

Operation is the hardest teacher: estimating DNN accuracy looking for mispredictions

Antonio Guerriero, Roberto Pietrantuono, Stefano Russo
DIETI, Università degli Studi di Napoli Federico II
Via Claudio 21, 80125 - Napoli, Italy
{antonio.guerriero, roberto.pietrantuono, stefano.russo}@unina.it

Abstract—Deep Neural Networks (DNN) are typically tested for accuracy relying on a set of unlabelled real world data (operational dataset), from which a subset is selected, manually labelled and used as test suite. This subset is required to be small (due to manual labelling cost) yet to faithfully represent the operational context, with the resulting test suite containing roughly the same proportion of examples causing misprediction (i.e., failing test cases) as the operational dataset.

However, while testing to estimate accuracy, it is desirable to also learn as much as possible from the failing tests in the operational dataset, since they inform about possible bugs of the DNN. A smart sampling strategy may allow to intentionally include in the test suite many examples causing misprediction, thus providing this way more valuable inputs for DNN improvement while preserving the ability to get trustworthy unbiased estimates.

This paper presents a test selection technique (DeepEST) that actively looks for failing test cases in the operational dataset of a DNN, with the goal of assessing the DNN expected accuracy by a small and “informative” test suite (namely with a high number of mispredictions) for subsequent DNN improvement. Experiments with five subjects, combining four DNN models and three datasets, are described. The results show that DeepEST provides DNN accuracy estimates with precision close to (and often better than) those of existing sampling-based DNN testing techniques, while detecting from 5 to 30 times more mispredictions, with the same test suite size.

Index Terms—Software testing, Artificial neural networks

I. INTRODUCTION

Deep Neural Networks (DNN) are today integral part of many applications, including safety-critical software systems such as in the medical [1] and autonomous driving domains [2]. Testing is a crucial activity in the development of such systems, for both quality/safety reasons to avoid DNN-caused catastrophic failures [3], [4], and for cost as well – very large samples may be needed to reliably test a DNN [5], [6].

A significant research effort is currently being put on DNN testing. A primary goal is to find adversarial examples causing *mispredictions*, namely to expose as many failing behaviours as possible [7]–[11]. Several structural coverage criteria have been proposed to drive the automated generation of test inputs and assess the test adequacy – neuron coverage [7], k-multisection neuron coverage, neuron boundary coverage [8], combinatorial coverage [10]. It has been argued that these criteria may be misleading, because of the low correlation between the number of misclassified inputs in a test set and its

coverage [12]. Wu *et al.* [13] and Kim *et al.* [14] recently considered discrepancy measures between the training/validation data and test data, to improve fault detection and to have coverage criteria better correlated to failure-inducing inputs.

The output of this type of failure-finding testing (and then debugging) process is an improved DNN, with higher accuracy. This resembles what is called *debug testing* in the traditional testing literature [15]. Beside differences in testing DNN-based and conventional software (e.g., the oracle definition), which make DNN testing problematic, a further issue is that the so-obtained testing results are not necessarily related to the accuracy actually experienced in operation. In fact, *testing data may be not representative of the actual operational context*. This may happen when test data are generated artificially (like in adversarial examples generation) or they differ significantly from input observed in the field. The number of mispredictions and/or the coverage achieved give only an “indirect” (and, for what said, inaccurate) measure of the expected accuracy in operation, and ultimately of the confidence that can be placed in the DNN.

A less investigated research path is testing a DNN with the explicit goal of providing a statistical estimate of its *expected accuracy* in the operational context. In software reliability engineering, this is a well established practice known as *operational testing* [16]. The objective is to assess how well a DNN will perform in the intended context using a small yet effective amount of test data. With a cost-effective accuracy estimate, testers can establish a threshold-based release criterion and correct or tune their artificial networks (e.g., adjust the DNN structure or hyper-parameters) until the criterion is met. The reference scenario is the following: a DNN model is trained with a set of data (*training dataset*) and is meant to operate in a given context (*operational context*). To test the model, an arbitrarily large set of operational data is available (*operational dataset*) containing examples whose correct label is unknown: the goal of the tester is to select a small subset of the operational data, to be manually labelled and used as test cases to estimate the accuracy of the model in operation. Due to manual labelling, testers typically look for a trade-off between a good estimate and the labelling cost. This problem has been recently faced by Li *et al.* [17]. Their selection criterion is to minimise cross-entropy between selected test data and the whole set of unlabelled operational data, so as to have a sample of tests representative of the operational context.

This work has been supported by the project COSMIC of UNINA DIETI.

As this approach does not explicitly look for failing examples, it does not give much room for improving the DNN accuracy. From this viewpoint, many tests selected are useless, as they are non-failing examples that serve only the purpose of having highly-confident estimates.

Given the above, we either have a test suite exposing failures but not representative (adversarial-like testing), or representative tests but few failures exposed (operational testing). This paper targets the drawbacks of both strategies together. We present DeepEST (Deep neural networks Enhanced Sampler for operational Testing), an operational testing method for DNN that builds test suites from an operational dataset so as to provide a close and efficient estimate of the expected accuracy *and* to find, at the same time, a high number of failing examples. Paraphrasing a famous aphorism by O. Wilde (“Experience is the hardest kind of teacher. It gives you the *test* first and the *lesson* afterward.”), DeepEST aims to sample, from operational data, many failing tests to learn from while building a set of tests suited for the DNN accuracy estimate.

With respect to pure operational testing, DeepEST is expected to provide DNN accuracy estimates that are more precise and more efficient (for number of tests required), plus the practical advantage of enabling debug-testing-like scenarios (improving accuracy through debugging/re-training). With respect to adversarial-like testing, DeepEST is able to provide accuracy estimates, enabling evaluation of alternative designs or DNN fine tuning, and establishing a release criterion directly related to what will be observed in operation.

Experiments with various DNN and datasets are presented, which assess DeepEST effectiveness, sensitivity to the number of tests to select from the operational dataset, and dependency on the dataset. As DeepEST can exploit various types of auxiliary information to select tests, the experimental study considers four variants. The results show that all the variants produce accuracy estimates similar to those of the state-of-the-art technique, while detecting from 5 to 30 times more mispredictions, with the same sample size.

The paper is structured as follows. Section II introduces sampling-based operational testing concepts used by DeepEST. Section III describes the DeepEST algorithm. Sections IV and V report the experimentation. Section VI discusses related work. Section VII concludes the paper.

II. SAMPLING-BASED OPERATIONAL TESTING OF DNNs

A. Operational testing of DNNs

In the traditional testing literature, *operational testing* refers to the family of techniques that use an operational profile to test a system to estimate its expected *reliability* (i.e., probability of not failing) in operation. Likewise, the primary goal of *DNN operational testing* is to estimate the expected *accuracy* (i.e., probability of not having mispredictions) in a given operational context [17].¹ Two main challenges arise.

Data skew. The idea of operational testing is that testing should not just care about exposing possible failures, but

should be able to spot those failure-causing inputs that are more likely to occur in operation. In case of a relevant mismatch between the pre-release test data and the post-release context, the system could be stimulated in operation with inputs never seen during testing, with unexpected failures.

Data skew is a concern for DNN, more than for traditional software systems. These are expected to work on a range of input data given to functionalities, and can be tested on a small carefully-selected sample of input data (e.g., obtained by input space partitioning). DNN are data-driven by nature; they are constructed around a training dataset, and generalising beyond what observed during training is hard. This data-driven approach is governed by a statistical process and, due also of the black-box nature of DNN, it is tricky to identify classes of inputs that homogeneously represent the expected behaviour in operation. For instance, white-box partitioning has been shown to be not clearly correlated to the failing behaviour [12]. Thus, a *drift* of the post-release operational context from the pre-release testing context is more likely to cause unexpected failures compared to traditional software.

The imitation bias of operational testing. The operational context drift is just a triggering condition for unexpected failures at runtime. The problem occurs because of the way in which operational testing is conducted. Operational testing selects a small data sample that can accurately represent the population; however, the mere imitation of the expected input can be inefficient, especially in highly reliable systems, because many failure-free tests are executed to get an acceptable estimate. A representative sample would roughly contain the same proportion of examples causing misprediction as the operational dataset. There are two problems with this: first, in highly accurate DNNs, the number of examples causing mispredictions is low, thus requiring other types of testing activities dedicated to detect mispredictions (e.g., through adversarial-like techniques) to possibly improve the accuracy. Second, just mimicking the expected usage is fine from the estimation point of view *as long as* the imitation is faithful. If this is not the case, the risk of overestimation increases: if we only aim at having a representative sample of tests, the actual experienced accuracy may be significantly smaller than the estimated one if the operational context drifts, because of the occurrence of unexpected mispredictions.

Thus, a conventional approach for DNN operational testing allows to obtain the desired accuracy estimate but, since it reveals few mispredictions, can be ineffective (leading to wrong estimates) and inefficient (requiring further separate testing activities to detect mispredictions). Recent results in operational testing for traditional software show that exposing failures and estimating reliability are not contrasting objectives [18]. A strategy that actively looks for failures (rather than just mimicking the expected usage) can lead to accurate and stable estimates and, at the same time, expose many failures. Our aim is exactly to spot failing examples in a DNN operational dataset, while preserving the ability to yield effective (small error) and efficient (low variance) estimates.

¹We use the terms misprediction and failure interchangeably for DNN.

B. Sampling-based testing

DeepEST uses statistical sampling, a natural way to cope with estimation problems: it serves to design sampling plans tailored for a population under study, providing effective and efficient estimators. In sampling-based testing [19], the sample is the set of n test cases $T=\{t_1, \dots, t_n\}$, having a binary outcome (pass/fail). Test outcomes are a series of independent Bernoulli random variables z_{t_i} such that $z_{t_i}=1$ if the DNN predicts the correct label for t_i , $z_{t_i}=0$ otherwise. The parameter of our interest is the DNN accuracy θ ; we aim to an estimate $\hat{\theta}$ with two desirable properties: *unbiasedness* – i.e., the expectation of the estimate $\mathbb{E}[\hat{\theta}]$ should be equal to the true value θ – and *efficiency* – for the given the sample size, the variance of $\hat{\theta}$ should be as low as possible (for a highly confident, stable estimate). The probability that $z_{t_i} = 1$ corresponds to the true (unknown) proportion: $\theta = \frac{\sum_{t=1}^N z_{t_i}}{N}$, with N being the population size (i.e., the size of the operational dataset).

Simple random sampling with replacement (SRSWR) is the baseline approach: an unbiased estimator of θ is the observed proportion of correct predictions over the number of trials n :

$$\hat{\theta}_{SRSWR} = \frac{\sum_{t=1}^n z_{t_i}}{n}. \quad (1)$$

Having assumed independent variables, the variance of $\hat{\theta}$ is:

$$V(\hat{\theta}_{SRSWR}) = \frac{\theta(1-\theta)}{n}. \quad (2)$$

An improvement is represented by simple random sampling *without replacement* (SRSWOR), namely, the same test case is not selected twice: this reduces the variance to:

$$V(\hat{\theta}_{SRSWOR}) = \frac{N-n}{N-1} \frac{\theta(1-\theta)}{n}. \quad (3)$$

While SRS keeps the mathematical treatment simple, it is unable to exploit additional information a tester might have.

Exploiting *auxiliary information* to modify the sampling scheme is what is done in sampling theory to get more efficient estimators [20]. The sampling is made proportionally to an auxiliary observable variable assumed to be related to the (unknown) quantity to estimate; the estimator is then adjusted to account for the non-uniform sampling, so as to preserve unbiasedness. For instance, stratified sampling is a strategy which uses knowledge about which sample units are expected to have homogeneous values, and selects units contributing more to lower the estimate’s variance.

Li *et al.* [17] present two sampling strategies for DNN, which are, to our knowledge, the only attempt to DNN operational testing: Confidence-based Stratified Sampling (CSS) and Cross Entropy-based Sampling (CES) – the latter being the authors’ proposal. Both strategies exploit auxiliary information to drive the sampling task.

In CSS, sampling is proportional to the confidence value provided by classifiers when predicting a label: examples with higher confidence are more likely to be selected as part of the test suite. This works well when the classifier is reliable, namely if examples with higher confidence are actually those

for which the prediction is more likely to be correct (in other words, the model is perfectly trained for the operation context). Whenever the operation context drifts from the training one, CSS exhibits poor performance.

CES attempts to overcome this limitation by using the output of the m neurons in the last hidden layer in the last hidden layer, assumed to be more robust to the operation context drift. It builds the test suite trying to minimize the average cross-entropy between the probability distribution of the m -dimensional representation of the output of neurons computed on the operational dataset and on the selected tests.

To pursue the double objective of sampling cases causing mispredictions and estimating accuracy efficiently, DeepEST adopts an *adaptive* and *without-replacement* sampling algorithm, described in Section III.

C. Auxiliary information

To look for mispredictions, the auxiliary information leveraged by DeepEST adaptive sampling should represent the belief about some factor(s) related to the model’s (in)accuracy. We consider two auxiliary variables, related to somehow opposite sources of information: the *confidence* value and the *distance* between the operational dataset example and the training dataset. The latter is based on the result by Kim *et al.* [14] that inputs more “distant” from training data are more likely to cause misprediction – they show that distance is correlated to mispredictions, which is what we look for. We borrow their distance metrics *Distance-based Surprise Adequacy* (DSA) and *Likelihood-based Surprise Adequacy* (LSA). These are computed by the *Activation Trace* (AT), namely a vector of activation values of each neuron of a certain layer corresponding to a certain input; we compute AT with reference to the last DNN activation layer. LSA is defined as the negative log of density (computed via Kernel Density Estimation). DSA is defined as the Euclidean distance between the AT of a new input and ATs observed during training. The variant using confidence is DeepEST_{C_S}; the variants using the distance metrics are DeepEST_{DSA} and DeepEST_{LSA}.

Performance of an auxiliary variable can strongly depend on the DNN and on the training/operation dataset: for instance, for a distance metric the belief that examples far from the one in the training set are more likely to cause a misprediction is not an absolute truth. In particular, if an example is very similar to many others in the training set according to the distance metric, but has a different label, distance will be not a good metric to select it. Similarly, the confidence is a good proxy *if* the DNN is well trained for the operational context. In general, relying on a single auxiliary variable may work well in some settings and bad in others. Combining multiple beliefs is a choice that is expected to improve the stability of results across multiple settings. Based on this, we define a further variant of our algorithm, named DeepEST_C, considering as auxiliary variable the combination of confidence and distance:

$$P = P_c \times (1 - P_d) \quad (4)$$

where P_c is the confidence value, P_d is the DSA value normalized $[0, 1]^2$ and P is probability of correct prediction. The intuition behind is that the probability of a correct prediction is related to both the confidence of the DNN and the *drift* of the example from what seen during training. In fact, P_c is related to the probability that the DNN does a correct prediction *according to what seen during training* (in other words: it is the probability of correct prediction with perfect training); P_d is a proxy for the probability of wrong prediction related to how far the example is from what seen during training – hence due to the *imperfection* of training. If confidence P_c is high *and* the example is close to the training dataset (i.e., P_d is small), there is a high chance of correct prediction.

III. THE DEEPEST ALGORITHM

Both CES and CSS select a sample of tests representative of the operational context. DeepEST takes a different approach, favouring the sampling of failing (namely, *misprediction-causing*) tests, then used to estimate the DNN accuracy. While the estimation ability demands for probabilistic (hence “sampling-based”) selection of examples, the very idea of DeepEST is that, since mispredictions are usually rare compared to correct examples, looking for failing tests is well handled by *adaptive sampling* [21]: the examples progressively selected may give hints about the probability of finding other failing examples, so as to adapt the sampling process accordingly. Given a samples size, adaptive sampling implicitly assumes that inputs of interest are not uniformly spread across the input space, and adopts a disproportional selection to spot them, counterbalanced by the estimator to preserve unbiasedness. In software testing, this is expected to give smaller variance compared to conventional sampling (e.g., SRS), since the failing inputs are usually not uniformly distributed. The DeepEST sampling algorithm is inspired to the Adaptive Web Sampling [22], a flexible design for sampling rare populations.

The above auxiliary variables are used to define a *weight* $w_{i,j}$ between any pair of examples i and j of the operational dataset, used to explore the example space adaptively. For instance, let us assume to use the normalised DSA *distance*: if the example i has distance P_{d_i} (representing the belief that i causes a misprediction), and a pre-defined threshold τ is exceeded (i.e., i has a sufficiently high distance compared to the others), then all the $w_{i,j}$ values ($\forall j$ of the operational dataset) are set to their distance P_{d_j} ; otherwise $w_{i,j} = 0$. This way, a strong-enough belief about example i causing misprediction entails the activation of all the weights toward i . The latter are used, as explained hereafter, for sampling, and makes the algorithm follow the distance criterion to spot potential clusters of failing examples.

The DeepEST sampling strategy is sketched in Algorithm 1. Assuming n examples to be selected from the operational dataset as test cases, the algorithm selects and executes one test case per step. The first input is selected by simple random

²In DeepEST_C, DSA is preferred to LSA since it has been shown to have better performance for the deeper layer [14].

Algorithm 1: DeepEST adaptive sampling

Data: OD : operational dataset; i : current sample; T : test suite; n : number of tests; r : probability of WBS

- 1 $T = \emptyset$;
- 2 $i = \text{SRS_sampling}(OD)$; // first sample by SRS
- 3 $OD = OD \setminus i$; // remove sample from dataset
- 4 $T = T \cup i$; // add to suite
- 5 $y_1 = \text{labelling_and_test_execution}(i)$;
- 6 **for** $k = 2; k \leq n; k++$ **do**
- 7 $rs = \text{random}(0, 1)$;
- 8 **if** $rs < r$ **then**
- 9 $i = \text{WBS_sampling}(OD)$; $OD = OD \setminus i$;
- 10 **else**
- 11 $i = \text{SRS_sampling}(OD)$; $OD = OD \setminus i$;
- 12 $T = T \cup i$;
- 13 $y_i = \text{labelling_and_test_execution}(i)$;
- 14 $z_k = \text{Equation 6}$;
- 15 $\hat{\theta} = \text{Equation 7}$; // compute final estimate

sampling, namely, initially all examples have equal probability of being selected. Then, one of two sampling schemes is used to select next test: *weight-based sampling* (WBS), or *simple random sampling* (SRS). Example i is selected from the operational dataset at step k with probability $q_{k,i}$ given by:

$$q_{k,i} = r \cdot \frac{\sum_{j \in s_k} w_{i,j}}{\sum_{h \notin s_k, j \in s_k} w_{h,j}} + (1-r) \cdot \frac{1}{N - n_{s_k}} \quad (5)$$

where:

- r : probability of using WBS (hence, probability of using SRS = $1-r$; r is set to 0.8 in our implementation);
- s_k : current sample (all examples selected up to step k);
- $w_{i,j}$: weight relating example j in s_k to example i ;
- n_{s_k} : size of the current sample s_k ;
- N : size of the operational dataset.

For both WBS and SRS, the selection is *without replacement*.³ WBS selects an example i proportionally to the sum of weights $w_{i,j}$ of already selected examples toward i – the chance of taking i depends on the current sample, favouring the identification of clusters of failing examples *if* the auxiliary variable (hence, $w_{i,j}$) is well-correlated with mispredictions. As this is not always the case, the WBS “depth” exploration is balanced by SRS, chosen with probability $(1-r)$, for a breadth exploration of the example space. This diversification in the search is useful to escape from unproductive cluster searches. The steps are repeated until the testing budget $n \leq N$ is over.

At step 1, the probability that a randomly selected example will cause a misprediction is estimated as the outcome y_1 of the first test (1 in case of misprediction, 0 otherwise).

³Without replacement sampling schemes generally give smaller variance than their with replacement counterpart on the same sample size [20].

TABLE I: Experimental DNN and datasets

	Subject		Classes	Training set size	Test set size	True accuracy
	DNN	Dataset				
S1	CN5	MNIST	10	60,000	10,000	0.9905
S2	LeNet5					0.9868
S3	CN12	CIFAR10	10	50,000	10,000	0.8066
S4	VGG16					0.9359
S5	VGG16	CIFAR100	100	50,000	10,000	0.7048

At step $k > 1$, example i , whose outcome is y_i , is selected with probability $q_{k,i}$ according to Eq. 5, and the estimator of the misprediction probability is that by Hansen-Hurwitz [23]:

$$z_k = \frac{1}{N} \left(\sum_{j \in s_k} y_j + \frac{y_i}{q_{k,i}} \right) \quad (6)$$

where the y_j values are the outcome of the tests already selected. z_k is an unbiased estimator of the expected misprediction probability at step k ; the final estimator of the expected accuracy of the DNN is 1 minus the average of the z_k values:

$$\hat{\theta} = 1 - \frac{1}{n} \left(y_1 + \sum_{k=2}^n z_k \right) \quad (7)$$

where n is the number of tests run.

IV. EVALUATION

A. Experimental subjects

The four variants of DeepEST are evaluated against the SRSWR scheme as baseline and the mentioned state-of-the-art technique CES. Five experiments are conducted with four DNN models and three popular datasets.

The datasets are MNIST, a dataset of handwritten digits [24]; CIFAR10, for image processing systems; and CIFAR100, similar to the previous one but with 100 classes [25]. The chosen DNN models are ConvNet5 (here simply CN5) and LeNet5 for MNIST classification; ConvNet12 (simply, CN12) and VGG16 for CIFAR10 classification; VGG16 for CIFAR100.⁴ Table I lists the five subjects (DNN-dataset pairs); the true accuracy is in the last column.

B. Research questions and experiment design

The evaluation answers the following research questions.

RQ1: Effectiveness. *How does DeepEST perform in finding inputs causing misprediction (i.e., failing examples) and simultaneously estimating a DNN operational accuracy?*

To gauge DeepEST ability to provide effective DNN accuracy estimates with few examples, while spotting a high number of failing inputs, we set to 200 the number of tests to select (then varied to answer RQ2) and repeat 30 times the execution of the 6 compared techniques on the 5 subjects.

As for evaluation metrics, we compute:

⁴CN5 and CN12 are calibrated in the same way as Li *et al.* [17]; LeNet5 is calibrated as Kim *et al.* [14]; for the VGG16 network we considered the weights at <https://github.com/geifmany/cifar-vgg>.

- The accuracy $\hat{\theta}_i$ at the i -th repetition, and then compute the Mean Squared Error (MSE) as $MSE(\hat{\theta}) = \frac{1}{30} \sum_{i=1}^{30} (\hat{\theta}_i - \theta)^2$, where θ is the true operational accuracy. Note that for unbiased estimators, MSE and variance can be considered indistinguishable. In fact, $MSE = Variance + Bias^2$ and $Bias(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta = 0$. The precision of the estimator is: $\pi(\hat{\theta}) = \frac{1}{MSE(\hat{\theta})}$, and the relative precision (or relative efficiency) of estimator A with respect to B is: $\pi_{A,B} = \frac{MSE(\hat{\theta}_B)}{MSE(\hat{\theta}_A)}$ ($\pi_{A,B} > 1$ means that A is better than B).
- The average number of failures ($\varphi = Mean(\varphi_i)$) with φ_i being the number of failures in repetition i . For comparison purpose, we consider the relative number of failures of technique A with respect to B : $\rho_{A,B} = \frac{\varphi_A}{\varphi_B}$ ($\rho_{A,B} > 1$ means that A is better than B).

RQ2: Sensitivity to sample size. *How does the performance of DeepEST vary with the sample size?*

It is important to figure out how performance varies with the number of test cases to select from the operational dataset, namely with sample size. Indeed, DeepEST aims to perform well especially with a small sample size, so as to yield precise estimates with relatively few examples to be manually labelled.

RQ3: Dataset influence. *How is DeepEST performance affected by the datasets?*

In DNN testing, results are often heavily dependent on the (training and operational) datasets. This RQ aims to figure out how these may influence the ability of the auxiliary variables (confidence, DSA, LSA) to discriminate failing examples, affecting the performance of DeepEST. To answer RQ3, the test set is completely labeled, so as to identify all failures.

C. Implementation

DeepEST is implemented mostly in Java.⁵ The implementation of the distance metric in DeepEST_{DSA} and DeepEST_{LSA} is the same used by Kim *et al.* [14]; we used their Python scripts to compute DSA and LSA values. These are computed considering the last *activation layer* of each DNN. The threshold τ needed for the weights definition is set as follows:

- DeepEST_{DSA}: $\tau = mean(DSA) + 2 \times Std(DSA)$;
- DeepEST_{LSA}: $\tau = mean(LSA) + Var(LSA)$.

The threshold for *confidence*, used by DeepEST_{CES} and DeepEST_C, is set to 0.7, assuming that lower confidence values are more related to misprediction (i.e., the weights are activated when the confidence is less than $\tau = 0.7$).

In CES, the selection exploits the output of the last hidden layer. We use the same configuration as the original article [17]. The size of the initial sample is $p=30$, enlarged by a group Q^* of $q=5$ examples at each step. The number of random groups from which Q^* is selected is $L = 300$. For CES and SRS, we used the Python scripts provided by Li *et al.* [17].

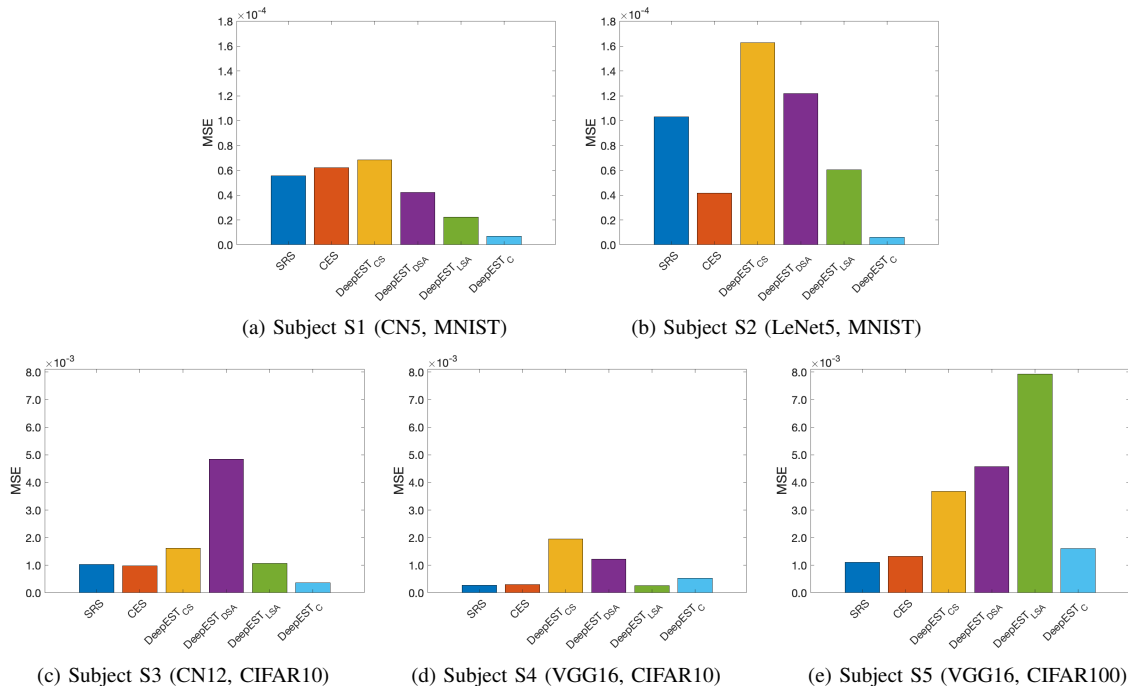


Fig. 1: RQ1 (effectiveness): Mean Squared Error of estimates

TABLE II: RQ1 (effectiveness): Mean and standard deviation (σ) of the number of failing examples detected

Subject	SRS		CES		DeepEST _{CS}		DeepEST _{DSA}		DeepEST _{LSA}		DeepEST _C		Total failing examples
	Mean # %	σ	Mean # %	σ	Mean # %	σ	Mean # %	σ	Mean # %	σ	Mean # %	σ	
S1 (CN5, MNIST)	1.90 <input type="text"/>	1.52	2.23 <input type="text"/>	1.57	44.57 <input type="text"/>	0.68	30.17 <input type="text"/>	4.03	22.43 <input type="text"/>	3.29	11.93 <input type="text"/>	3.06	95
S2 (LeNet5, MNIST)	2.77 <input type="text"/>	2.06	2.03 <input type="text"/>	1.16	65.13 <input type="text"/>	0.86	37.57 <input type="text"/>	5.12	24.10 <input type="text"/>	3.74	24.07 <input type="text"/>	4.15	132
S3 (CN12, CIFAR10)	37.93 <input type="text"/>	6.47	33.77 <input type="text"/>	3.92	106.10 <input type="text"/>	7.26	98.33 <input type="text"/>	4.62	38.07 <input type="text"/>	6.61	70.77 <input type="text"/>	7.09	1,934
S4 (VGG16, CIFAR10)	12.60 <input type="text"/>	3.33	13.03 <input type="text"/>	3.49	85.00 <input type="text"/>	4.86	16.80 <input type="text"/>	3.02	14.77 <input type="text"/>	3.14	24.90 <input type="text"/>	4.97	641
S5 (VGG16, CIFAR100)	57.33 <input type="text"/>	6.53	55.77 <input type="text"/>	6.62	131.17 <input type="text"/>	7.35	78.20 <input type="text"/>	6.62	67.47 <input type="text"/>	3.68	108.70 <input type="text"/>	4.45	2,952

V. RESULTS

A. RQ1: effectiveness

Figure 1 plots the MSE of the estimated accuracy. The techniques exhibit comparable performances, with DeepEST_C being the best one for 3 of the 5 subjects. CES has good performance in terms of MSE, it is the second technique in 3 cases. Considering the single variables: confidence, DSA or LSA lead to results slightly more variable over the subjects – an aspect explored in RQ3. The SRS case is interesting, too: it is never the worst approach and is the best in one case.

Table II reports the average number and the standard deviation of the failing examples detected. All variants of DeepEST identify many more failing examples than SRS and CES, even up to a factor of 30x (DeepEST_{CS} vs CES for subject S2) and reaching in some cases (S1 and S2, DeepEST_{CS})

almost 50% of the total number of failing examples in the datasets (last column). The DeepEST algorithm leverages the adaptive sampling to spot clusters of failing examples with relatively few tests (set to 200 for RQ1). Its performance varies depending on the auxiliary information used, but it is always remarkably better than SRS and CES.

Among the DeepEST variants, confidence (DeepEST_{CS}) turns out to be the most effective auxiliary variable in detecting failures, showing the best performance for all datasets and models, followed by DSA. CES and SRS select the lowest number of mispredicted examples, and are close to each other.

Considering both the failure detection ability and the estimate accuracy, DeepEST_C – that combines confidence and DSA, the two best auxiliary variables for failing examples detection – gives a good trade-off, since it provides stable (across subjects) and close-to-true estimates of the accuracy, with many more detected failing examples than CES and SRS.

⁵A replication package is at: <https://github.com/dessertlab/DeepEST>

TABLE III: Pairwise comparison of techniques. A value of $\rho_{R,C} > 1$ means the technique on the row has a greater precision than that on the column. Similarly for the relative number of detected failures $\pi_{R,C}$

(a) Subject S1 (CN5, MNIST)							(b) Subject S2 (LeNet5, MNIST)								
	<i>row vs col</i>	DeepEST				SRS		<i>row vs col</i>	DeepEST				SRS		
		<i>CS</i>	<i>DSA</i>	<i>LSA</i>	<i>C</i>				$\rho_{R,C}$	$\pi_{R,C}$	$\rho_{R,C}$	$\pi_{R,C}$			
DeepEST	CES	$\frac{\rho_{R,C}}{\pi_{R,C}}$	0.0501 1.0987	0.0740 0.6790	0.0996 0.3564	0.1872 0.1136	1.1754 0.8929	CES	$\frac{\rho_{R,C}}{\pi_{R,C}}$	0.0312 3.9040	0.0541 2.9238	0.0844 1.4523	0.0845 0.1473	0.7349 2.4766	
	<i>CS</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	1.4773 0.6180	1.9866 0.3244	3.7346 0.1034	23.4561 0.8127	<i>CS</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	1.7338 0.7489	2.7026 0.3720	2.7064 0.0377	23.5422 0.6344	
	<i>DSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	1.3447 0.5249	2.5279 0.1673	15.8772 1.3150	<i>DSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	1.5588 0.4967	1.5609 0.0504	13.5783 0.8470	
	<i>LSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	1.8799 0.3187	11.8070 2.5054	<i>LSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	1.0014 0.1014	8.7108 1.7053	
	<i>C</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	-	6.2807 7.8624	<i>C</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	-	8.6988 16.8140	
(c) Subject S3 (CN12, CIFAR10)							(d) Subject S4 (VGG16, CIFAR10)								
	<i>row vs col</i>	DeepEST				SRS		<i>row vs col</i>	DeepEST				SRS		
		<i>CS</i>	<i>DSA</i>	<i>LSA</i>	<i>C</i>				$\rho_{R,C}$	$\pi_{R,C}$	$\rho_{R,C}$	$\pi_{R,C}$			
DeepEST	CES	$\frac{\rho_{R,C}}{\pi_{R,C}}$	0.3183 1.6669	0.3434 4.9728	0.8870 1.0967	0.4772 0.3741	0.8902 1.0522	CES	$\frac{\rho_{R,C}}{\pi_{R,C}}$	0.1533 6.5937	0.7758 4.1503	0.8826 0.8772	0.5234 1.7910	1.0344 0.9106	
	<i>CS</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	1.0790 2.9833	2.7872 0.6579	1.4993 0.2244	2.7970 0.6312	<i>CS</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	5.0595 0.6294	5.7562 0.1330	3.4137 0.2716	6.7460 0.1381	
	<i>DSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	2.5832 0.2205	1.3895 0.0752	2.5923 0.2116	<i>DSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	1.1377 0.2114	0.6747 0.4315	1.3333 0.2194	
	<i>LSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	0.5379 0.3411	1.0035 0.9594	<i>LSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	0.5930 2.0417	1.1720 1.0380	
	<i>C</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	-	1.8656 2.8126	<i>C</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	-	1.9762 0.5084	
(e) Subject S5 (VGG16, CIFAR100)							(f) Number of wins and losses								
	<i>row vs col</i>	DeepEST				SRS		<i>wins</i> (<i>r vs c</i>)	CES	DeepEST				SRS	Total wins
		<i>CS</i>	<i>DSA</i>	<i>LSA</i>	<i>C</i>					<i>CS</i>	<i>DSA</i>	<i>LSA</i>	<i>C</i>		
DeepEST	CES	$\frac{\rho_{R,C}}{\pi_{R,C}}$	0.4252 2.7761	0.7131 3.4999	0.8266 3.5473	0.5130 1.2076	0.9727 0.8323	CES	-	0	0	0	0	0	0/25
	<i>CS</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	1.6773 1.2607	1.9442 1.2778	1.2067 0.4350	2.2878 0.2998	<i>CS</i>	0	-	2	1	0	0	3/25
	<i>DSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	1.1591 1.0136	0.7194 0.3450	1.3640 0.2378	<i>DSA</i>	1	0	-	1	0	1	3/25
	<i>LSA</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	0.6207 0.3404	1.1767 0.2346	<i>LSA</i>	2	0	0	-	0	3	5/25
	<i>C</i>	$\frac{\rho_{R,C}}{\pi_{R,C}}$	-	-	-	-	1.8959 0.6892	<i>C</i>	3	0	2	2	-	3	10/25
							SRS	1	0	0	0	0	-	1/25	
							Total losses	7/25	0/25	4/25	4/25	0/25	7/25		

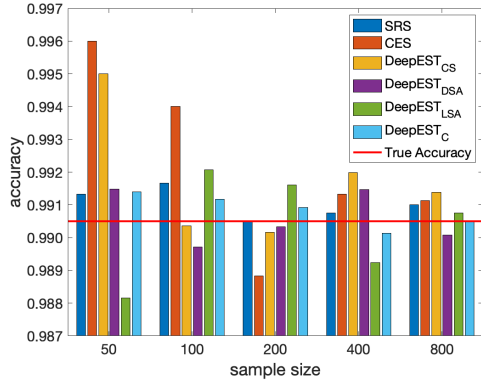
Tables III(a)–(e) show the results of the pairwise comparison of the techniques. For the four DeepEST variants, rows and columns headings list (in blue) the name of the auxiliary variables. The evaluation metrics are the ratio ρ of the failing examples and the relative precision π of the estimators. Values of ρ or π greater (lower) than 1 mean the technique on the row (column) has better performance. If a technique is better than the other in a pair for both metrics (values colored in the table), we say that it *wins*.

Table III(f) summarizes the number of wins (and losses). CES never wins over other approaches, while DeepEST_C wins against CES 3 out of 5 times. SRS never wins over DeepEST, while it wins *vs* CES considering VGG16 on CIFAR100. DeepEST_C never loses and collects the highest number of wins (10), showing the best trade-off among accuracy of the estimation and number of detected failing examples.

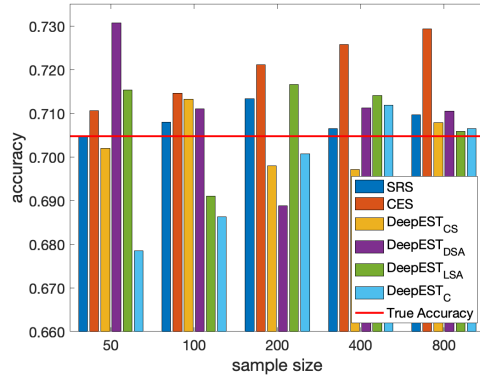
The choice of the DeepEST variant may be determined by which auxiliary information can be collected. We see that DeepEST_C, exploiting a combination of two variables, has good and more stable results in terms of MSE than the other variants, at the expense of a slight decrease of detected failures. Single auxiliary variables are more sensitive to the specific dataset/model pair (e.g., confidence works well if the DNN is reliable), but expose more mispredictions. Confidence has the advantage of not requiring knowledge of the hidden layers and is easier to compute. When no information is available or easily computable, SRS could be a good low cost solution.

B. RQ2: sensitivity to sample size

To answer this RQ, experiments are run with the sample sizes 50, 100, 200, 400, 800, considering the subject with the highest accuracy (S1), and the one with the lowest accuracy

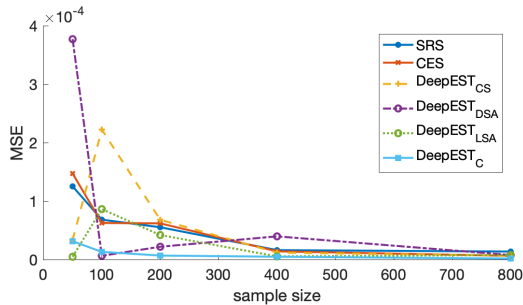


(a) Subject S1 (CN5, MNIST)

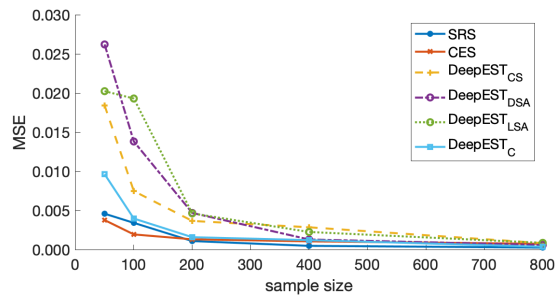


(b) Subject S5 (VGG16, CIFAR100)

Fig. 2: RQ2 (sensitivity to sample size): Accuracy for the most (a) and the least (b) accurate subjects

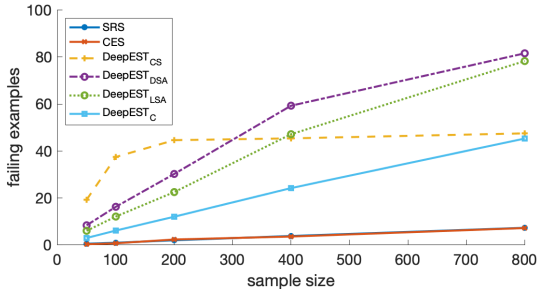


(a) Subject S1 (CN5, MNIST)

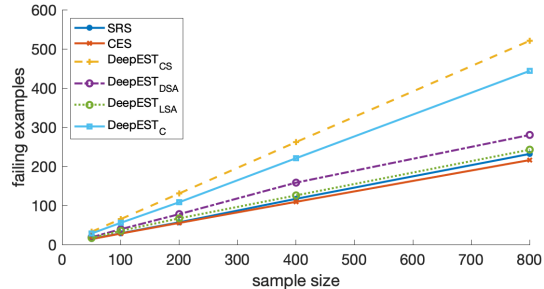


(b) Subject S5 (VGG16, CIFAR100)

Fig. 3: RQ2 (sensitivity to sample size): MSE for the most (a) and the least (b) accurate subjects



(a) Subject S1 (CN5, MNIST)



(b) Subject S5 (VGG16, CIFAR100)

Fig. 4: RQ2 (sensitivity to sample size): Number of failing examples for the most (a) and the least (b) accurate subjects

(S5), so as to analyze how DeepEST performs when there are very few and many failing examples in the dataset, respectively. Figure 2 shows the mean values of the estimates' accuracy over repetitions. Figures 3 and 4 plot the MSE and the mean number of detected failures, respectively. Expectedly, increasing the sample size all techniques exhibit a decreasing trend in MSE and an increasing trend in failing examples.

For the subject with highest accuracy (S1), we observe the following. DeepEST_C shows very good performance for MSE (Fig. 3(a), Fig. 1), and it detects on average about six times the failures of SRS and CES (Table II). The advantage for MSE is more pronounced with the smallest sample sizes, which make DeepEST_C particularly suited when the number of examples

to select and label is very small. For larger sample sizes, the MSE is similar but the advantage of DeepEST_C over SRS and CES is very pronounced for detected failures (Fig. 4(a)). DeepEST_{CS} is the best among all techniques to detect failures for small sample sizes; for sizes 400 and 800, the best is DeepEST_{DSA} (Fig. 4(a)). Although the estimates are unbiased (hence, they tend to the true value), if we look at the mean estimates over the repetitions (Figure 2(a)), CES shows bad performance with up to 200 test cases. SRS works well with low budget; its good results may be influenced by the very low number of failures: 18/30 repetitions show 0 failures and 100% accuracy. It is interesting to note that in most cases CES and SRS overestimate the true accuracy – an undesired

property, especially for critical systems. This is related to the low number of failures detected, as discussed in Section II.

For the subject with lowest accuracy (S5), CES and SRS outperform DeepEST_C for small sample sizes (50 and 100); from 200 to 800, the estimation by CES starts diverging, while SRS and DeepEST keep good performance. The tendency to overestimate the accuracy by CES and SRS is confirmed. Performance in failing examples detection is always clearly in favour of all DeepEST variants. Performance in estimation accuracy is almost specular to what observed with the most accurate model. A reason is that DeepEST is a sampling techniques particularly suited for rare populations, which is not the case of S5. As for the ability to detect failing examples, *confidence* is the best auxiliary variable for DeepEST for subject S5: it presents the best values in all configurations.

C. RQ3: dataset influence

We have seen in the experiments for RQ1 and RQ2 that no single auxiliary variable performs best in all situations. For instance, we can consider the confidence a good auxiliary information for subject S1, and a bad choice for S5. This may depend on several factors: assuming a perfect training, the confidence could be affected by a bias in the training set; or, with a perfect training set, a wrong training phase (e.g., due to overfitting) could generate mispredictions with high confidence. In some cases the operational dataset could contain examples very similar (i.e., small distance) to those in the training set but with a different label, affecting the discriminative power of the DSA and LSA metrics. To analyze how the three datasets influence the ability of auxiliary variables to discriminate failing examples (impacting the performance observed in RQ1-RQ2), we consider the subjects S1 and S5, as for RQ2, plus the VGG16 DNN for CIFAR10.

Figures 5, 6, and 7 show the logistic regression for the three datasets. The curves fit the probabilities for the outcome fail/pass to the three predictors: confidence, DSA, and LSA. Consider MNIST, for which CN5 reaches the highest accuracy among all subjects; the probability for a test to fail is very low for values of confidence between 0.7 and 1 (Figure 5). This is clearly not the case for CIFAR10 (Figure 6) and, especially, CIFAR 100 (Figure 7). In the latter, there is a high chance of misprediction even with high values of confidence; this could be due to a high skew between training and test data.

The discriminating power of DSA and LSA is clearly greater for MNIST, as the high slope of the S-shaped curves in Figures 5(b)-(c) suggests, compared to corresponding ones in Figures 6 and 7. In this case, it actually happens that the farthest examples have highest probability to be related to a failure, with a sharp increase after $DSA \approx 0.9$ and $LSA \approx 300$. This means that if in operation there are (a lot of) examples far from what observed in training, re-training with a more representative dataset can be useful to improve the accuracy. This behaviour is not observed for CIFAR10 and CIFAR100, and DSA and LSA do not seem to be effective in dividing the two sets of examples. In CIFAR10, the DSA line is more horizontal, meaning that the DSA value does not reflect well

the failure probability. The scarcely discriminative power of the auxiliary variables in CIFAR10 and CIFAR100 partially explains the smaller gain of DeepEST over CES and SRS (especially in terms of MSE); nevertheless, its adaptivity allows spotting many failing examples even in these conditions.

Finally, it is interesting to highlight the performance of DeepEST_{CS} (based on confidence) on MNIST in RQ2. The saturation in its failure detection ability (Figure 4(a)) can be explained observing that, after a number of tests able to spot failures looking at low-confidence examples, the few remaining ones with high confidence are selected with low probability; the sharper discrimination made by DSA and LSA determines a high detection ability even when few failing examples remain. In summary, whenever a tester has good belief/evidence about the appropriateness of one of the above auxiliary variables, it is a good choice to select the specific DeepEST variant; if not, the combined variant DeepEST_C has shown to give the best trade-offs in all five experimented cases.

D. Threats to validity

A threat to the internal validity comes from the selection of the experimental subjects. To favor the repeatability of the experiments under different possibly influencing factors, we have used publicly available networks and pre-trained models, to avoid incorrect implementation. The configuration of parameter r in DeepEST and a different setting of thresholds may also affect the results (in terms of efficiency of the estimator), hence a fine-tuning is suggested before applying the method to other dataset-DNNs. Although the code developed was carefully inspected, a common threat is the correctness of the scripts to collect data and compute the results.

The choice of the sample size influences the effectiveness too. We ran a sensitivity analysis to show that DeepEST is still effective (compared to both CES and SRS) under five (from 50 to 800) values of the sample size, but different values could yield different results. Threats to external validity depend on both the number of models and datasets considered for experimentation. We strived to control this threat considering different widely used DNN and datasets. Although the results may change with different subjects, the diversity and significance of the chosen subjects give confidence to the general considerations. Replicability of the experiments on other subjects is to further mitigate this threat.

VI. RELATED WORK

Testing of DNN is a hot research topic. Zhang et al. [26] present a survey on Machine Learning testing, identifying three main families of techniques: mutation testing [27], metamorphic testing [9], [28], and cross referencing [7], [29]. The former two are for test generation: they generate adversarial examples causing mispredictions. Cross-referencing can highlight the most interesting test cases when different implementations of a system disagree. Most of these techniques are meant for fault detection, rather than for accuracy estimate.

Testing DNN for accuracy is the focus of Li et al. [17], who presented CES – which DeepEST has been compared against

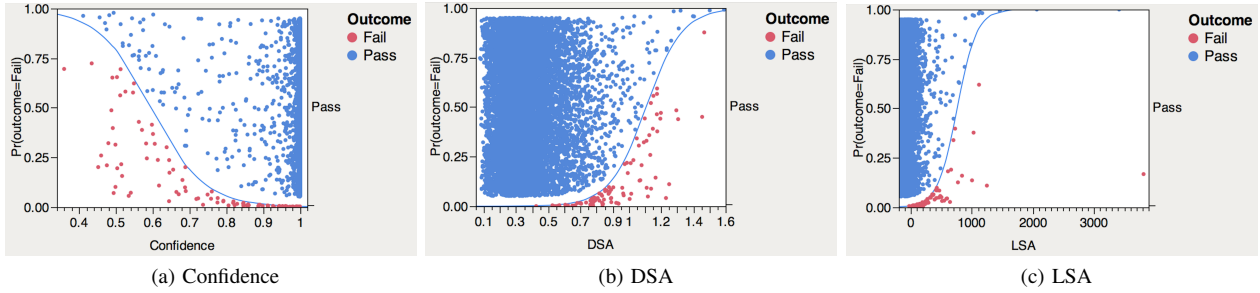


Fig. 5: RQ3 (dataset influence): MNIST samples distribution

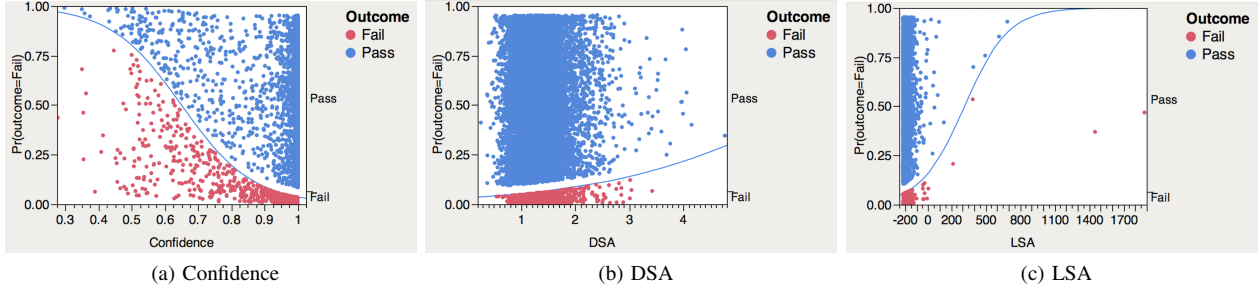


Fig. 6: RQ3 (dataset influence): CIFAR10 samples distribution

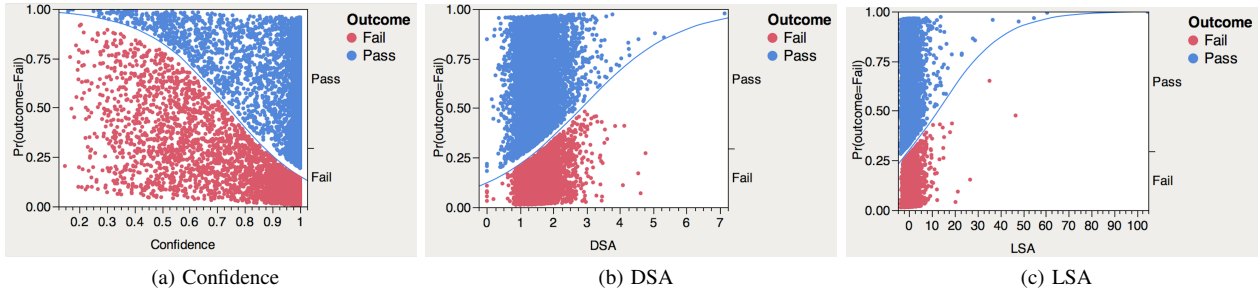


Fig. 7: RQ3 (dataset influence): CIFAR100 samples distribution

- as the first approach using operational testing to this aim. DeepEST differs from CES in several aspects, the key ones being the sampling algorithm and the used auxiliary variables, which are conceived to improve failing examples detection while preserving the accuracy estimation power.

Operational testing, where tests are derived according to the operational profile, is an established practice to estimate software reliability [30]. It was the core technique of Cleanroom software engineering [31]–[34] and of the Software Reliability Engineering Test process [30]. Cai *et al.* developed *Adaptive Testing*, still based on the operational profile, but foreseeing adaptation in assigning tests to partitions [35]–[38]. Recently, Pietrantuno *et al.* stressed the use of unequal probability sampling to improve the estimation efficiency [39], to this aim formalizing several sampling schemes [19]. Our proposal goes along this direction: we do not sample tests representative of the operational dataset, but we “alter” the selection and counterbalance the uneven selection in the estimator. Unequal probability, without-replacement and adaptive sampling are the key concepts we borrowed for operational testing of DNNs.

VII. CONCLUSION

Testing the accuracy of a DNN with operational data aims at precise estimates with small test suites, due to the cost for manually labelling the selected test cases. This effort motivates to pursue also the goal of exposing many mispredictions with the same test suite, so as to improve the DNN after testing. With these two concurrent goals, we presented DeepEST, a technique to select failing tests from an operational dataset, while ultimately yielding faithful estimates of the accuracy of a DNN under test.

We evaluated experimentally four variants of DeepEST, based on various types of auxiliary information that its adaptive sampling strategy can leverage. The results with four DNN models and three popular datasets show that all DeepEST variants provide accurate estimates, compared to existing sampling-based DNN testing techniques, while generally much outperforming them in exposing mispredictions. Practitioners may choose the appropriate variant, depending on the characteristics of their operational dataset and on which auxiliary information is available or they can collect.

REFERENCES

- [1] Ziad Obermeyer and Ezekiel J. Emanuel. Predicting the future — big data, machine learning, and clinical medicine. *New England Journal of Medicine*, 375(13):1216–1219, 2016. PMID: 27682033.
- [2] Mariusz Bojarski *et al.* End to End Learning for Self-Driving Cars. arXiv:1604.07316, 2016.
- [3] Amir Efrati. Uber finds deadly accident likely caused by software set to ignore objects on road. *The Information*, May 7, 2018.
- [4] Jack Stewart. Tesla’s autopilot was involved in another deadly car crash. [Online]. Available: <https://www.wired.com/story/tesla-autopilot-self-driving-crash-california/>, 2018.
- [5] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. DeepTest: Automated testing of deep-neural-network-driven autonomous cars. In *40th International Conference on Software Engineering*, ICSE, pages 303–314. ACM, 2018.
- [6] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. DeepFace: Closing the gap to human-level performance in face verification. In *IEEE Conference on Computer Vision and Pattern Recognition*, CVPR, pages 1701–1708. IEEE, 2014.
- [7] Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana. Deepxplore: Automated whitebox testing of deep learning systems. *Communications of the ACM*, 62(11):137–145, 2019.
- [8] Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, Jianjun Zhao, and Yadong Wang. DeepGauge: Multi-Granularity Testing Criteria for Deep Learning Systems. In *33rd ACM/IEEE International Conference on Automated Software Engineering*, ASE, pages 120–131. ACM, 2018.
- [9] Mengshi Zhang, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. DeepRoad: GAN-Based Metamorphic Testing and Input Validation Framework for Autonomous Driving Systems. In *33rd ACM/IEEE International Conference on Automated Software Engineering*, ASE, pages 132–142. ACM, 2018.
- [10] Lei Ma, Fuyuan Zhang, Minhui Xue, Bo Li, Yang Liu, Jianjun Zhao, and Yadong Wang. Combinatorial testing for deep learning systems. arxiv.org/abs/1806.07723, 2018.
- [11] Augustus Odena and Ian Goodfellow. TensorFuzz: Debugging neural networks with coverage-guided fuzzing. In *36th International Conference on Machine Learning*, volume 97 of *Proc. of Machine Learning Research*, 2019.
- [12] Zenan Li, Xiaoxing Ma, Chang Xu, and Chun Cao. Structural coverage criteria for neural networks could be misleading. In *41st International Conference on Software Engineering: New Ideas and Emerging Results*, ICSE-NIER, pages 89–92. IEEE, 2019.
- [13] Weibin Wu, Hui Xu, Sanqiang Zhong, Michael R. Lyu, and Irwin King. Deep Validation: Toward Detecting Real-World Corner Cases for Deep Neural Networks. In *49th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, DSN, pages 125–137. IEEE, 2019.
- [14] Jinhan Kim, Robert Feldt, and Shin Yoo. Guiding Deep Learning System Testing Using Surprise Adequacy. In *41st International Conference on Software Engineering*, ICSE, pages 1039–1049. IEEE, 2019.
- [15] Phyllis G. Frankl, Richard G. Hamlet, Bev Littlewood, and Lorenzo Strigini. Evaluating testing methods by delivered reliability. *IEEE Transactions on Software Engineering*, 24(8):586–601, 1998.
- [16] Michael R. Lyu, editor. *Handbook of Software Reliability Engineering*. McGraw-Hill, Inc., Hightstown, NJ, USA, 1996.
- [17] Zenan Li, Xiaoxing Ma, Chang Xu, Chun Cao, Jingwei Xu, and Jian Lü. Boosting Operational DNN Testing Efficiency through Conditioning. In *Proc. of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ESEC/FSE, pages 499–509. ACM, 2019.
- [18] Domenico Cotroneo, Roberto Pietrantuono, and Stefano Russo. RELAI testing: a technique to assess and improve software reliability. *IEEE Transactions on Software Engineering*, 42(5):452–475, 2016.
- [19] Roberto Pietrantuono and Stefano Russo. On adaptive sampling-based testing for software reliability assessment. In *27th International Symposium on Software Reliability Engineering*, ISSRE, pages 1–11. IEEE, 2016.
- [20] Sharon L. Lohr. *Sampling: Design and Analysis*. Duxbury Press, 2009.
- [21] Steven K. Thompson. *Sampling, Third Edition*. John Wiley & Sons, Inc., 2012.
- [22] Daniel G. Horvitz and Donovan J. Thompson. A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260):pp. 663–685, 1952.
- [23] Morris H. Hansen and William N. Hurwitz. On the Theory of Sampling from Finite Populations. *The Annals of Mathematical Statistics*, 14(4):333–362, 1943.
- [24] 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.
- [25] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical Report TR-2009, University of Toronto, 2009.
- [26] Jie M. Zhang, Mark Harman, Lei Ma, and Yang Liu. Machine learning testing: Survey, landscapes and horizons. *IEEE Transactions on Software Engineering*, pages 1–37, 2020.
- [27] Lei Ma, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Felix Juefei-Xu, Chao Xie, Li Li, Yang Liu, Jianjun Zhao, and et al. DeepMutation: Mutation Testing of Deep Learning Systems. In *29th International Symposium on Software Reliability Engineering*, ISSRE, pages 100–111. IEEE, 2018.
- [28] Xiaoyuan Xie, Joshua W. K. Ho, Christian Murphy, Gail Kaiser, Baowen Xu, and Tsong Yueh Chen. Testing and validating machine learning classifiers by metamorphic testing. *Journal of Systems and Software*, 84(4):544–558, 2011.
- [29] Siwakorn Srisakaokul, Zhengkai Wu, Angello Astorga, Oreoluwa Alebiosu, and Tao Xie. Multiple-implementation testing of supervised learning software. In *AAAI Workshops*. Association for the Advancement of Artificial Intelligence, 2018.
- [30] John D. Musa. Software reliability-engineered testing. *Computer*, 29(11):61–68, 1996.
- [31] Harlan D. Mills, Michael Dyer, and Richard C. Linger. Cleanroom software engineering. *IEEE Software*, 4(55):19–24, 1987.
- [32] P. Allen Currit, Michael Dyer, and Harlan D. Mills. Certifying the reliability of software. *IEEE Transactions on Software Engineering*, SE-12(1):3–11, 1986.
- [33] Richard H. Cobb and Harlan D. Mills. Engineering software under statistical quality control. *IEEE Software*, 7(6):45–54, 1990.
- [34] Richard C. Linger and Harlan D. Mills. A case study in cleanroom software engineering: the IBM COBOL Structuring Facility. In *12th International Computer Software and Applications Conference*, COMP-SAC, pages 10–17. IEEE, 1988.
- [35] Junpeng Lv, Bei-Bei Yin, and Kai-Yuan Cai. On the asymptotic behavior of adaptive testing strategy for software reliability assessment. *IEEE Transactions on Software Engineering*, 40(4):396–412, 2014.
- [36] Junpeng Lv, Bei-Bei Yin, and Kai-Yuan Cai. Estimating confidence interval of software reliability with adaptive testing strategy. *Journal of Systems and Software*, 97:192–206, 2014.
- [37] Kai-Yuan Cai, Chang-Hai. Jiang, Hai Hu, and Cheng-Gang Bai. An experimental study of adaptive testing for software reliability assessment. *Journal of Systems and Software*, 81(8):1406–1429, 2008.
- [38] Kai-Yuan Cai, Yong-Chao Li, and Ke Liu. Optimal and adaptive testing for software reliability assessment. *Information and Software Technology*, 46(15):989–1000, 2004.
- [39] Roberto Pietrantuono and Stefano Russo. Probabilistic sampling-based testing for accelerated reliability assessment. In *International Conference on Software Quality, Reliability and Security*, QRS, pages 35–46. IEEE, 2018.