

# Are Labels Always Necessary for Classifier Accuracy Evaluation?

Weijian Deng Liang Zheng  
Australian National University  
{firstname.lastname}@anu.edu.au

## Abstract

To calculate the model accuracy on a computer vision task, e.g., object recognition, we usually require a test set composing of test samples and their ground truth labels. Whilst standard usage cases satisfy this requirement, many real-world scenarios involve unlabeled test data, rendering common model evaluation methods infeasible. We investigate this important and under-explored problem, Automatic model Evaluation (AutoEval). Specifically, given a labeled training set and a classifier, we aim to estimate the classification accuracy on unlabeled test datasets. We construct a meta-dataset: a dataset comprised of datasets generated from the original images via various transformations such as rotation, background substitution, foreground scaling, etc. As the classification accuracy of the model on each sample (dataset) is known from the original dataset labels, our task can be solved via regression. Using the feature statistics to represent the distribution of a sample dataset, we can train regression models (e.g., a regression neural network) to predict model performance. Using synthetic meta-dataset and real-world datasets in training and testing, respectively, we report a reasonable and promising prediction of the model accuracy. We also provide insights into the application scope, limitation, and potential future direction of AutoEval.

## 1. Introduction

Model evaluation is an indispensable step in almost every computer vision task. Using a test set that is unseen during training, the goal of evaluation is to estimate a model's (hopefully) unbiased accuracy when deployed in real-world scenarios. In most cases, we are provided with a labeled test set, allowing us to calculate the accuracy of a model by comparing the predicted labels with the ground truth labels (e.g., Fig. 1(a)). In the community, there are many well-established benchmarks (e.g., ImageNet [7] and COCO [26]) that provide various types of evaluation metrics. For example, top-1 error, commonly used in image classification, indicates whether the predicted class is the

same as the ground truth. There are some other metrics such as mean average precision in object detection [26] and panoptic quality [22] in panoptic segmentation.

Compared with the evaluation on these benchmarks, evaluating model performance for real-world deployment is not that straightforward. Often, real-world data follow distributions that differ from the original training distribution. In this case, a model's performance on the test sets in a benchmark may not reflect that achieved during deployment. If we still need to have an estimation of the model's accuracy in this scenario, we have to re-evaluate it on the real-world data. However, we often face scenarios where annotations of test samples are not provided. Furthermore, it can be very complex and expensive to manually gather labels. Even if acquired, these samples may only cover a very limited set of conditions, adding bias to the evaluated performance. For example, it is very expensive to annotate test samples for license plate recognition systems; even label is gathered for every car, it still can not capture the diversity of real-world circumstances such as lighting and weather condition. This raises an interesting question: *can we estimate model performance on a test set without test labels?*

To answer this question, this paper introduces the Automatic model Evaluation (AutoEval) problem. Given a classifier trained on a training set, the goal is to estimate its accuracy on an unlabeled test set. Here, we introduce an example in Fig. 1(b). Given a digit classifier trained on MNIST [23], we want to predict the classification accuracy on a test set *without* ground truths. This problem is challenging, as a test set contains many images, and each image has varied and rich visual contents. However, by visually inspecting the obvious differences between test and training sets, we can infer that the accuracy on the test set is low.

From this observation, we study AutoEval by considering the distribution difference between training and test sets and how it effects classifier accuracy. Existing literature gives us important hints. Data distributions can be represented by first and second-order statistics of the mean vector of output image feature representations [36, 32, 15]. For example, distribution difference can be estimated via Fréchet Distance (FD) [12] or maximum mean discrepancy (MMD)

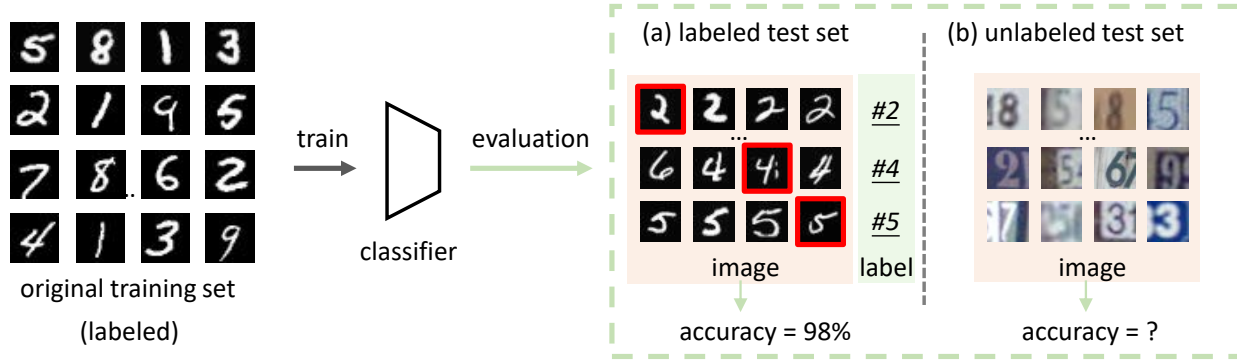


Figure 1. Problem illustration. Given a classifier trained on a training set, we can gain a hopefully unbiased estimate of its real-world performance by evaluating it on an unseen labeled test dataset, as shown in (a). However, in many real-world deployment scenarios, we are presented with unlabeled test datasets (b), and as such are unable to evaluate our classifier using common metrics. This inspired us to explore the problem of Automatic model Evaluation.

	Image Classification	AutoEval
Sample	Image	Dataset (sample set)
Label	Sample class ground truth	Accuracy of model on sample set
Train Set	Set of labeled images	Set of synthetic labeled sample sets (meta set)
Test Set	Set of unseen labeled images	Set of unseen labeled real-world datasets
Loss	Class cross-entropy	Predicted accuracy RMSE
Task	Classify images	Predict accuracy of model from statistics of dataset

Table 1. Analogies between standard image classification terms and their AutoEval equivalents. The analogy shows that the image classification is an image-based task, while the AutoEval problem in this work is dataset-based.

metric [15]. In addition, domain adaptation literature shows that a smaller distribution difference leads to higher target domain accuracy and implies that a large domain gap causes a low test accuracy [14, 37, 38].

In this work, we explicitly show that there is a very strong negative correlation between accuracy and distribution difference (the Spearman’s Rank Correlation [35] is  $-0.9$ ). This observation indicates that it is feasible to estimate classifier accuracy with distribution statistics. With this, we attempt to quantitatively estimate the test accuracy by studying the underlying relationship between dataset distribution and classifier performance. We propose to learn this relationship via a meta-dataset (dataset of datasets). We use the terms meta set and meta-dataset interchangeably. Unlike most existing datasets that treat each image as a sample, we focus on the dataset level: in the meta-dataset, each dataset is treated as a sample, which we term “sample set”. The analogy between standard image classification and Au-

toEval task is shown in Table 1. The sample sets should possess an appropriate number of images, exhibit a diverse spread of distributions, and in the case of image classification, have the same set of classes.

It is difficult to collect sufficient real-world sample sets that meet the above mentioned requirements. In this work, we propose to construct the meta set by data synthesis. Every sample set in the meta set is generated from a seed set that follows the same distribution as the original training set. This is achieved via various geometric and photo-metric transformation operations on the seed set, including blurring, background substitution, foreground rotation, translation, *etc.* Note that, the synthetic sample sets are fully labeled because they are transformed versions of the seed set. Using these labels, we can obtain the recognition accuracy of the classifier on each sample set. Sample set  $i$  can thus be denoted by  $(\mathbf{f}_i, a_i)$ , where  $a_i$  is a recognition accuracy, and  $\mathbf{f}_i$  is the vector representation of the dataset, *e.g.*, the mean vector of image features in this dataset. With this meta set denoted as  $\{(\mathbf{f}_i, a_i)\}, i = 1, \dots, N$ , where  $N$  is the number of sample sets, we can train a regression model that takes input as the  $\mathbf{f}$  of a sample set and returns the predicted classifier accuracy on this set.

In summary, the main contributions of this paper include:

- We introduce the AutoEval task, aiming to estimate the recognition accuracy of a trained classifier on a test set *without* any human annotated label.
- We validate the feasibility of estimating classifier accuracy from dataset-level feature statistics. With this, we propose to learn an accuracy regression model from a synthetic meta-dataset (a dataset comprised of many datasets) and obtain promising accuracy predictions for real-world test datasets.

## 2. Automatic Model Evaluation

We are interested in predicting the recognition accuracy of a trained classifier on an unlabeled test set.

### 2.1. Problem Definition

We first define a labeled dataset,  $D^l = \{(\mathbf{x}_i, y_i)\}$  where  $i \in [1, \dots, M]$ ,  $\mathbf{x}_i$  is an image,  $y_i$  is its class label, and  $M$  is the number of images. Consider a source domain  $\mathcal{S}$ , from which we sample an original training dataset  $\mathcal{D}_{ori}$ . We use  $\mathcal{D}_{ori}$  to train a classifier  $f_\theta : \mathbf{x}_i \rightarrow \hat{y}_i$ , which is parameterized by  $\theta$  and maps an image  $\mathbf{x}_i$  to its predicted class  $\hat{y}_i$ . Given  $D^l$ , we obtain its classification accuracy by comparing the class predictions  $\hat{y}_i$  with the ground truth  $y_i$  to obtain accuracy,

$$a_{standard} = \frac{\sum_{i=1}^M \mathbb{I}[\hat{y}_i == y_i]}{M}, \quad (1)$$

where  $\mathbb{I}[\cdot]$  is an indicator function returning 1 if argument is true and 0 otherwise.

In AutoEval, given  $f_\theta$  and an unlabeled dataset  $\mathcal{D}^u = \{\mathbf{x}_i\}$  for  $i \in [1, \dots, M]$ , we use an accuracy predictor  $A : (f_\theta, \mathcal{D}^u) \rightarrow a$ , which outputs an estimated classifier accuracy  $a \in [0, 1]$  on this test set,

$$a_{auto} = A(f_\theta, \mathcal{D}^u). \quad (2)$$

Note that in image classification,  $\mathcal{D}_{ori}$  and  $\mathcal{D}^u$  share the same label space.

### 2.2. An Intuitive Solution

We first present an intuitive solution to the AutoEval problem, which is not learning based. This solution is motivated by the pseudo labeling strategy in many vision tasks [17, 41, 28]. The basic assumption is: if a class prediction is made with a high softmax output score, this prediction is likely to be correct. Formally, let us consider a  $K$ -way classification problem. When feeding a test image  $x_i$  to a trained classifier  $f_\theta$ , we obtain  $\mathbf{s}_i \in \mathbb{R}^K$ , which is the output of the softmax layer. The  $k$ -th entry in  $\mathbf{s}_i$  characterizes the probability of  $x_i$  belonging to class  $k$ . The  $\ell_1$  norm  $\|\mathbf{s}_i\|_1 = 1$ . If the maximum entry of  $\mathbf{s}_i$  is greater than a threshold  $\tau$ , image  $x_i$  is considered to be correctly classified. The accuracy predictor is written as,

$$a_{max} = A_{max}(f_\theta, \mathcal{D}^u) = \frac{\sum_{i=1}^M \mathbb{I}[\max(\mathbf{s}_i) > \tau]}{M}, \quad (3)$$

where  $M$  is the number of images in  $\mathcal{D}^u$ . We will evaluate  $A_{max}$  in the experiment and show that it does not work consistently well across datasets.

## 3. Methods

### 3.1. Formulation

Motivated by the implications in domain adaptation, we propose to address AutoEval by measuring the distribution difference between the original training set and the test set, and explicitly learning a mapping function from the distribution shift to the classifier accuracy.

Under this consideration, we formulate AutoEval as a dataset-level regression problem. In this problem, we view a dataset as a sample, and its label is the recognition accuracy on the dataset itself. Given  $N$  sample sets, we denote the  $j$ -th sample set  $\mathcal{D}_j$  as  $(\mathbf{f}_j, a_j)$ , where  $\mathbf{f}_j$  is the representations for  $\mathcal{D}_j$ , and  $a_j \in [0, 1]$  is the recognition accuracy of classifier  $f_\theta$  on  $\mathcal{D}_j$ . We aim to learn a regression model (accuracy predictor), written as,

$$a_j = A(\mathbf{f}_j). \quad (4)$$

We use a standard squared loss function for this model,

$$\mathcal{L} = \frac{1}{N} \sum_{j=1}^N (\hat{a}_j - a_j)^2, \quad (5)$$

where  $\hat{a}_j$  is the predicted accuracy of the  $j$ -th sample set  $\mathcal{D}_j$ , and  $a_j$  is the ground truth classifier accuracy of  $\mathcal{D}_j$ .

During testing, we extract the dataset representation  $\mathbf{f}^u$  for unlabeled test set  $\mathcal{D}^u$ , and obtain estimated classification accuracy using  $a = A(\mathbf{f}^u)$ .

To learn regression models defined in Eq. 4 and Eq. 5, we need to specify the design of 1) dataset representation  $\mathbf{f}_j$ , 2) regression model  $A$ , and 3)  $N$  sample sets (meta-dataset).

### 3.2. Regression Model and Dataset Representation

**Linear regression.** We first introduce a simple linear regression model,

$$a_{linear} = A_{linear}(\mathbf{f}) = w_1 f_{linear} + w_0, \quad (6)$$

where  $f_{linear} \in \mathbb{R}$  is the representation of sample set  $\mathcal{D}$ , and  $w_0, w_1 \in \mathbb{R}$  are parameters of this linear regression model. Based on the intuition that the domain gap impacts classifier accuracy, we define  $f_{linear}$  as the quantified domain gap between dataset  $\mathcal{D}$  and the original training set  $\mathcal{D}_{ori}$ . Specifically, we use the Fréchet distance [12] to measure the domain gap, and thus,

$$f_{linear} = \text{FD}(\mathcal{D}_{ori}, \mathcal{D}) = \|\boldsymbol{\mu}_{ori} - \boldsymbol{\mu}\|_2^2 + \text{Tr}(\boldsymbol{\Sigma}_{ori} + \boldsymbol{\Sigma} - 2(\boldsymbol{\Sigma}_{ori}\boldsymbol{\Sigma})^{\frac{1}{2}}), \quad (7)$$

where  $\boldsymbol{\mu}_{ori}$  and  $\boldsymbol{\mu}$  are the mean feature vectors of  $\mathcal{D}_{ori}$  and  $\mathcal{D}$ , respectively.  $\boldsymbol{\Sigma}_{ori}$  and  $\boldsymbol{\Sigma}$  are the covariance matrices of  $\mathcal{D}_{ori}$  and  $\mathcal{D}$ , respectively. They are calculated from the image features in  $\mathcal{D}_{ori}$  and  $\mathcal{D}$ , which are extracted using the classifier  $f_\theta$  trained on  $\mathcal{D}_{ori}$ . Other measurements of the domain gap can also be used, such as MMD [15].

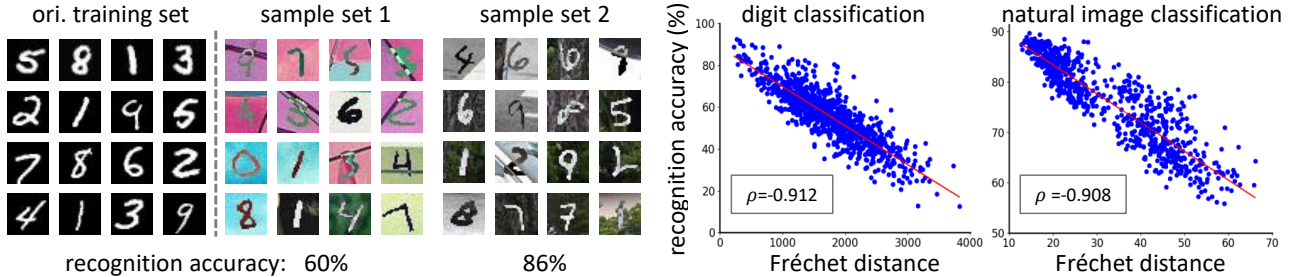


Figure 2. Relationship between the distribution shift and classifier accuracy on digits and natural image classification. Each point represents a sample set of the meta set. The Spearman’s Rank Correlation ( $\rho$ ) [35] between distribution shift and classifier accuracy is around  $-0.91$  on two scenarios, indicating they have a very strong negative correlation. The red straight line is fit with robust linear regression [18].

**Proof of concept.** Given a meta set and a classifier trained on the training dataset  $\mathcal{D}_{ori}$  from a source domain  $\mathcal{S}$ , we study the relationship between classifier’s accuracy and distribution shift. In Fig. 2, we show the accuracy as a function of the distribution shift. The distribution shift is measured by Fréchet distance (FD) with the features extracted from the trained classifier. In practice, we use the activations in the penultimate of the classifier as features.

In both digits and natural image classification, we observe a very strong negative correlation between accuracy and distribution shift in both digits and natural image classification: the Spearman’s Rank Correlation ( $\rho$ ) [35] is about  $-0.91$ . That is, the classifier tends to achieve a low accuracy on the sample set which has a high distribution shift with training set  $\mathcal{D}_{ori}$ . This indicates it is feasible to learn a regression model to predict classifier performance based on distribution difference between training and test sets.

**Neural network regression.** Besides the linear regression, we also propose a neural network regression model,  $a_{neural} = A_{neural}(f_{neural})$ , which has the same formulation as Eq. 4. In practice, we use a simple fully connected neural network for regression. The input of the model is the dataset representation  $f_{neural}$ , and the output is the estimated classifier accuracy  $a_{neural}$ .

With the observation in the proof of concept, we propose to use distribution-related statistics to represent a dataset. In this work, we use its first-order and second-order feature statistics, *i.e.*, mean vector and covariance matrix. Moreover, we also include a 1-dim FD score as an auxiliary information to the representation. Compared with linear regression, the neural network regression has a richer dataset representation. The dataset representation is written as,

$$f_{neural} = [f_{linear}; \boldsymbol{\mu}; \boldsymbol{\Sigma}], \quad (8)$$

where  $f_{linear} \in \mathbb{R}$  is the Fréchet distance between  $\mathcal{D}$  and  $\mathcal{D}_{ori}$ ,  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  are calculate the same way as Eq. 7. Covariance  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  is very high-dimensional, making training difficult. Dimension reduction is thus necessary. Specifically, we calculate  $\boldsymbol{\sigma}$  by taking a weighted summation of



Figure 3. Visual examples of transformations. Here we show autoContrast, rotation, translation, brightness, and color. For other used transformations, we refer readers to [6].

each row of  $\boldsymbol{\Sigma}$  to produce a single vector, using learned column specific coefficients that are shared across all rows. For example, if the feature extracted from  $f_{\theta}$  is  $d$ -dim, the dimensionality of  $f_{neural}$  is  $1 + 2d$ .

### 3.3. Constructing Training Meta-dataset

**Meta-datasets for training.** The regression model (Eq. 4, Eq. 5, Eq. 8) takes the dataset representation as input and outputs a classification accuracy. To train it, we need to prepare a meta-dataset in which each sample is a dataset. In classification, the diversity of the samples in the training set should ideally be sufficient such that test scenario is represented in its distribution. In this work, we seek to create a diverse meta set that (hopefully) contains the test distributions. To construct such a meta set, we should collect sample sets that are 1) large in number, 2) diverse in the data distribution, and 3) have the same label space with the training set. There are very few real-world datasets that satisfy these requirements, so we resort to data synthesis.

For each classification task (digits or natural images), we synthesize sample sets from a single seed dataset. The

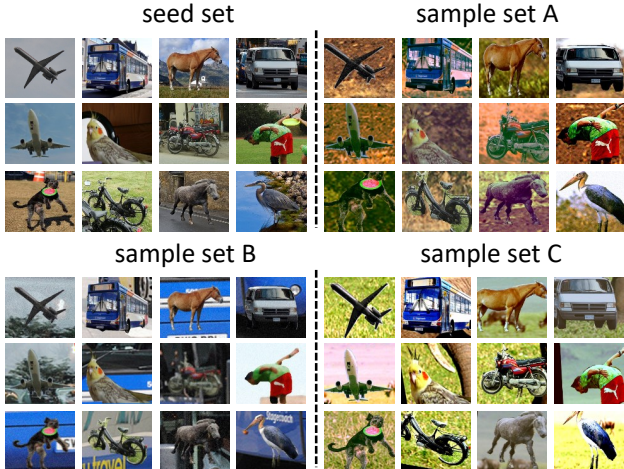


Figure 4. The seed set and examples of three sample sets. The seed set is from the same distribution with the original training set; they share the same classes but do not have image overlap. The sample sets are generated from the seed by background replacement and image transformations. The sample sets exhibit distinct data distributions, but inherit the foreground objects from the seed, and thus are fully labeled. Many sample sets form a meta-dataset from which an accuracy regression model is trained.

seed  $\mathcal{D}_s$  is sampled from source domain  $\mathcal{S}$ , and thus has the same distribution as  $\mathcal{D}_{ori}$ . Given  $\mathcal{D}_s$ , we apply various visual transformation and obtain  $N$  different sample sets  $\mathcal{D}_j, j = 1, \dots, N$ . Since  $\mathcal{D}_s$  is fully labeled, these sample sets inherit the labels from  $\mathcal{D}_s$ .

To create a sample set  $\mathcal{D}_j$ , we adopt a two-step procedure: perform background change, and then image transformations. In the first step, we keep the foreground / object unchanged and replace the background. For each sample set, we randomly select an image from the COCO dataset [26], from which we randomly crop a patch and use it as the background. The patch scale and position in that image are both random. In the second step, for the background-replaced images, we use six image transformations defined in [6], including autoContrast, rotation, color, brightness, sharpness, and translation. Examples of some transformations are shown in Fig. 3. For each sample set, we randomly select and combine three out of the six transformations, with the magnitude of each transformation being random on per-sample basis. As such, each sample set is generated by background replacement and a combination of three image transformations. Fig. 4 presents examples of sample sets in natural image classification, where background replacement can be observed. In the supplementary materials, we present the detailed transformation parameters and more visual examples of the training meta set. Note that a sample set inherits all the image labels from the seed set and is fully labeled. As such, we can calculate the recognition accuracy



Figure 5. Sample images from real-world test datasets, including SVNH, USPS, Pascal, Caltech and ImageNet. The former two are for digit classification, and the latter three are for natural image classification. We predict the classifier accuracy on these datasets.

of classifier  $f_\theta$  on each sample set. Sample set  $\mathcal{D}_j$  can be denoted as  $(f_j, a_j)$ , which is used as a training sample to optimize the regression model.

**Real-world datasets for testing.** This is an early attempt for the AutoEval problem. To our knowledge, we could only find few real-world datasets that have different distributions but contain the same classes. To clarify the AutoEval problem, we conduct extensive analyses with these dataset. For digits classification, we use USPS [19] and SVHN [30], both with 10 classes. For natural image classification, we use three existing datasets, *i.e.*, PASCAL [13], Caltech [16], and ImageNet [7], all with 12 classes. Details of the test meta sets are provided in Sec. 4.1.

## 4. Experiment and Analysis

### 4.1. Experimental Settings

We study the AutoEval problem on two classification tasks: digit classification and natural image classification.

**Digit classification.** The original training set contains all the training images of MNIST. We use the testing images of MNIST as the seed to generate the training meta set. Because MNIST images are binary, the foreground can be separated from the background. When generating a meta set, we randomly select an image from the COCO training set, and the background of each image is replaced with a random patch of the sampled COCO image. Then, we apply three out of six image transformations to images. We generate 3,000 sample sets, of which we use 2,000 and 1,000 for the training and the validation meta set, respectively. Moreover, we use two real datasets for testing, *i.e.*, USPS [19] and SVHN [30] datasets.

**Natural image classification.** We use COCO [26] training set as the original training set, and COCO validation set as the seed set to build meta set. When generating meta set for training, we use instance mask annotations of COCO validation set to get foreground regions. Similar to digit classification, for each sample set, we replace the background with a random patch of an image from COCO test set. We then use image transformations to introduce more visual changes. We create 1,600 sample sets from the seed set, of which we use 1,000 and 600 for the training and the val-

Train Set	Digits			Natural images			
	SVHN	USPS	RMSE↓	Pascal	Caltech	ImageNet	RMSE↓
Unseen Test Set	SVHN	USPS	RMSE↓	Pascal	Caltech	ImageNet	RMSE↓
Ground-truth accuracy	25.46	64.08	-	86.13	93.40	88.83	-
Predicted score ( $\tau = 0.7$ )	10.09	43.60	18.11	88.34	93.28	90.17	1.49
Predicted score ( $\tau = 0.8$ )	7.97	37.22	22.66	84.32	90.78	86.50	2.28
Predicted score ( $\tau = 0.9$ )	7.03	32.94	25.59	78.61	87.71	81.33	6.96
Linear reg.	26.28	50.14	9.87	83.87	79.77	83.19	8.62
Neural network reg.	27.52	64.11	1.46	87.76	89.39	91.82	3.04

Table 2. Method comparison in predicting classification accuracy. Results on digit classification (SVHN and USPS datasets) and natural image classification (Pascal, Caltech, and ImageNet) are shown. We compare three methods, *i.e.*, predicted score based (Section 2.2), linear regression and neural network regression (Section 3.2). For each dataset, we report the estimated classification accuracy (%). For both digit and natural image classification, RMSE (%) is reported. The original training sets are MNIST and COCO, respectively. The ground-truth recognition accuracy (%) is presented.

ification meta set, respectively. In testing, we use PASCAL [13], Caltech [16], and ImageNet [7]. For each dataset, we select images of 12 common classes, *i.e.*, aeroplane, bike, bird, boat, bottle, bus, car, dog, horse, monitor, motorbike, and person. We reduce the “person” class to 600 images to balance the overall number of images per class.

**Classifier architecture.** For digit classification, we use LeNet-5 [23] as classifier. Since the all images are mapped to the RGB space, we modify the number of input channel of LeNet-5 to 3. For natural image classification, we use the ResNet-50 pretrained on ImageNet [7] which is adapted to the 12-way classification.

**Metrics.** This paper estimates the recognition accuracy of a model on a test set. To evaluate the performance of such estimate, we use root mean squared error (RMSE) and mean absolute error (MAE) as metrics. RMSE measures the average squared difference between the estimated classifier accuracy and ground-truth accuracy. MAE measures the average magnitude of the errors. Small RMSE and MAE correspond to good predictions and vice versa.

## 4.2. Classifier Accuracy Prediction

This paper introduces three possible methods to estimate the recognition accuracy, including the confidence-based method, linear regression and neural network regression. We report the estimations of these methods in Table 2. For the predicted score based method, three thresholds (*i.e.*,  $\tau = 0.7, 0.8$  and  $0.9$  in Eq. 3) are used.

**The predicted score based method is very sensitive to the threshold.** Under a specific threshold ( $\tau = 0.7$ ), this method makes accuracy prediction on natural image datasets (RMSE=1.49%), but its prediction quality drops significantly (from 1.49% to 6.96%) when we increase value of  $\tau$  to 0.9. What is more, its performance is very poor when considering the digit classification task. Under two values of  $\tau$ , the RMSE is consistently high, *i.e.*, 22.66% and 25.59%, respectively. Note that, it is infeasible to select the optimal threshold because 1) test labels are unavailable and

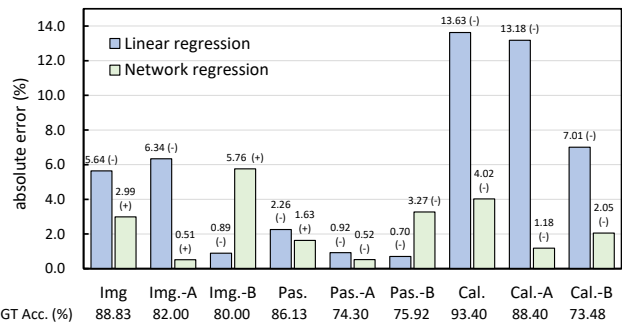


Figure 6. Comparing linear regression and neural network regression when test data undergo new image transformations such as Cutout [9, 43], Shear, Equalize and ColorTemperature [6]. The transformed datasets are denoted by “-A” and “-B”. We report the absolute error (%) of predictions and the ground truth accuracy is also shown below each dataset. (-) / (+) means the predicted accuracy is lower / higher than the ground-truth accuracy, respectively.

2) the test domain keeps changing. Our method does not depend on such a hyper-parameter and yields much more stable results. That said, it would be interesting to address this drawback in the context of AutoEval.

**Regression methods achieve better predictions than predicted score based method.** In digit datasets, the RMSE values of linear regression and neural network regression are 9.87% and 1.46%, respectively. A similar trend can be observed in natural image datasets. Their RMSE scores are generally lower and more stable than the predicted-based method. This indicates the effectiveness of learning-based methods: the distribution difference between the original training and test sets is a critical feature.

**Neural network regression is generally better than linear regression.** As shown in Table 2, the neural network regression is more accurate than linear regression in both digit and natural image datasets. For example, RMSE of the former is 8.41% lower than the latter on digit datasets. In fact, the RMSE of neural network regression is as small as 1.46%: the predicted classifier accuracy is very close to the ground truth accuracy.

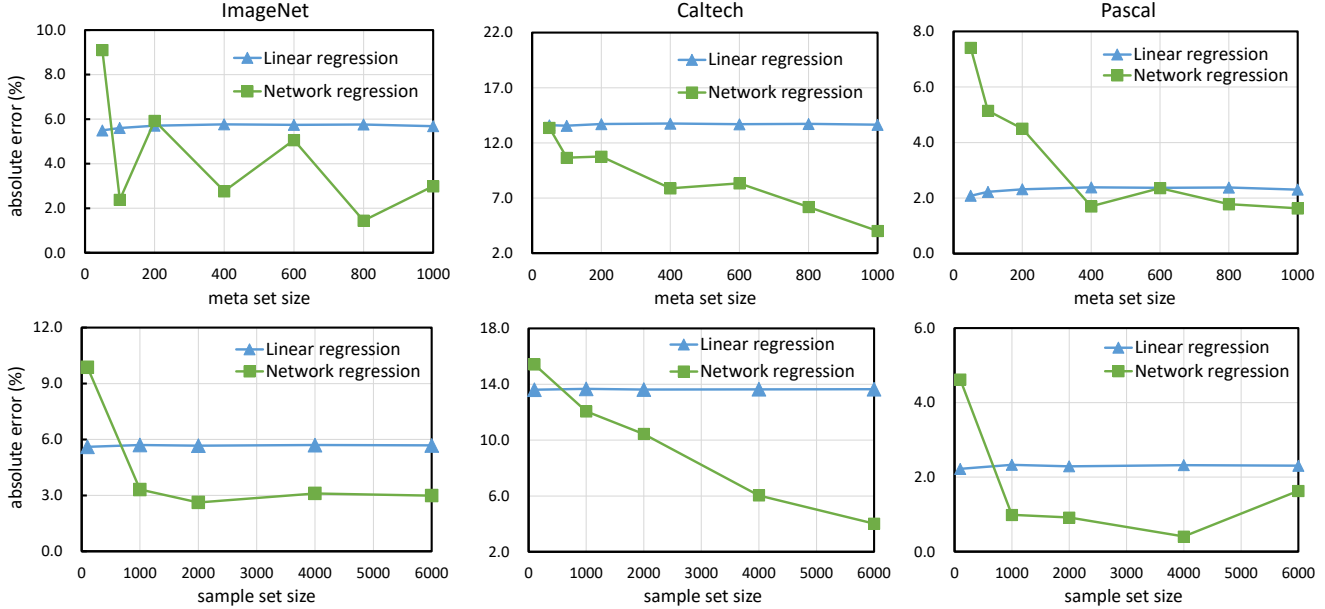


Figure 7. The impact of meta set size (first row) and sample set size (second row) on the performance of regression methods. We report the absolute errors (%) between estimated classifier accuracy and the ground-truth accuracy. We observe that linear regression is relatively stable with different sample set and meta set size. In comparison, neural network needs more and larger sample sets for training.

We note that linear regression is significantly inferior to neural network regression on Caltech datasets, where linear regression gives errors higher than 10%. Caltech is an interesting dataset. Its images have relatively simple backgrounds and salient foregrounds, implying that they are “easy” to classify. However, such simple background contrasts significantly with the original training set (COCO), so the FD score between Caltech and COCO is very large. Only looking at the FD score, linear regression tends to predict low accuracy on Caltech. In comparison, neural network regression considers the data statistics of Caltech, such that it can make more accurate prediction. Furthermore, the meta-dataset might already contain sample sets with such “simple backgrounds” (large FD), and high recognition accuracy. Under such circumstances, the network has learned to overrule the large FD and instead resort to the “simple background” when making predictions.

**Robustness of regression models.** To further examine the two regression methods, we perform image transformations to real datasets (ImageNet, Pascal and Caltech) and assess the performance of the two regression methods on these “edited real-world datasets”. Note that image transformations we use here are *not* applied in meta-dataset generation. Thus, this experiment assesses some generalization ability of the regression methods. From Fig. 6, we first observe the ground truth accuracy on the edited datasets is lower than that on the original sets. It suggests that the image transformations are introducing visual differences that hinder the classifier performance. The results show that the two regres-

sion methods could also achieve reasonably good estimated results. For example, linear regression makes promising predictions on 6 out of 9 datasets (It has the same issue discussed above on the 3 Caltech sets). Our network regression model gives lower errors on all 9 datasets. This suggests that our network can learn from diverse and various sets of the meta-set to make an accurate performance prediction.

### 4.3. Analysis of the Training Meta-Dataset

The synthetic meta-dataset is a key component of our system, allowing us to obtain labeled samples sets in a large scale. We analyze its impact on the regression methods from two aspects, *i.e.*, meta set size and sample set size.

**Meta set size.** We first study the impact of meta set size on the regression methods. Meta set contains training samples/datasets for regression models. In Fig. 7 (first row). We observe the results of linear regression are relatively stable with different meta set sizes. It can achieve good performance even with 50 sample sets. This is because linear regression only has two parameters (Eq. 6), which can be learned with relatively few samples [18]. In comparison, neural network cannot achieve good results when the number of sample sets is small. When provided adequate sample sets, the neural network can learn effectively from rich and diverse sample datasets and surpasses the linear regression.

**Sample set size.** By default, the number of images in each sample set is equal to that of seed  $\mathcal{D}_s$ . We study the impact of sample set size on the regression methods. In the experiment, we set the meta set size 1000, and vary the sample set

size. In Fig. 7 (second row), we observe linear regression is stable under different sample set sizes. In comparison, the neural network needs more images in each sample set for training. We think more images in each sample set make the distribute-related representations more accurate. This is beneficial for regression learning of network.

## 5. Related Work

**Model generalization prediction.** There are some works develop complexity measurements on training sets and model parameters to predict generalization error [20, 2, 4, 20, 31, 39]. Corneanu *et al.* [4] use the persistent topology measures to predict the performance gap between training and testing error, even without the need of any testing samples. Jiang *et al.* [20] introduce a measurement of layer-wise margin distributions for generalization ability. Neyshabur *et al.* [31] develop bounds on the generalization gap based on the product of norms of the weights across layers. Moreover, the agreement score of several classifiers' predictions can be used for estimation [29, 34, 33, 11, 21]. Our work differs significantly: we focus on the measuring statistics related to test sets for prediction.

**Out-of-distribution (OoD) detection.** This task [10, 17, 24, 40, 25] considers the distribution of test samples. Specifically, this task aims to detect test samples that follow a distribution different from the training distribution. This has been studied from different views, such as anomaly detection [1], open-set recognition [3], and rejection [5]. For example, Hendrycks *et al.* [17] use probabilities output from a softmax classifier as indicator to find out-of-distribution samples. While this task attempts to detect abnormal testing samples, our work considers the overall statistics of all test samples to predict classifier accuracy.

**Unsupervised Domain adaptation.** Our work also relates to unsupervised domain adaptation, where the goal is to use labeled source samples and unlabeled target samples to learn a model that can generalize well on the target dataset [27, 38, 44, 8]. Many moment matching schemes have been studied for this task [36, 27, 38, 32, 36, 42]. Long *et al.* [27] and Tzeng *et al.* [38] utilize the maximum mean discrepancy (MMD) metric [15] to learn a shared feature representation. In this work, we study the underlying relationship between the model performance and the distribution shift. By leveraging dataset level statistics, we are able to accurately predict model performance on unlabeled test sets.

## 6. Conclusions and Perspectives

This paper investigates the problem of predicting classifier accuracy on test sets without ground-truth labels. It has the potential to yield significant practical value, such as predicting system failure in unseen real-world environments. Importantly, this task requires us to derive similarities and

representations on the dataset level, which is significantly different from common image-level problems. We make some tentative attempts by devising two regression models which directly estimate classifier accuracy based on overall distribution statistics. We build a dataset of datasets (meta-dataset) to train the regression model. We show that the synthetic meta-dataset can cover a good range of data distributions and benefit AutoEval on real-world test sets. For the remainder of this section, we discuss the limitations, potential, and interesting aspects of AutoEval.

**Application scope.** Our system assumes that variations in the real-world cases can be approximated by the image transformations in the training meta set. With various and diverse sample sets, our system learns to make promising predictions for novel environments. However, if the test datasets exhibit some very special patterns or conditions, the system might not be able to work. An example is that the test dataset has an entirely different set of classes, and this test distribution cannot be approximated by the meta-dataset in our work. Under this circumstance, our trained models will still give an estimated accuracy, which is clearly incorrect. On a related extreme case, the test dataset might only contain ambiguous and adversarial samples, meaning that the test accuracy could be as poor as random. Such cases are not included in meta-dataset, either. Potentially, the above two issues could be alleviated by including such cases into the meta-dataset with a specific dataset design. Another option is to use out-of-distribution detection techniques to help detect and reject such cases.

**Dataset Representation.** Our work relates to an interesting research problem: how to represent a dataset? This problem is more challenging than describing a single image because a dataset contains much more information. This work uses distribution-related feature statistics (mean and covariance) to characterize a classification dataset. We believe there are other potential representations for better representing a dataset. On the other hand, it would be interesting to study the representation in other tasks (*e.g.*, object detection and semantic segmentation), where global feature statistics might not be suitable to characterize a dataset.

**Similarities between datasets.** We measure dataset similarity using the FD score. However, this problem is as challenging as dataset representation, especially when we aim to connect the similarity with test accuracy. This problem will benefit the domain adaptation field, where more precise domain bias measurement and its connection to target set accuracy will significantly help algorithm design.

## Acknowledge

This work was supported by the ARC Discovery Early Career Researcher Award (DE200101283) and the ARC Discovery Project (DP210102801). We thank all anonymous reviewers and ACs for their constructive comments.



## References

- [1] J Andrews, Thomas Tanay, Edward J Morton, and Lewis D Griffin. Transfer representation-learning for anomaly detection. *Journal of Machine Learning Research*, 2016. 8
- [2] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. *arXiv preprint arXiv:1802.05296*, 2018. 8
- [3] Abhijit Bendale and Terrance Boulton. Towards open world recognition. In *Proc. CVPR*, pages 1893–1902, 2015. 8
- [4] Ciprian A Corneanu, Sergio Escalera, and Aleix M Martinez. Computing the testing error without a testing set. In *Proc. CVPR*, pages 2677–2685, 2020. 8
- [5] Corinna Cortes, Giulia DeSalvo, and Mehryar Mohri. Learning with rejection. In *International Conference on Algorithmic Learning Theory*, pages 67–82. Springer, 2016. 8
- [6] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Proc. CVPR*, pages 113–123, 2019. 4, 5, 6
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proc. CVPR*, pages 248–255, 2009. 1, 5, 6
- [8] Weijian Deng, Liang Zheng, Qixiang Ye, Yi Yang, and Jianbin Jiao. Similarity-preserving image-image domain adaptation for person re-identification. *arXiv preprint arXiv:1811.10551*, 2018. 8
- [9] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 6
- [10] Terrance DeVries and Graham W Taylor. Learning confidence for out-of-distribution detection in neural networks. *arXiv preprint arXiv:1802.04865*, 2018. 8
- [11] Pinar Donmez, Guy Lebanon, and Krishnakumar Balasubramanian. Unsupervised supervised learning i: Estimating classification and regression errors without labels. *Journal of Machine Learning Research*, 11(4), 2010. 8
- [12] DC Dowson and BV Landau. The fréchet distance between multivariate normal distributions. *Journal of multivariate analysis*, 12(3):450–455, 1982. 1, 3
- [13] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge 2007 (voc2007) results. 2007. 5, 6
- [14] Yaroslav Ganin and Victor S. Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proc. ICML*, pages 1180–1189, 2015. 2
- [15] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander J. Smola. A kernel method for the two-sample-problem. In *Proc. NeurIPS*, pages 513–520, 2006. 1, 2, 3, 8
- [16] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. Technical Report 7694, California Institute of Technology, 2007. 5, 6
- [17] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016. 3, 8
- [18] Peter J Huber. *Robust statistics*, volume 523. John Wiley & Sons, 2004. 4, 7
- [19] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994. 5
- [20] Yiding Jiang, Dilip Krishnan, Hossein Mobahi, and Samy Bengio. Predicting the generalization gap in deep networks with margin distributions. In *Proc. ICLR*, 2019. 8
- [21] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *ICLR*, 2020. 8
- [22] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *Proc. CVPR*, pages 9404–9413, 2019. 1
- [23] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 1, 6
- [24] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *Proc. NeurIPS*, pages 7167–7177, 2018. 8
- [25] KIMIN LEE, Kibok Lee, Honglak Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. In *Proc. ICLR*, 2018. 8
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proc. ECCV*, pages 740–755. Springer, 2014. 1, 5
- [27] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. Learning transferable features with deep adaptation networks. In *Proc. ICML*, pages 97–105, 2015. 8
- [28] Fan Ma, Deyu Meng, Qi Xie, Zina Li, and Xuanyi Dong. Self-paced co-training. In *ICML*, pages 2275–2284. JMLR.org, 2017. 3
- [29] Omid Madani, David Pennock, and Gary Flake. Co-validation: Using model disagreement on unlabeled data to validate classification algorithms. In *Proc. NeurIPS*, volume 17, pages 873–880, 2004. 8
- [30] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisaccho, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 5
- [31] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. In *Proc. NeurIPS*, pages 5947–5956, 2017. 8
- [32] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proc. ICCV*, pages 1406–1415, 2019. 1, 8
- [33] Emmanouil Platanios, Hoifung Poon, Tom M Mitchell, and Eric J Horvitz. Estimating accuracy from unlabeled data: A probabilistic logic approach. In *Proc. NeurIPS*, pages 4361–4370, 2017. 8
- [34] Emmanouil Antonios Platanios, Avinava Dubey, and Tom Mitchell. Estimating accuracy from unlabeled data: A bayesian approach. In *Proc. ICML*, pages 1416–1425, 2016. 8
- [35] Charles Spearman. The proof and measurement of association between two things. 1961. 2, 4

- [36] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. In *Proc. AAAI*, 2016. 1, 8
- [37] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proc. CVPR*, pages 2962–2971, 2017. 2
- [38] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. 2, 8
- [39] Thomas Unterthiner, Daniel Keysers, Sylvain Gelly, Olivier Bousquet, and Ilya Tolstikhin. Predicting neural network accuracy from weights. *arXiv preprint arXiv:2002.11448*, 2020. 8
- [40] Apoorv Vyas, Nataraj Jammalamadaka, Xia Zhu, Dipankar Das, Bharat Kaul, and Theodore L Willke. Out-of-distribution detection using an ensemble of self supervised leave-out classifiers. In *Proc. ECCV*, pages 550–564, 2018. 8
- [41] Weichen Zhang, Wanli Ouyang, Wen Li, and Dong Xu. Collaborative and adversarial network for unsupervised domain adaptation. In *Proc. CVPR*, pages 3801–3809, 2018. 3
- [42] Zhen Zhang, Mianzhi Wang, Yan Huang, and Arye Nehorai. Aligning infinite-dimensional covariance matrices in reproducing kernel hilbert spaces for domain adaptation. In *Proc. CVPR*, pages 3437–3445, 2018. 8
- [43] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. *arXiv preprint arXiv:1708.04896*, 2017. 6
- [44] Yiming Zuo, Weichao Qiu, Lingxi Xie, Fangwei Zhong, Yizhou Wang, and Alan L Yuille. Craves: Controlling robotic arm with a vision-based economic system. In *Proc. CVPR*, pages 4214–4223, 2019. 8