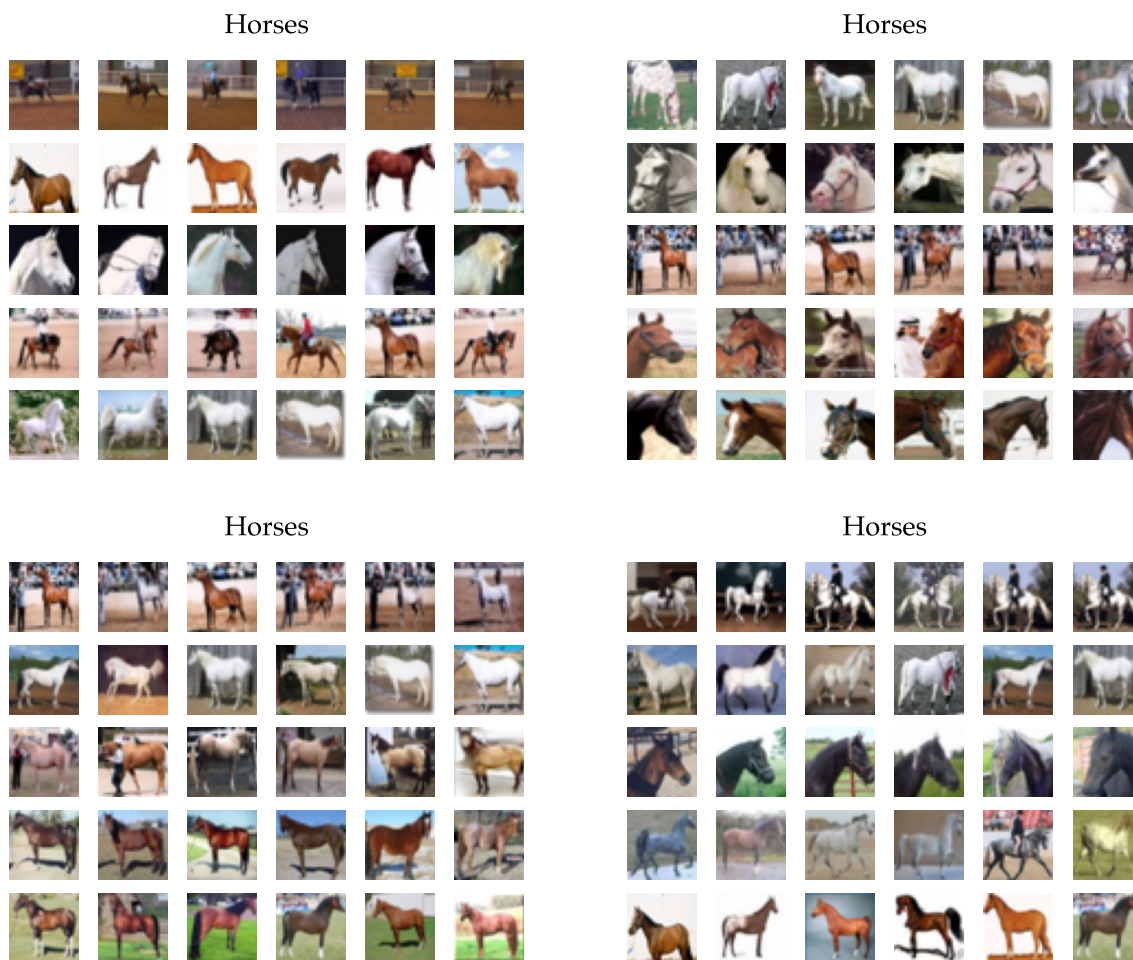


**Figure I.1:** Omitted spectral clustering results for datamodels computed with $\alpha = 10\%$ (top left), 20% (top right), 50% (bottom left), and 75% (bottom left).
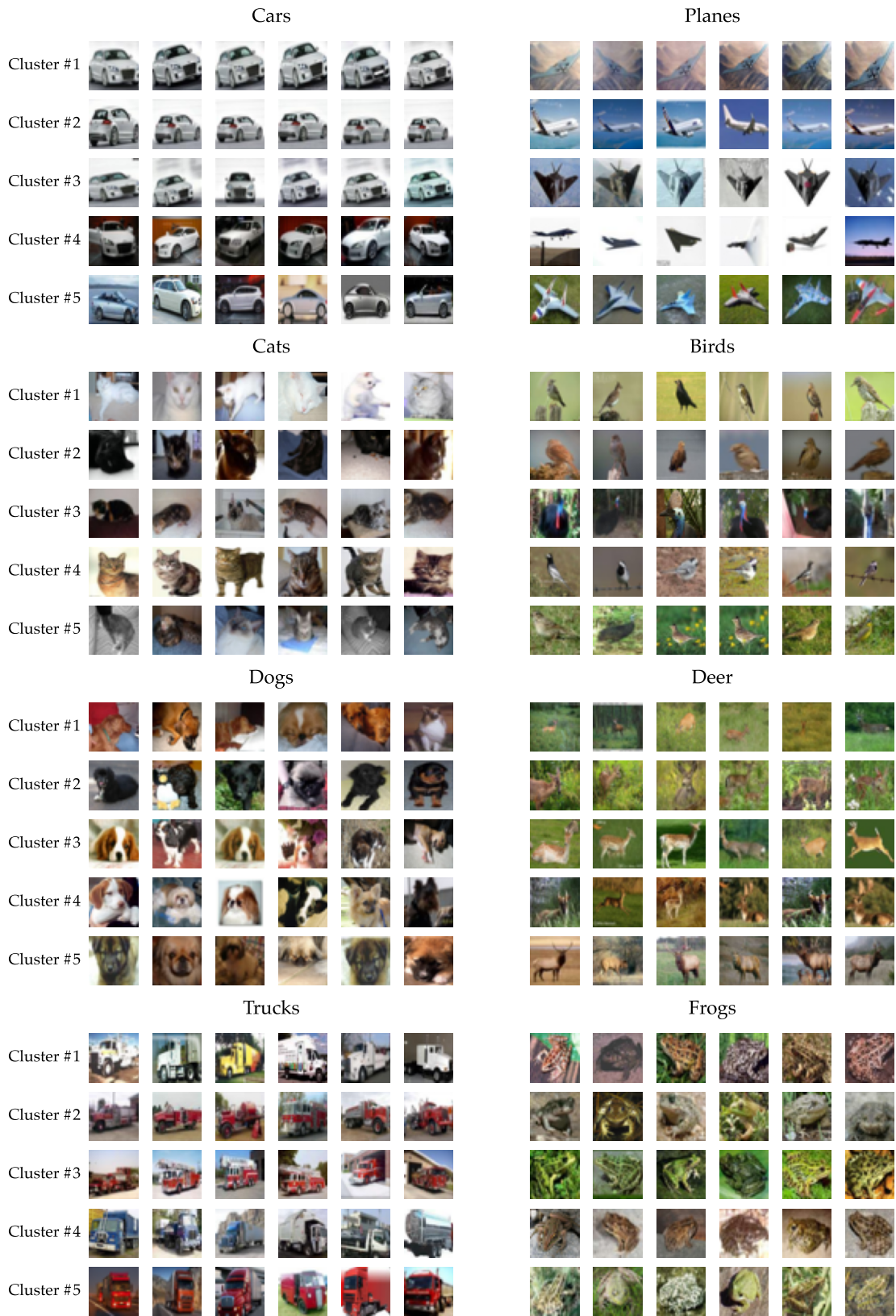
**Cars**

Cluster #1
Cluster #2
Cluster #3
Cluster #4
Cluster #5

**Planes**

Cluster #1
Cluster #2
Cluster #3
Cluster #4
Cluster #5

**Cats**

Cluster #1
Cluster #2
Cluster #3
Cluster #4
Cluster #5

**Birds**

Cluster #1
Cluster #2
Cluster #3
Cluster #4
Cluster #5

**Dogs**

Cluster #1
Cluster #2
Cluster #3
Cluster #4
Cluster #5

**Deer**

**Trucks**

Cluster #1
Cluster #2
Cluster #3
Cluster #4
Cluster #5

**Frogs**

**Figure I.2:** Omitted spectral clustering results for classes other than those in the main paper (Figure 13).

# J  PCA on Datamodel Embeddings

## J.1  CIFAR

**Setup.**  For the PCA experiments, we use datamodels for the training and test sets estimated with $\alpha = 0.5$ unless mentioned otherwise.

**Effective dimensionality.**  In Figure J.1, we compare the effective dimensionality of datamodel embeddings with that of a deep representation pretrained on CIFAR-10.
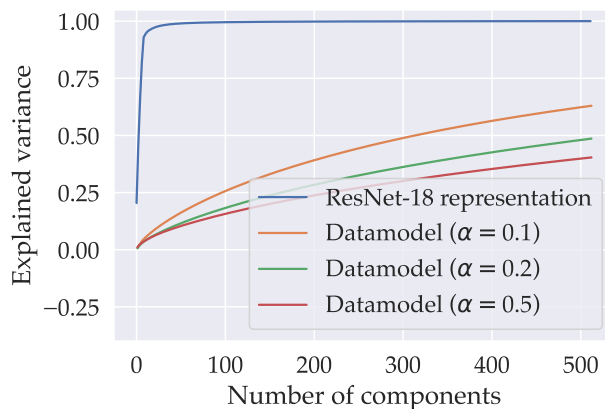


**Figure J.1: Datamodel embeddings have a higher effective dimension than deep representations**.  For different embeddings, we plot the cumulative fraction of variance explained by the top $k$ components while varying $k$.  For a network layer based embedding, 95% of the variation in embedding space is captured by the first 10 principal components; meanwhile, datamodel embeddings need up to 500 components to capture even half of the variance.  Here, we use a ResNet-18 model instead of ResNet-9 as it has more features in the representation layer (512 vs. 128); the plot looks similar for ResNet-9.

**Analyzing model-faithfulness.**  To see whether PCA directions reflect model behavior, we look at how "removing" different principal components affect model predictions.  More precisely, we remove training examples corresponding to:

- Top $k$ most positive coordinates of the principal component vector
- Top $k$ most negative coordinates of the principal component vector

Then, for each principal component direction considered, we measure their impact on three groups of held-out samples:

- The top 100 examples by most positive projection on the principal component
- The bottom 100 examples by most negative projection on the principal component
- The full test set

For each of these groups, we measure the mean change in margin after removing different principal component directions. Our results (Figure J.2) show that:

- Removing the most positive coordinates of the PC decreases margin on the test set examples with the most positive projections on the PC and increases margin on the examples with the most negative projections on the PC.

- Removing the most negative coordinates of the PC has the opposite effect, increasing margin on the positive projection examples and decreasing margin on the negative projection examples.
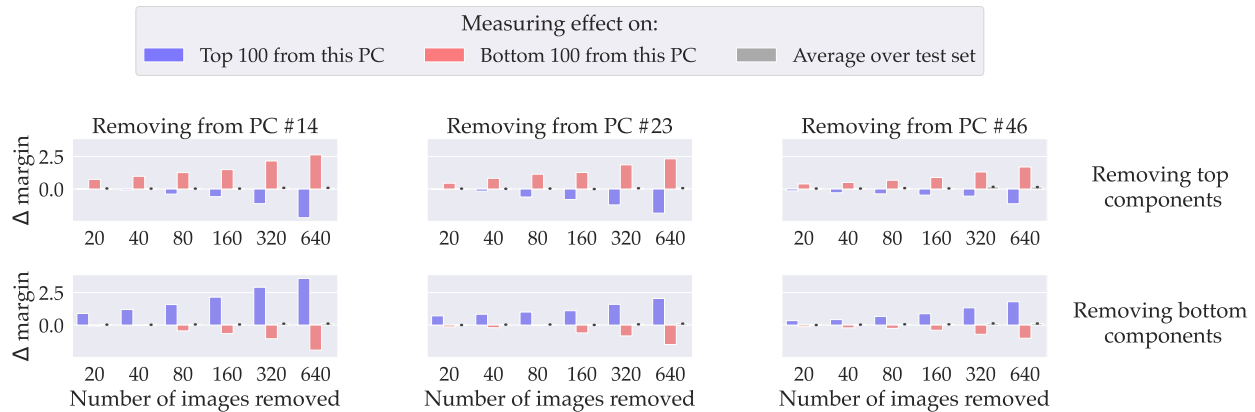
**Figure J.2: PCA directions generalize and capture "orthogonal" directions.** For each of the three principal components (PCs) above (randomly chosen from top 50), we consider the counterfactual of removing the training examples corresponding to the top or bottom $k$ coordinates in the PC, and measure its average effect on different groups: (red) test examples with the highest projections on the PC; (blue) test examples with the lowest projections on the PC; and (grey) the entire test set. The direction of the effect is consistent with the datamodel embeddings; removing top (resp. bottom) coordinates decrease (resp. increase) the average margin on test examples whose embeddings are most aligned with the PC. Moreover, the negligible impact over the test set in aggregate shows that the different PCs, which are orthogonal in the embedding space (by definition), are also approximately "orthogonal" in terms of their effect on model predictions.

- Increasing the size of each removed set increases the effect magnitude.

- Removing PC's have negligible impact on the aggregate test set, indicating that the impact of different PC's are roughly "orthogonal," as one would expect based on the orthogonality of the PCs.

- Lastly, Figure J.3 shows that datamodels can accurately predict the counterfactual effect of the above removed groups, similarly as in Figure 8.
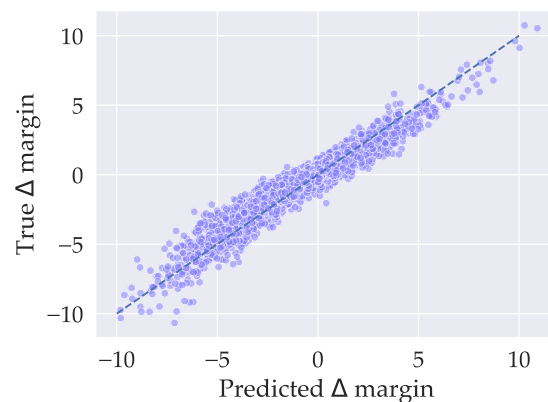


**Figure J.3: Datamodels predict the effect of "removing" principal components**. Each point corresponds to a *PCA counterfactual*: removing training examples with the largest weights in the principal component (i.e., top-$k$ most positive or negative coordinates), and evaluating on test examples whose embeddings most align with the PC (e.g. smallest cosine distance). The $y$-coordinate of each point represents the *ground-truth* counterfactual effect (evaluated by retraining $T = 20$ times). The $x$-coordinate of each point represents the *datamodel-predicted* value of this quantity.

**Representation baseline.** In Figure J.4 we show the top principal components computed using a representation embedding.

**Additional results.** In Figure J.5 we show additional PCs from a datamodel PCA.
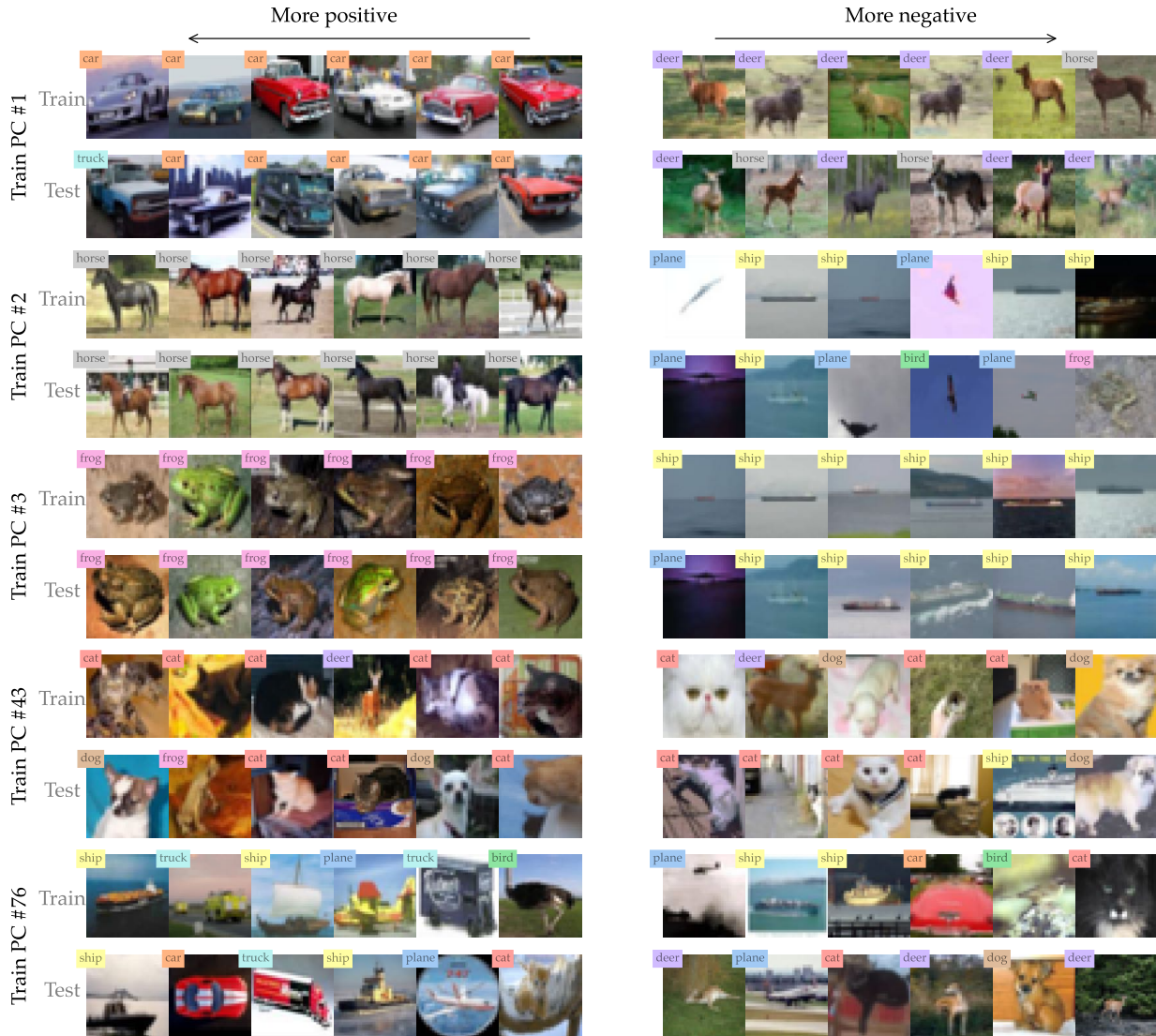


**Figure J.4: Representation-based baseline for PCA.** Visualization of highest magnitude images along top principal components of *representation* embeddings for CIFAR-10. In each row $i$, on the left we show the images with the highest normalized projections onto $v_i$, and on the right the images with the lowest projections. The PCs seem less coherent than those obtained from running PCA on datamodel embeddings.
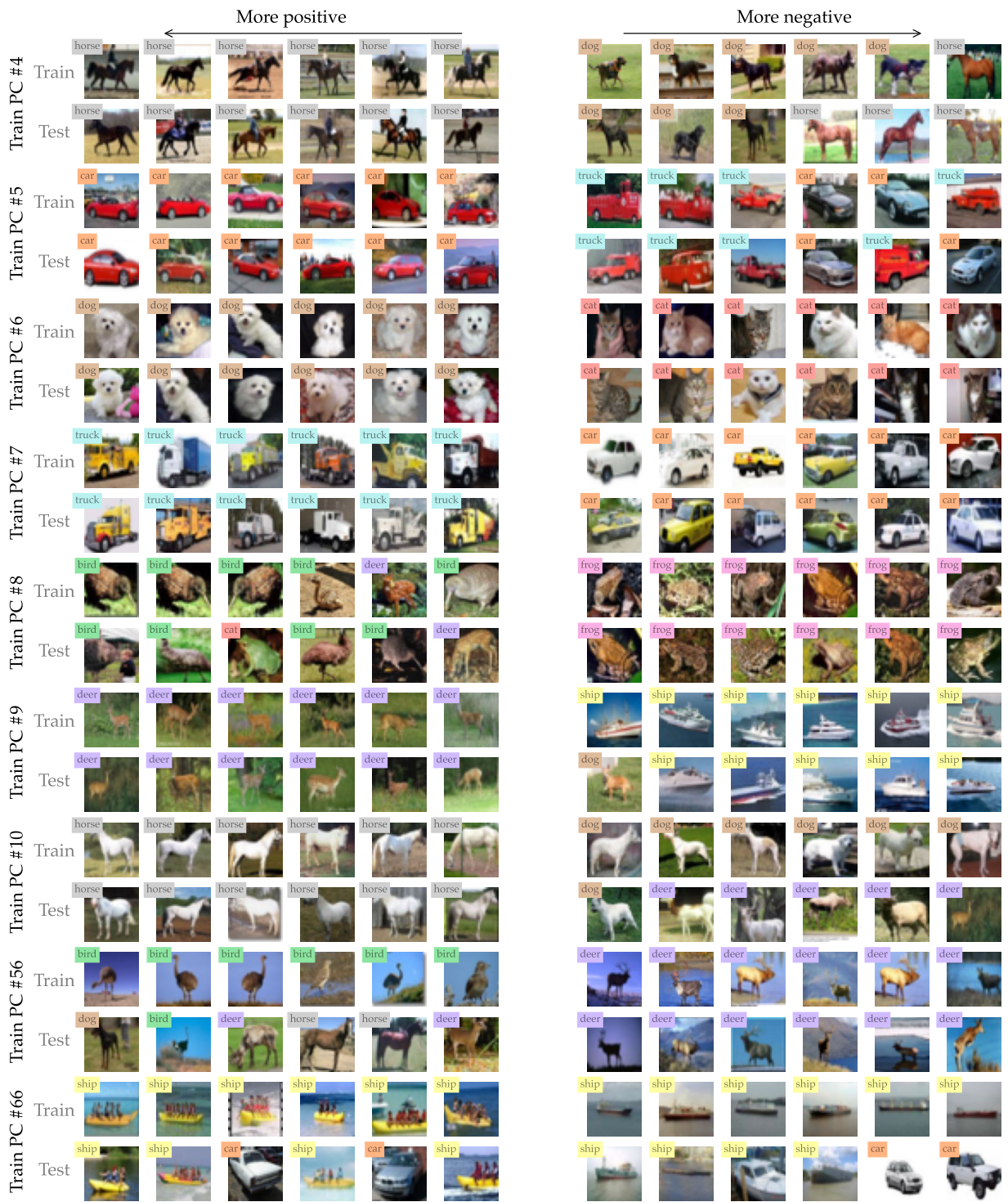
**Figure J.5:** The remainder of the top 10 PCA directions and two selected directions.
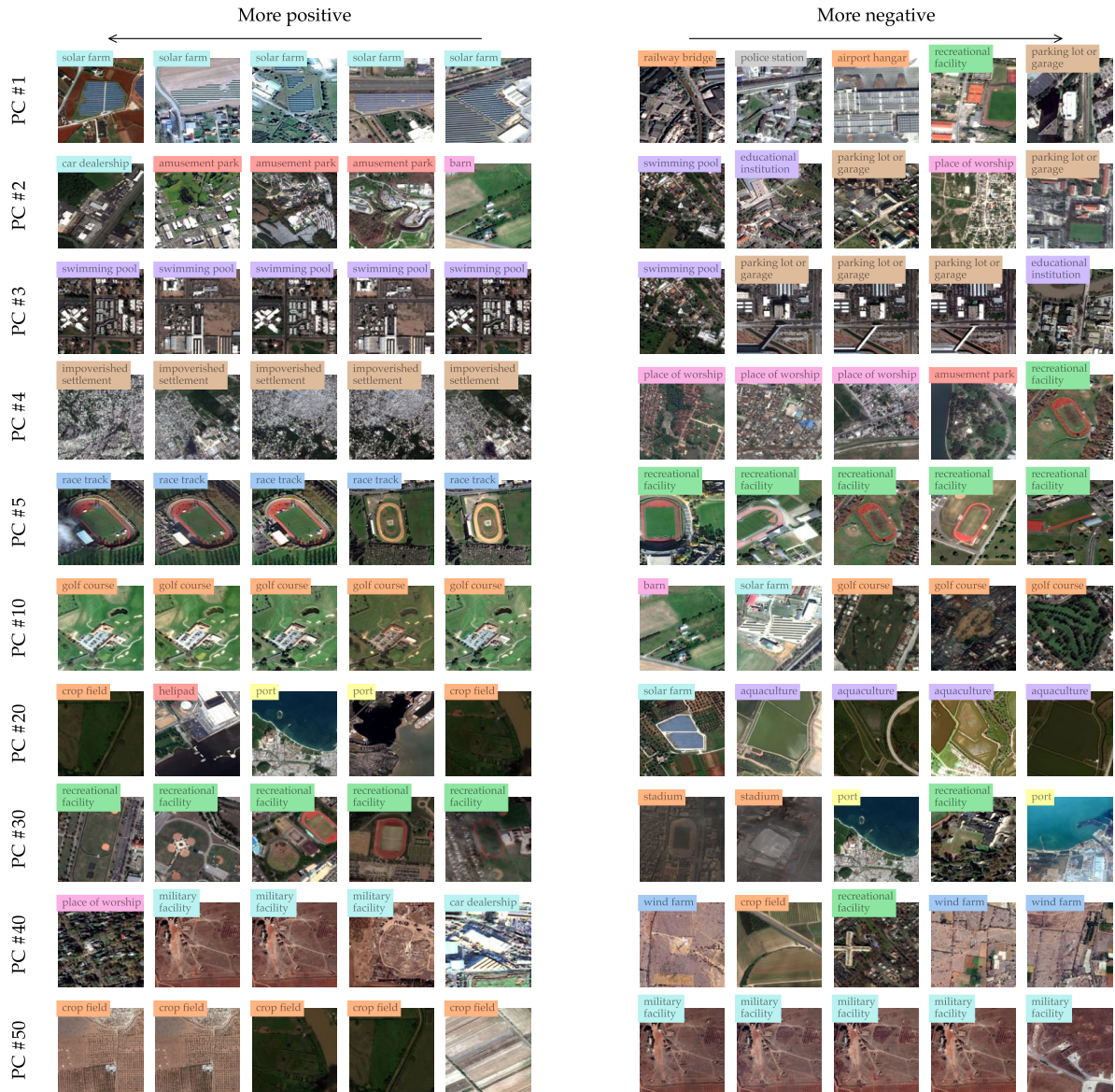
## J.2 FMoW



**Figure J.6:** FMoW top PCA components, using $\alpha = 20\%$. Top 5 and 5 selected from the top 50.

# K Connection between Influence Estimation and Datamodels

## K.1 Proof of Lemma 1

**Lemma 1.** *Fix a training set S of size n, and a test example x. For $i \in [m]$, let $S_i$ be a random variable denoting a random 50%-subset of the the training set S. Let $\boldsymbol{w}_{infl} \in \mathbb{R}^n$ be the estimated empirical influences (15) onto x estimated using the sets $S_i$. Let $\boldsymbol{w}_{OLS}$ be the least-squares estimator of whether a particular model will get image x correct, i.e.,*

$$\boldsymbol{w}_{OLS} := \arg\min_w \frac{1}{m} \sum_{i=1}^{m} \left( w^\top \boldsymbol{z_i} - \mathbf{1}\{\text{model trained on } S_i \text{ correct on } x\} \right)^2, \qquad \text{where } z_i = 2 \cdot \mathbf{1}_{S_i} - \mathbf{1}_n.$$

*Then, as $m \to \infty$,*

$$\left\| (1 + 2/n)\boldsymbol{w}_{OLS} - \frac{1}{2}\boldsymbol{w}_{infl} \right\|_2 \to 0.$$

*Proof.* For convenience, we introduce the $m \times n$ binary *mask matrix* $A$ such that $A_{ij}$ is an indicator for whether the $j$-th training image was included in $S_i$. Note that $A$ is a random matrix with fixed row sum of $n/2$. Next, we define the *output vector* $y \in \{0,1\}^m$ that indicates whether a model trained on $S_i$ was correct on x. Finally, we introduce the *count matrix* $C = \mathrm{diag}(\mathbf{1}^\top A)$, i.e., a diagonal matrix whose entries are the columns sums of $A$, e.g. the number of times each example appears across $m$ different masks.

We begin with $\boldsymbol{w}_{OLS}$. Consider the $n \times n$ matrix $\Sigma = \frac{1}{m} Z^\top Z = \frac{1}{m}(2 \cdot A - \mathbf{1}_{m \times n})^\top (2 \cdot A - \mathbf{1}_{m \times n})$. The diagonal entries of this matrix are $\Sigma_{ii} = 1$ (due to $A$ having constant row sum), while the off-diagonal is

$$\Sigma_{ab} = \frac{1}{m} \sum_{i=1}^{m} \begin{cases} +1 & \text{if training image } x_a, x_b \in S_i \text{ or } x_a, x_b \notin S_i \\ -1 & \text{otherwise.} \end{cases}$$

Since $\Sigma$ has bounded entries ($|\Sigma_{ab}| \leq 1$), we have that for fixed $n$, $\lim_{m \to \infty} \Sigma = \mathbb{E}[\Sigma]$, and in particular

$$\Sigma_{ab} \to \mathbb{P}(x_a, x_b \in S_i \text{ or } x_a, x_b \notin S_i) - (1 - \mathbb{P}(x_a, x_b \in S_i \text{ or } x_a, x_b \notin S_i))$$

$$\mathbb{P}(x_a, x_b \in S_i \text{ or } x_a, x_b \notin S_i) = 2 \cdot \left( \frac{\frac{n}{2}}{n} \cdot \frac{\frac{n}{2} - 1}{n} \right) = \frac{1}{2} - \frac{1}{n}$$

$$\text{Thus, } \Sigma_{ab} \to -\frac{1}{2n}.$$

Now, using the Sherman-Morrison formula,

$$\Sigma^{-1} = \frac{n}{n+2} \left( I + \frac{2}{n} \mathbf{1}_{n \times n} \right)$$

By construction, the row sums of $Z = 2 \cdot A - \mathbf{1}_{m \times n}$ are 0, and so $\mathbf{1}_{n \times n} \cdot Z^\top = 0$. Thus,

$$\boldsymbol{w}_{OLS} = (Z^\top Z)^{-1} Z^\top y = \frac{1}{m} \left( \frac{1}{m} Z^\top Z \right)^{-1} Z^\top y = \frac{1}{m} \cdot \frac{n}{n+2} Z^\top y.$$

We now shift our attention to the empirical influence estimator $\boldsymbol{w}_{infl}$. Using our notation, we can rewrite the (vectorized) empirical influence estimator (15) as:

$$\boldsymbol{w}_{infl} = C^{-1} A^\top y - (m \cdot I_n - C)^{-1} (\mathbf{1}_{m \times n} - A)^\top y$$

$$= \left( C^{-1} - (m \cdot I_n - C)^{-1} \right) A^\top y - (m \cdot I_n - C)^{-1} \mathbf{1}_{m \times n}^\top y$$

$$= m \cdot C^{-1} (m \cdot I_n - C)^{-1} A^\top y - (m \cdot I_n - C)^{-1} \mathbf{1}_{m \times n}^\top y$$

$$= (m \cdot I_n - C)^{-1} \left( m \cdot C^{-1} A^\top - \mathbf{1}_{m \times n}^\top \right) y.$$

Now, as $m \to \infty$ for fixed $n$, the random variable $mC^{-1}$ converges to $2 \cdot I$ with probability 1. Thus,

$$m \cdot AC^{-1} - \mathbf{1}_{m \times n} \to 2 \cdot A - \mathbf{1}_{m \times n},$$

and the empirical influence estimator $\boldsymbol{w}_{infl} \to \frac{2}{m} Z^\top y$, which completes the proof. □

| Algorithm | # models ($m$) | Output type | Spearman $r$ | MSE | AUC | Difference |
|---|---|---|---|---|---|---|
| Diff. of means | 25,000 | Correctness | 0.028 | N/A | 0.529 | |
| Diff. of means | 100,000 | Correctness | 0.053 | N/A | 0.555 | Under $\rightarrow$ Over-determined |
| Diff. of means | 100,000 | Margin | 0.213 | 2.052 | 0.653 | Output type |
| LASSO | 100,000 | Margin | 0.320 | 1.382 | 0.724 | Explicit datamodel |

Table K.1: **Disentangling the effect of different factors in datamodel performance.** Each row shows a different estimator for datamodels. We begin with the empirical influence (or difference of means) on correctness computed with 25,000 models, which is in the overparameterized regime (as there are $d = 50,000$ variables). Then, we increase the number of models to an underparameterized regime. Next, we change the output type from correctness to margins. Lastly, we change the estimation algorithm from difference of means (which is approximately equivalent to OLS, as shown in Appendix K.1) to LASSO. Each of these changes brings about significant gains in the signal captured by datamodels, as measured by Spearman rank correlation, MSE, or AUC.

## K.2 Evaluating influence estimates as datamodels

Lemma 1 suggests that we can re-cast empirical influence estimates as (rescaled) datamodels fit with least-squares loss. Under this view, (i.e., ignoring the difference in conceptual goal), we can differentiate between explicit datamodels and those arising from empirical influences along three axes:

- **Estimation algorithm**: Most importantly, datamodels *explicitly* minimize the squared error between true and predicted model outputs. Furthermore, datamodels as instantiated here use (a) a sparsity prior and (b) a bias term which may help generalization.

- **Scale**: Driven by their intended applications (where one typically only needs to estimate the highest-influence training points for a given test point), empirical influence estimates are typically computed with relatively few samples (i.e., $m < d$, in our setting) [FZ20]. In contrast, we find that for datamodel loss to plateau, one needs to estimate parameters using a much larger set of models.

- **Output type**: Finally, datamodels do not restrict to prediction of a binary correctness variable—in this paper, for example, for deep classification models we find that *correct-class margin* was best both heuristically and in practice.

In this section, we thus ask: how well do the rescaled datamodels that arise from empirical influence estimates predict model outputs? We address this question in the context of the three axes of variation described above. In order to make results comparable across different outputs types (e.g., correctness vs. correct-class margin), we measure correlation (in the sense of Spearman [Spe04]) between the predicted and true model outputs, in addition to MSE where appropriate. To ensure a conservative comparison, we also measure performance as a predictor of *correctness*. In particular, we treat $w^\top \mathbf{1}_{S_i}$ as a continuous predictor of the binary variable $\mathbf{1}\{$model trained on $S_i$ is correct on $x\}$, and compute the AUC of this predictor (intuitively, this should favor empirical influence estimates since they are computed using correctnesses directly).

In Table K.1 we show the difference between empirical influence estimates (first row) and our final datamodel estimates (last row), while disentangling the effect of the three axes above using the rows in between. As expected, there is a vast difference in terms of correlation between the original empirical influence estimates and explicit datamodels. We further illustrate this point in Figure K.1, where we show how the correlation, MSE, and AUC vary with $m$ for both empirical influence estimates and datamodels, as well as an intermediate estimator that uses the estimation procedure of empirical influence estimates but replaces correctness with margin.

## K.3 Testing Lemma 1 empirically

In this section, we visualize the performance of empirical influence estimates ([FZ20]) as datamodels. In Figures K.2a and K.2b we plot the distributions of $w_{infl}^\top \mathbf{1}_{S_i} | y_i$ for different CIFAR-10 test examples; Figure
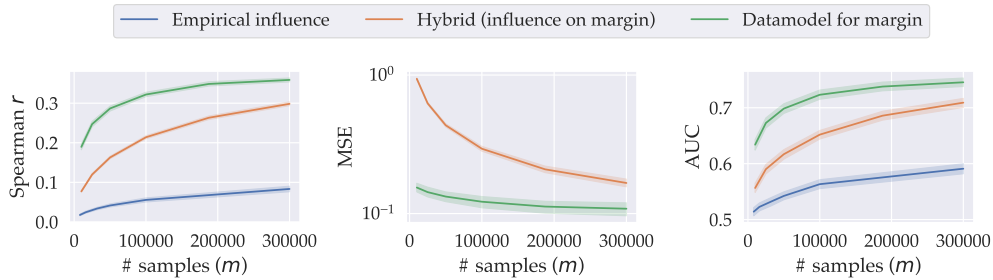
**Figure K.1: Datamodels have significantly better sample complexity than empirical influences.** We compare three estimators—empirical influence, empirical influence on margins, and $\ell_1$-regularized linear regression on margins (datamodels)—across a wide range of sample sizes on three different metrics. All metrics are averaged over the entire test set (i.e. over 10,000 datamodels). For MSE, we only show the estimators on margins as different output types are incomparable. Across all metrics, datamodels capture significantly more signal than empirical influences using the same number of samples. Conversely, datamodels need far fewer samples to reach the same level of performance.

K.2a shows these "conditional prediction distributions" for subsets $S_i$ that were used to estimate the empirical influence, while Figure K.2b shows the corresponding distributions on held-out (unseen) subsets $S_i$. The figures suggest that (i) indeed, empirical influences are somewhat predictive of the correctness $y_i$, (ii) their predictiveness increases as number of samples $m \to \infty$ but is still rather low, and (c) a significant amount of the prediction error is generalization error, as the train predictions in Figure K.2a are significantly better-separated than the heldout predictions in Figure K.2b.
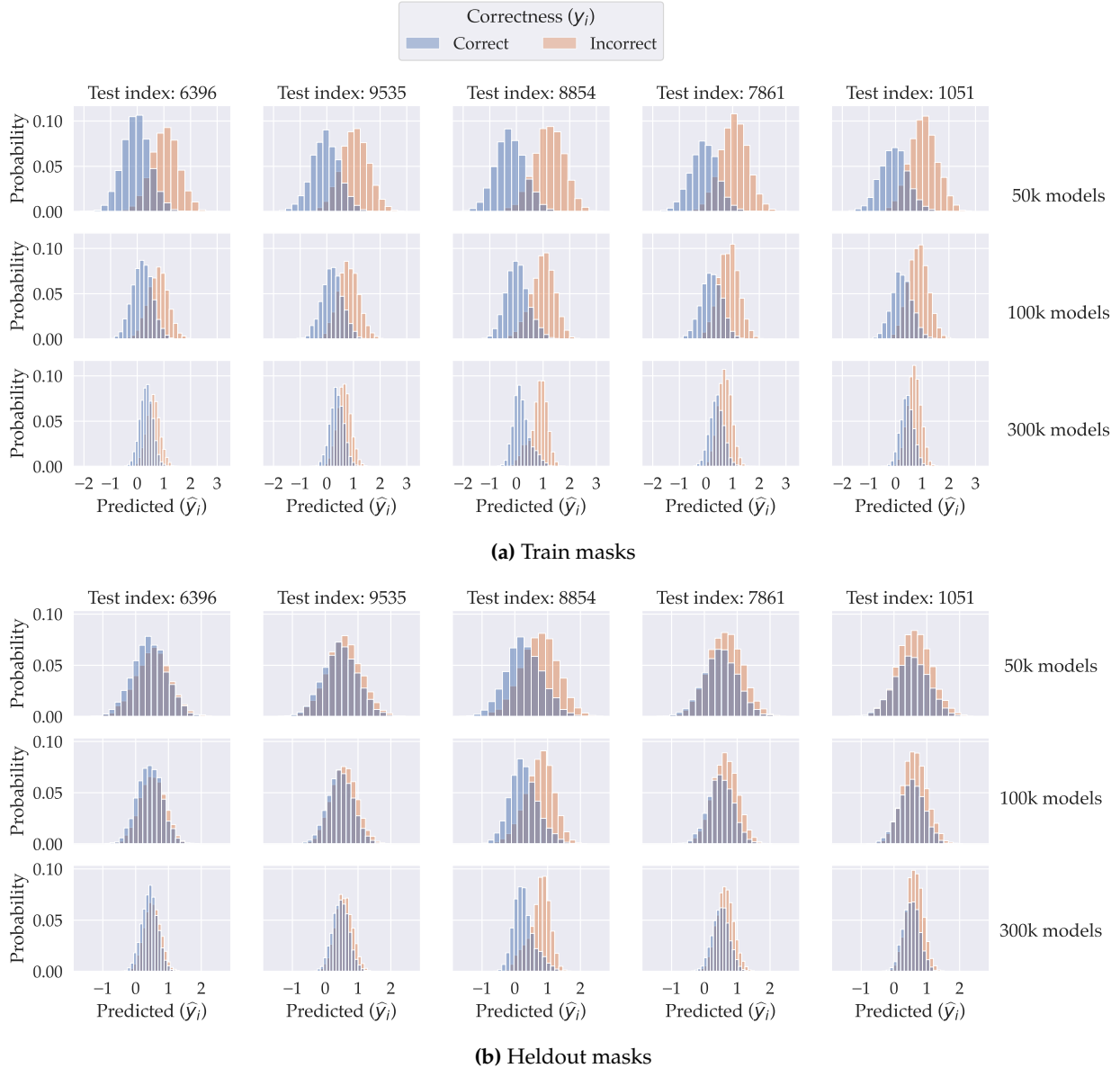
**(a)** Train masks



**(b)** Heldout masks

**Figure K.2: Empirical influence estimates are (weak) datamodels.** Each histogram illustrates the performance of empirical influences when the output of the corresponding datamodel is used as a statistic to distinguish between correct and incorrect predictions on the target example. Empirical influences can predict correctness on the "train set" of subsets (i.e. the masks used to estimate them), but suffer from significant generalization error when evaluated on a held-out set of subsets.

## K.4 View of empirical influences as a Taylor approximation

Lemma 1 shows that we can interpret empirical influences as (rescaled) estimates of the weights of a *linear* datamodel. Here, we give an alternative intuition for why this is the case, even though the definition of empirical influence does not explicitly assume linearity anywhere: we show that the influences define a first-order Taylor approximation of the multilinear extension $f$ of our target function $F$ of interest, where the influences (approximately) correspond to first-order derivatives of $f$.

Recall that we want to learn some output of interest $F : 2^T \to \mathbb{R}$, say the probability of correctness on a test example $z$, as a function of the examples $S \subset T$ included in the training set. We first extend this function continuously so that we can take its derivatives. The multilinear extension [Owe72] of set function $F$ to the domain $[0,1]^n$ ($|T| = n$) is given by:

$$f(x) = \sum_{S \subseteq T} F(S) \prod_{i \in S} x_i \prod_{i \notin S} (1 - x_i) \tag{20}$$

$f(x)$ also has an intuitive interpretation: it is the expected value of $F(S)$ when $S$ is chosen by including each $x_i$ in the input with probability $x_i$.

Next, we take the derivative of $f$ w.r.t. to the input $x_i$:

$$\frac{\partial f}{\partial x_i} = \underbrace{\sum_{S \subseteq T, i \in S} F(S) \prod_{j \in S, j \neq i} x_j \prod_{j \notin S} (1 - x_j)}_{\underset{s_j \sim Bern(x_j), s_i = 1}{\mathbb{E}} F(S)} - \underbrace{\sum_{S \subseteq T, i \notin S} F(S) \prod_{j \in S} x_j \prod_{j \notin S, j \neq i} (1 - x_j)}_{\underset{s_j \sim Bern(x_j), s_i = 0}{\mathbb{E}} F(S)}$$

Note that because $f$ is multilinear, the derivative w.r.t. to $x_i$ is constant in $x_i$, but not w.r.t. to other $x_j$. Now, observe that the above expression evaluated at $x_j = \alpha$, $\forall x_j$ corresponds approximately[22] to $\alpha$-subsampled influence $\theta_i$, of $i$ on $F$: the first term corresponds (using our earlier interpretation) to the expectation of $F(S)$ conditional on $S$ including $i$, and the second to that conditional on $S$ excluding $i$.

Finally, the first-order Taylor approximation of $f$ around an $x$ is given as:

$$f(x) \approx F(\emptyset) + \sum_i \frac{\partial f}{\partial x_i} \cdot x_i \approx F(\emptyset) + \sum_i \theta_i \cdot x_i$$

where $\theta_i$ are the empirical influences.

**The role of $\alpha$.** The above perspective provides an alternative way to think the role of the sampling fraction $\alpha$. The weights $\theta_i$ depend on the regime we are interested in; if we use $\alpha$-subsampled influences, then we are effectively taking a local linear approximation of $f$ in the regime around $\vec{x} = \alpha \cdot \vec{1}$.

**Remark.** Though we include the exposition above for completeness, this is a classical derivation that has appeared in similar form in prior works [Owe72]. Another connection is that *Shapley value* is equivalent to the integral of $f$ along the "main diagonal" of the hypercube; it is effectively empirical influences averaged uniformly over the choice of $\alpha$.

---

[22]There are two sources of approximation here. First, the $\alpha$-subsampling used in our datamodel definition is defined globally (e.g. $\alpha$ fraction of entire train set), which is different from the i.i.d. $Bern(\alpha)$ sampling that is considered here. Second, we only observe noisy versions of $F(S)$.