

SPARDACUS SafetyCage: A new misclassification detector

Anonymous Full Paper
Submission 29

Abstract

Given the increasing adoption of machine learning techniques in society and industry, it is important to put procedures in place that can infer and signal whether the prediction of an ML model may be unreliable. This is not only relevant for ML specialists, but also for laypersons who may be end-users. In this work, we present a new method for flagging possible misclassifications from a feed-forward neural network in a general multi-class problem, called SPARDA-enabled Classification Uncertainty Scorer (SPARDACUS). For each class and layer, the probability distribution functions of the activations for both correctly and wrongly classified samples are recorded. Using a Sparse Difference Analysis (SPARDA) approach, an optimal projection along the direction maximizing the Wasserstein distance enables p -value computations to confirm or reject the class prediction. Importantly, while most existing methods act on the output layer only, our method can in addition be applied on the hidden layers in the neural network, thus being useful in applications, such as feature extraction, that necessarily exploit the intermediate (hidden) layers. We test our method on both a well-performing and underperforming classifier, on different datasets, and compare with other previously published approaches. Notably, while achieving performance on par with two state-of-the-art-level methods, we significantly extend in flexibility and applicability. We further find, for the models and datasets chosen, that the output layer is indeed the most valuable for misclassification detection, and adding information from previous layers does not necessarily improve performance in such cases.

1 Introduction

A crucial consideration when deploying a machine learning (ML) model in real-life applications is the ability to infer how reliable the predictions are. As an example, consider a model used to detect a hazardous situation within an industrial facility. First, it is important that the model can capture an unsafe situation as it happens to then signal the operators. Second, to achieve trust, the model should not raise false alarms too frequently, or future warnings lose credibility. Thus, to deploy the ML model in a real-

world setting, the reliability of its predictions need to be considered in some way prior to making actual decisions (such as stopping the production line).

In most neural network (NN) classifiers, a softmax activation function is utilized on the output layer to interpret each value as the probability of belonging to a particular class. The class prediction for any sample is most often equal to the class with the largest softmax probability in the output layer.

It has been shown that misclassifications may arise even when the largest softmax probability is close to one [1]. Nonetheless, a pattern was discovered where the maximum softmax probability tended to be smaller for incorrectly classified samples than for correctly classified samples. This discovery was used to make a simple threshold-based misclassification detector, see [2], called Maximum Softmax Probability (MSP) Detector.

In [3], a method named DOCTOR for misclassification detection was proposed, based on an approximation of the misclassification probability, $\text{Pe}(x)$, for a particular sample x by only using the softmax output layer values, $P_{\hat{Y}|X}(c | x)$, for each class c of total C classes. In particular $\text{Pe}(x) \approx 1 - \sum_{c=1}^C P_{\hat{Y}|X}^2(c | x)$. The method flags a prediction as untrustworthy whenever the odds of a misclassification event is larger than some threshold.

From our literature search, it is apparent that the DOCTOR and MSP-detector methods represent the current state of the art of misclassification detection, which we use for comparison to SPARDACUS.

The *SafetyCage*, introduced in [4], is another misclassification detector. This statistical framework collects the pre-activation vector in each layer, and assumes the corresponding multivariate probability density function (PDF) of correctly predicted in-distribution samples to follow a Gaussian distribution. These PDFs, per class, are fitted to the training data of the NN model. To infer the uncertainty of a class prediction, the Mahalanobis distance, inspired by the approach described in [5], is used to measure the likeliness that the pre-activation values, in each layer, are generated from the fitted Gaussian distribution of correctly predicted samples. This Mahalanobis-based SafetyCage was tested on two feed-forward neural networks trained on benchmark datasets MNIST and CIFAR-10, respectively. The classifier trained on the MNIST dataset had an accuracy of 0.93, whereas the one trained on CIFAR-10

096 had an accuracy of 0.48. It was observed that for
097 the well-performing MNIST model, the multivariate
098 Mahalanobis SafetyCage was able to detect and flag
099 60% of the wrong classifications. On the other hand,
100 for the CIFAR-10 model with poor performance, the
101 SafetyCage was no better than random guessing. Af-
102 ter closer inspection, the assumption of Gaussianity
103 of the pre-activation vector for the CIFAR-10 model
104 was not accurate [4].

105 While the Mahalanobis-based SafetyCage only
106 uses samples that are correctly predicted, we pro-
107 pose the SPARDACUS method for misclassification
108 detection that uses both correctly and wrongly clas-
109 sified samples. We compare the results of SPAR-
110 DACUS to the previous SafetyCage, the DOCTOR
111 method, and the MSP-detector method.

112 We note that a task related to the detection of
113 misclassifications is what is called *out-of-distribution*
114 (OOD) detection, where the aim is to detect when-
115 ever an input sample is inherently different from
116 the data used during training of the model, and
117 hence the corresponding prediction should not be
118 trusted. However, the most insidious misclassifica-
119 tions happen with in-distribution data, for which
120 the model would be assumed to work correctly, and
121 not OOD-data. Indeed, in [6] it is shown that the
122 best OOD-detector is not always the best at de-
123 tecting NN classification errors. The authors in [6]
124 further emphasize that if the focus is on use of NNs
125 in safety-critical applications, misclassification de-
126 tection should be the paramount focus, and not OOD-
127 detection. For these reasons, this work focuses on
128 misclassification detection using in-distribution data,
129 and will not draw a comparison to OOD-detection
130 methods.

131 2 Methods

132 2.1 SPARDACUS

133 Consider the function $\mathcal{F}_{i,l}$ which corresponds to the
134 PDF that generates the activation values at layer
135 l for a sample correctly predicted as class i by the
136 NN classifier. Conversely, let $\mathcal{G}_{i,l}$ correspond to the
137 PDF that generates the activation values for incor-
138 rectly classified samples. Given these two PDFs,
139 one may infer which distribution a new test sample
140 x' belongs to by designing a decision procedure to
141 flag wrongly classified samples. In principle this is
142 a binary classification problem to which we could
143 apply any ML method to *predict* if x' is correctly
144 or wrongly classified; but if $\mathcal{F}_{i,l}$ and $\mathcal{G}_{i,l}$ could be
145 approximated directly, statistical tests backed by
146 solid theory would become available. A direct appli-
147 cation of this procedure is however challenging, since
148 the dimensionality of the PDFs is linked to the size
149 of the NN classifier’s layers; i.e. large layers imply
150 high-dimensional PDFs. To combat this, we project

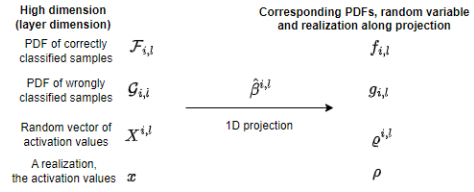


Figure 1. Notation in high dimension and correspond-
ing notation in dimension 1 after the random vector $X^{i,l}$
is projected onto the one-dimensional subspace defined
by $\hat{\beta}^{i,l}$ for class i and layer l .

the data for each layer and class to one dimension, 151
effectively collapsing the multi-dimensional PDFs 152
into one-dimensional ones. The goal is to obtain two 153
PDFs along this 1D-projection, denoted $f_{i,l}$ and $g_{i,l}$, 154
which are minimally overlapping. The PDFs are 155
estimated using the same training data the classifier 156
is constructed from. 157

To this end, Mueller et al. [7] propose an 158
approach referred to as a Sparse Differences 159
Analysis (SPARDA), which given observed samples 160
from two multivariate PDFs, searches for the 161
optimal projection maximising the Wasserstein 162
distance between the projected 1D empirical 163
distribution functions (ECDFs). This is a non- 164
smooth, non-concave optimization problem. We 165
apply the fastSPARDA optimization algorithm 166
available at [https://bitbucket.org/jwmueller/](https://bitbucket.org/jwmueller/principal-differences-analysis/src/master/) 167
[principal-differences-analysis/src/master/](https://bitbucket.org/jwmueller/principal-differences-analysis/src/master/). 168
The optimization problem includes a regularization 169
parameter λ to induce sparsity in the projection. We 170
denote the projection direction given by SPARDA 171
as $\hat{\beta}^{i,l}$, and $\rho^{i,l}$ as the projected random variable 172
of activations along $\hat{\beta}^{i,l}$. Figure 1 summarises the 173
projection operation and notation. 174

If a new sample x' is predicted to be a member 175
of class i , one can infer at any given layer l whether 176
it is more likely generated from $f_{i,l}(\rho)$ or $g_{i,l}(\rho)$ by 177
inspecting the observed value ρ along the projection. 178
However, $f_{i,l}(\rho)$ and $g_{i,l}(\rho)$ are unknown. To over- 179
come this, we fit each PDF as a Gaussian mixture 180
model due to its flexibility and computational effi- 181
ciency. A typical situation with overlapping can be 182
seen in Figure 2 showing histograms of $f_{7,-1}$ and 183
 $g_{7,-1}$ on training data and test data for the class 184
(digit) 7 from MNIST, with $l = -1$ indicating the 185
output layer. 186

187 2.2 Inference

With inspiration from the likelihood ratio test 188
we define the statistic, $S^{i,l}(\rho^{i,l})$, using the two afore- 189
mentioned PDFs for a random variable $\rho^{i,l}$ with 190
observed values along the projection $\hat{\beta}_{i,l}$: 191

$$S^{i,l}(\rho^{i,l}) = -\ln \left(\frac{f_{i,l}(\rho^{i,l})}{g_{i,l}(\rho^{i,l})} \right). \quad (1) \quad 192$$

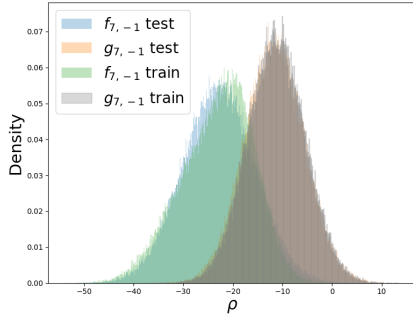


Figure 2. Illustration of the PDFs for $f_{7,-1}$ and $g_{7,-1}$ on both the training data and the test data. This represents the digit 7, output layer, and the MNIST model.

193 Moreover, consider the statistic $S_C^{i,l}(\varrho_f^{i,l})$ where
 194 $\varrho_f^{i,l} \sim f_{i,l}(\rho)$. The corresponding PDF of $S_C^{i,l}$, de-
 195 noted $h_C(s_C^{i,l})$, is the PDF we expect for samples
 196 with correct class prediction i . The larger the value
 197 of $s_C^{i,l}$, the more unlikely it is generated from $h_C(s_C^{i,l})$.
 198 We define the following hypothesis test for a test
 199 sample x' with class prediction \hat{y}' :

$$H_0: S^{\hat{y}',l} \sim h_C(s^{\hat{y}',l}) \quad \text{vs.} \quad H_1: S^{\hat{y}',l} \not\sim h_C(s^{\hat{y}',l}). \quad (2)$$

200 A corresponding p -value for an observed value
 201 $s^{\hat{y}',l}$ can be computed as $p_{S_C} = P(S_C^{\hat{y}',l} \geq s^{\hat{y}',l}) =$
 202 $1 - H_C(s^{\hat{y}',l})$ with $H_C(s^{\hat{y}',l})$ the CDF of $S_C^{\hat{y}',l}$. Recall
 203 that a p -value under the null hypothesis will follow
 204 a uniform probability distribution ($p_{S_C}^{H_0} \sim U(0, 1)$).
 205 A small p -value indicates the unlikelihood that $s^{\hat{y}',l}$
 206 is generated from $f_{\hat{y}',l}(\rho)$. Rather, it is an indica-
 207 tion that $s^{\hat{y}',l}$ is generated from $g_{\hat{y}',l}(\rho)$. The null
 208 hypothesis can then be rejected, ultimately flagging
 209 a prediction as wrong, for a p -value less than a
 210 predefined significance level α_C .

212 Notice that one can also define the statistic
 213 $S_W^{i,l}(\varrho_g^{i,l})$ where $\varrho_g^{i,l} \sim g_{i,l}(\rho)$. Using this statistic
 214 instead, one can construct in the same way a hypoth-
 215 esis test with significance level α_W . See Appendix
 216 A for more information.

217 We emphasize that, unlike the Mahalanobis SafetyCage [4], the SPARDACUS SafetyCage does not rely on pre-activation values, as there is no longer need to assume Gaussianity.

2.2.1 Estimation of p -values from the S statistics

223 As the PDFs $f_{i,l}$ and $g_{i,l}$ are estimated as a Gaus-
 224 sian mixture model, the corresponding PDFs of $S_C^{i,l}$
 225 and $S_W^{i,l}$ are not always easily accessible. An easy
 226 way to estimate the PDFs of $S_C^{i,l}$ and $S_W^{i,l}$ is via
 227 Monte Carlo simulation. After a specified number of
 228 repeated generations from $f_{i,l}$ and $g_{i,l}$, we can com-

pute the ECDFs of $S_C^{i,l}$ and $S_W^{i,l}$, denoted $\hat{F}_C^{i,l}(s_{i,l})$ 229
 and $\hat{F}_W^{i,l}(s_{i,l})$. From this we can estimate the corre- 230
 sponding p -values as: 231

$$\hat{p}_{S_C}^l(s_{\hat{y}',l}) = 1 - \hat{F}_C^{\hat{y}',l}(s_{\hat{y}',l}), \quad 232$$

and 233

$$\hat{p}_{S_W}^l(s_{\hat{y}',l}) = \hat{F}_W^{\hat{y}',l}(s_{\hat{y}',l}). \quad 234$$

Whether to use $S_C^{i,l}$ or $S_W^{i,l}$ to get the most ro- 235
 bust results will be based on the accuracies of the 236
 estimated PDFs $f_{i,l}$ and $g_{i,l}$. 237

The notation so far has been focused on a partic- 238
 ular layer l . By our method, each layer l provides a 239
 p -value. These p -values can be combined into one, 240
 using any p -value combination test. To this end, 241
 we will apply the Cauchy combination test as it is 242
 robust for statistically correlated p -values [8]. We 243
 denote \hat{p}_{S_C} and \hat{p}_{S_W} the final p -values when using ei- 244
 ther statistic S_C or S_W , respectively. Algorithms for 245
 the method is given in B.1 and B.2 in the Appendix 246
 where we separate between the training phase of the 247
 SPARDACUS method, and the subsequent detection 248
 procedure. 249

3 Results 250

We will evaluate our method by using the same 251
 setup as in [4], where a feed-forward neural network 252
 is trained on MNIST, yielding a well-performing 253
 model (accuracy of 0.98), and on CIFAR-10 yielding 254
 a poor-performing model (accuracy of 0.48). The 255
 neural network consists of two hidden layers, with 256
 256 and 128 neurons respectively, and ReLu activa- 257
 tion functions, along with an output layer featuring 258
 ten neurons with Softmax activation. Training uti- 259
 lized the standard Adam optimizer. Once the model 260
 is trained, our method estimates the projections $\hat{\beta}_{i,l}$ 261
 and corresponding PDFs $f_{i,l}$, $g_{i,l}$ for all classes and 262
 layers (except the input layer, i.e. the images). On 263
 a held-out test data disjoint from the data the NN 264
 model was trained on, we evaluate our misclassifi- 265
 cation detector for different values of α_C and α_W . 266
 Additionally, we present tables showing the preci- 267
 sion, recall, specificity, negative predictive value as 268
 well as the MCC for the best performing α value. 269

We advocate for using the Matthews correlation 270
 coefficient (MCC) as the most meaningful perfor- 271
 mance metric to evaluate any binary classification 272
 model [9]. The MCC ranges from -1, meaning the 273
 classifier is always wrong, to 1 meaning a perfect 274
 classifier. A coin tossing classifier with a 50 % chance 275
 to assign a prediction to either of the two classes will 276
 give MCC = 0 [9]. Also by definition, a classifier that 277
 predicts only one class every time will give MCC = 278
 0 (see [9] for more details). We compare our method 279
 with the MSP-detector method introduced in [2], the 280
 DOCTOR method [3], and the SafetyCage method 281
 introduced in [4]. The threshold in SPARDACUS 282

283 and in [4] can be interpreted as the same, as they
 284 both are based on p -values. The threshold for the
 285 MSP-detector method is with respect to the prob-
 286 ability of the prediction, while for the DOCTOR
 287 method it is with respect to the estimated odds of
 288 misclassification.

289 For the Mahalanobis SafetyCage and the SPAR-
 290 DACUS method, we evaluate the results for different
 291 layer aggregations (combining p -values from differ-
 292 ent layers) to investigate the information that is
 293 covered in the different layers. Specifically, we in-
 294 vestigate when only applying the output layer (out),
 295 the penultimate layer (pen) or all hidden layers plus
 296 the output layer (all).

297 3.1 SPARDACUS Results

298 In all results to come, the SPARDA projection $\hat{\beta}_{i,l}$
 299 was computed by setting $\lambda = 0$, imposing no regular-
 300 ization, and using the fastSPARDA algorithm. We
 301 used Monte Carlo simulations to estimate the CDFs
 302 of the S statistics by generating 1 million samples.

303 Tables 1 and 2 show the best results, ranked with
 304 respect to the MCC, for the four methods on the
 305 MNIST model and CIFAR-10 model respectively
 306 when the optimal threshold as well as the param-
 307 eters and PDFs for the Mahalanobis SafetyCage and
 308 SPARDACUS methods are estimated on the training
 309 data. The methods are evaluated on the unseen
 310 test data.

311 Notice that for all methods, only using the output
 312 layer turned out to give the largest MCC value. The
 313 same applies to the SafetyCage Mahalanobis, an
 314 aspect that was not investigated in the original paper
 315 [4] where only hidden layers up to the penultimate
 316 layer were investigated.

317 We also show the precision, recall, specificity, and
 318 negative predictive value (NPV) for each case. While
 319 the outputs of the MSP and DOCTOR methods
 320 are deterministic, the 1D-projections calculated in
 321 SPARDACUS using the SPARDA algorithm, result
 322 in stochasticity. For this method we therefore show
 323 the result of five runs in terms of mean and standard
 324 deviation. Notice the small variation in MCC. After
 325 closer inspection, this is due to nearly equal 1D-
 326 projection.

327 In Figures 3 and 4, for each of the misclassification
 328 detectors, we show how the MCC varies for differ-
 329 ent threshold values on the test data for MNIST
 330 and CIFAR-10 respectively. We see that the SPAR-
 331 DACUS method is most sensitive to the choice of
 332 threshold with respect to performance, with a sharp
 333 peak for the optimal threshold. At the same time,
 334 if we compare the MCC values at the peaks for
 335 each method, the SPARDACUS method achieves
 336 by a small margin the highest MCC values for both
 337 datasets. See Table C.1 in Appendix C where the
 338 MCC values at the peaks for the SPARDACUS

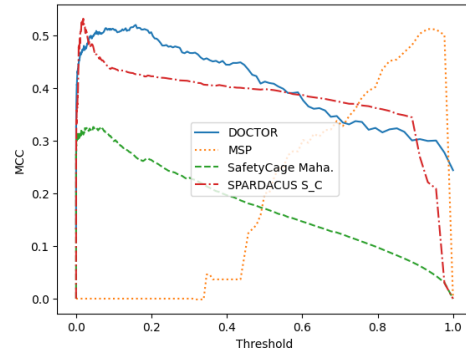


Figure 3. Plots showing thresholds vs MCC for the misclassification detectors for the test data with the MNIST model.

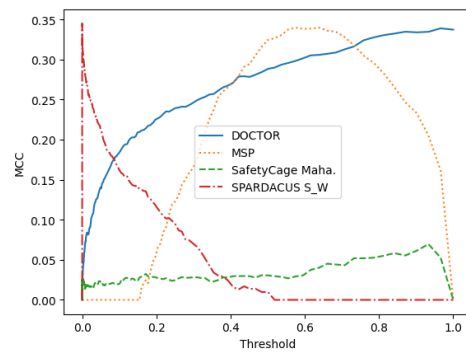


Figure 4. Plots showing thresholds vs MCC for the misclassification detectors for the test data with the CIFAR-10 model.

method are extracted. Here, we also include results
 when applying different sets of layers for the SPAR-
 DACUS method which confirms that only using the
 output layer yields the highest MCC values.

In practice we do not know in advance what the
 optimal threshold is. To estimate the threshold on
 the same data that the ML model is trained on can
 make sense in terms of utilizing all data available,
 particularly for data-driven methods such as the Ma-
 halanobis SafetyCage and the SPARDACUS, where
 also PDFs are fitted to data. However, in certain
 situations the training data is not available, and
 instead new data must be collected during deploy-
 ment. Moreover, using the training data may lead to
 overfitting and less generalization capabilities. For
 comparison, we include the scenario where the pa-
 rameters needed for each misclassification detection
 method, including the threshold value, are estimated
 based on data never used by the ML model. We
 do this by randomly splitting the test data (10 000
 samples) equally in two subsets, estimate the param-
 eters on the first subset, and evaluate the detection
 methods on the second subset. To account for differ-
 ence in performance with respect to the splitting
 of data, we repeat the process five times each with
 random splitting of the test data. The results are

Method	S	L	Threshold	Prec.	Recall	Spec.	NPV	MCC
DOCTOR	—	out	1.13E-01	0.409	0.662	0.978	0.992	0.507
MSP	—	out	9.46E-01	0.408	0.658	0.978	0.992	0.504
SPARDACUS	S_C	out	9.41E-03	0.457	0.554	0.985	0.990	0.491 \pm 0.010
SafetyCage Maha.	—	out	8.98E-03	0.213	0.518	0.957	0.988	0.310

Table 1. Results ordered by the MCC score for each method on the MNIST dataset when the threshold is optimized on the training data, and evaluated on the test data. Due to Stochasticity in the SPARDACUS method we do five runs, and present the performance metrics with the respect to the average, and additionally the MCC with \pm standard deviation. The SafetyCage Mahalanobis method is taken from [4]. The MSP-detector and DOCTOR methods utilise a threshold T , while SPARDACUS and SafetyCage have a significance level α .

Method	S	L	Threshold	Prec.	Recall	Spec.	NPV	MCC
SPARDACUS	S_W	out	3.67E-06	0.646	0.783	0.550	0.707	0.343 \pm 0.003
MSP	—	out	6.27E-01	0.636	0.807	0.515	0.718	0.338
DOCTOR	—	out	9.99E-01	0.623	0.858	0.462	0.744	0.337
SafetyCage Maha.	—	out	8.98E-03	0.213	0.51	0.15	0.57	0.070

Table 2. Same procedure and results as for Table 1, however with respect to the CIFAR-10 dataset.

365 given in Tables 3 and 4.

366 Based on the experiments and following results
 367 we see that the SPARDACUS method is superior to
 368 the SafetyCage Mahalanobis method from [4]. More-
 369 over, compared to the MSP-detector and DOCTOR
 370 methods, the SPARDACUS method using the S_C
 371 and S_W is essentially on par.

372 As expected, all methods performed much better
 373 on the well-performing model, due to there being
 374 more useful information encoded in the activation
 375 values that could be extracted and used to flag in-
 376 correct predictions. Interestingly, the output layer
 377 was by far the most useful layer, which aids to show
 378 how this is a markedly different problem compared
 379 to OOD where the output layer can be less infor-
 380 mative [5]. An interesting observation can be made
 381 regarding the SPARDACUS method yielding the
 382 best results when evaluating only the output layer.
 383 When inspecting the fitted projection at the output
 384 layer, $\hat{\beta}_{i,-1}$, and specifically looking at using S_C
 385 for the MNIST model, we see that for every pre-
 386 dicted class, the associated projection vector has its
 387 maximal-value element in the position corresponding
 388 to the predicted class itself. In fact, on average the
 389 maximum value along the class prediction dimen-
 390 sion was 2.88 times larger than the second largest
 391 element. This shows that the projection vector for a
 392 specific class is heavily dominated by the dimension
 393 along the class itself. This is in fact in correspon-
 394 dence with the MSP-detector method, where the
 395 corresponding projection would have zero-elements
 396 along all dimensions except for the class dimension.
 397 This shows that, in the special case of being applied
 398 to only the output layer, SPARDACUS can be seen
 399 as an *extension* of the MSP-detector method.

400 In all cases, S_C performed better for the well-
 401 performing model, while S_W was better for the poor-
 402 performing model.

4 Discussion and Conclusion 403

404 In this work, we presented a method to infer whether
 405 a particular sample is wrongly classified by an
 406 underlying NN model. SPARDACUS is based on a
 407 SPARDA projection maximizing the Wasserstein
 408 distance of the PDFs of the samples that were cor-
 409 rectly and wrongly classified, and a hypothesis test
 410 inspired by the likelihood-ratio test. SPARDACUS
 411 can be applied at any stage of an NN classifier, and
 412 with an easy extension it could also draw information
 413 from any arbitrary combination of layers.

414 We tested SPARDACUS on two simple NN classi-
 415 fiers, one well-performing trained on MNIST and
 416 one poorly performing trained on CIFAR-10. The
 417 results have further been compared with three pre-
 418 existing methods from the literature: The DOCTOR
 419 method, the Mahalanobis-based method presented
 420 in [4], and the MSP-detector method [2].

421 Our results show that SPARDACUS significantly
 422 outperforms the Mahalanobis SafetyCage [4], where
 423 in particular, detection performance on the inaccu-
 424 rate classifier for CIFAR-10 has been improved by
 425 an order of magnitude. The MSP, DOCTOR and
 426 SPARDACUS are essentially performing at com-
 427 parable levels with respect to the MCC values as
 428 shown in Tables 1, 2, 3 and 4. It is worth not-
 429 ing how SPARDACUS is an extension of the MSP-
 430 detector, being in principle (for large λ) equivalent
 431 to it when only considering the output layer. How-
 432 ever, in contrast to the MSP and DOCTOR method,
 433 SPARDACUS can draw information from not only
 434 the output layer, but also hidden layers in the NN
 435 classifier. For the particular examples of datasets/-
 436 models shown in this work, including information
 437 from the hidden layers has not shown a definite ad-
 438 vantage when compared to only using the output
 439 layer’s information. Nonetheless, we regard it as

Method	S	L	Threshold	Prec.	Recall	Spec.	NPV	MCC
MSP	—	out	9.33E-01	0.434	0.629	0.980	0.991	0.509 ± 0.0159
DOCTOR	—	out	1.24E-01	0.415	0.657	0.979	0.992	0.507 ± 0.023
SPARDACUS	S_C	out	3.19E-02	0.373	0.700	0.974	0.994	0.496 ± 0.023
SafetyCage Maha.	—	out	2.26E-02	0.173	0.613	0.934	0.991	0.295 ± 0.027

Table 3. Results ordered by the MCC score for each misclassification detection method on the MNIST dataset when the threshold is optimized on half the test data (5000 samples), and the methods are evaluated on the other half. The threshold and performance metrics are presented with the average value based on five random splits of the test data. The variation in MCC is additionally presented in terms of ± standard deviation.

Method	S	L	Threshold	Prec.	Recall	Spec.	NPV	MCC
DOCTOR	—	out	9.55E-01	0.618	0.862	0.443	0.755	0.338 ± 0.012
MSP	—	out	5.98E-01	0.645	0.772	0.554	0.701	0.336 ± 0.009
SPARDACUS	S_W	out	3.20E-06	0.649	0.749	0.579	0.690	0.333 ± 0.013
SafetyCage Maha.	—	out	6.53E-01	0.524	0.634	0.401	0.528	0.043 ± 0.015

Table 4. Same procedure and results as for Table 3, however with respect to the CIFAR-10 dataset.

440 an interesting research path to investigate whether
441 this is always the case, or whether modifications are
442 needed to take more advantage of the information in
443 the hidden layers. Of all investigated methods, the
444 performance of the SPARDACUS method is most
445 sensitive to the choice of threshold. Clearly, SPAR-
446 DACUS is more involved than the MSP-detector
447 and DOCTOR methods both in terms of theoretical
448 background and computational complexity. Even
449 though MSP, DOCTOR and SPARDACUS are close
450 in performance with respect to the MCC, we regard
451 SPARDACUS to have three advantages: Firstly,
452 SPARDACUS is most flexible. The flexibility that
453 the three methods share is the choice of threshold.
454 Other than that the MSP and DOCTOR method is
455 static in terms of deterministic computations (in the
456 output layer only), while the SPARDACUS method
457 can be investigated further in several ways. Here,
458 we list some aspects to investigate: 1) We have
459 used the Wasserstein distance, but one may substi-
460 tute it with other statistically-relevant distances. 2)
461 The SPARDA projection can be investigated fur-
462 ther using other optimization algorithms. 3) The
463 importance from different layers can be weighted
464 during training, e.g. by deploying Cauchy weights
465 in the Cauchy combination test. 4) Experiment-
466 ing with the λ regularization parameter, or other
467 significant statistics in place of $S^{i,l}$ is possible. Sec-
468 ondly, SPARDACUS can be refitted and updated for
469 newly labelled data of misclassifications, in terms
470 of projections and PDFs, while the MSP-detector
471 and DOCTOR methods can only be updated with
472 respect to the threshold. Thirdly, by being able to
473 inspect not only the output layer, but also the inner
474 working of the neural network in terms of previous
475 layers, we believe that the framework SPARDACUS
476 builds on can help us not only answer the question
477 *whether* a classification is false, but also *why* it is
478 false.

References

- [1] I. J. Goodfellow, J. Shlens, and C. Szegedy. *Explaining and Harnessing Adversarial Examples*. Mar. 2015. URL: <http://arxiv.org/abs/1412.6572>.
- [2] D. Hendrycks and K. Gimpel. *A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks*. Oct. 2018. URL: <http://arxiv.org/abs/1610.02136>.
- [3] F. Granese, M. Romanelli, D. Gorla, C. Palamidessi, and P. Piantanida. “DOCTOR: A Simple Method for Detecting Misclassification Errors”. In: *Advances in Neural Information Processing Systems*. Vol. 34. Curran Associates, Inc., 2021, pp. 5669–5681.
- [4] P. V. Johnsen and F. Remonato. “SafetyCage: A misclassification detector for feed-forward neural networks”. In: *Northern Lights Deep Learning Conference 2024*. 2023. URL: <https://openreview.net/forum?id=cWSk611xGo>.
- [5] K. Lee, K. Lee, H. Lee, and J. Shin. *A Simple Unified Framework for Detecting Out-of-Distribution Samples and Adversarial Attacks*. 2018.
- [6] J. Guerin, K. Delmas, R. Ferreira, and J. Guiochet. “Out-of-Distribution Detection Is Not All You Need”. In: *Proceedings of the AAAI Conference on Artificial Intelligence 37.12* (June 2023), pp. 14829–14837. ISSN: 2374-3468. DOI: [10.1609/aaai.v37i12.26732](https://doi.org/10.1609/aaai.v37i12.26732).
- [7] J. W. Mueller and T. Jaakkola. “Principal Differences Analysis: Interpretable Characterization of Differences between Distributions”. In: *Advances in Neural Information Processing Systems*. Ed. by C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R.

515 Garnett. Vol. 28. Curran Associates, Inc.,
516 2015. URL: [https://proceedings.neurips.
517 cc / paper _ files / paper / 2015 / file /
518 83fa5a432ae55c253d0e60dbfa716723-Paper.
519 pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/83fa5a432ae55c253d0e60dbfa716723-Paper.pdf).

520 [8] Y. Liu and J. Xie. “Cauchy combination test:
521 a powerful test with analytic p-value calcula-
522 tion under arbitrary dependency structures”.
523 In: *Journal of the American Statistical Asso-
524 ciation* 115.529 (2020), pp. 393–402. DOI: [10.
525 1080/01621459.2018.1554485](https://doi.org/10.1080/01621459.2018.1554485).

526 [9] D. Chicco and G. Jurman. “The advantages
527 of the Matthews correlation coefficient (MCC)
528 over F1 score and accuracy in binary classifica-
529 tion evaluation”. In: *BMC Genomics* 21.1 (Jan.
530 2020), p. 6. DOI: [10.1186/s12864-019-6413-
531 7](https://doi.org/10.1186/s12864-019-6413-7).

532 A The statistic S_W

533 In a similar fashion as for the statistic S_C , define
534 the statistic

$$535 S_W^{\hat{y}',l} = -\ln \left(\frac{f_{\hat{y}',l}(\varrho_{\hat{y}',l}^{\hat{y}',l})}{g_{\hat{y}',l}(\varrho_{\hat{y}',l}^{\hat{y}',l})} \right), \quad (3)$$

536 where $\varrho_{\hat{y}',l}^{\hat{y}',l} \sim g_{i,l}(\rho)$. Using this statistic, one can
537 construct in the same way as for the statistic S_C ,
538 develop a hypothesis test:

$$539 H_0: S_W^{\hat{y}',l} \sim h_W(s^{\hat{y}',l}) \quad \text{vs.} \quad H_1: S_W^{\hat{y}',l} \not\sim h_W(s^{\hat{y}',l}). \quad (4)$$

540 The null hypothesis in this case shows that the
541 sample x' is wrongly classified, and the p -value is
542 computed as $p_{S_W} = P(S_W^{\hat{y}',l} \leq s^{\hat{y}',l})$, this time as a
543 right-sided test since a small observed $s^{\hat{y}',l}$ now indi-
544 cates the sample x' is correctly classified. The null
545 hypothesis is rejected for a predefined significance
546 level α_w .

547 B Algorithms

548 Algorithm B.1 outlines the training phase of the
549 SPARDACUS method where the 1D projections are
550 estimated followed by PDFs fitted to data along
551 the projections for correctly and wrongly classified
552 predictions. Algorithm B.2 outlines the deployment
553 phase of the SPARDACUS method for a new incom-
554 ing data point $\{x', \hat{y}'\}$.

Algorithm B.1 Training phase of SPARDACUS

Consider a dataset $D = \{x_k, y_k, \hat{y}(x_k)\}_{k=1}^N$ of in-
put samples x_k , true class labels $y_k \in \{i\}_{i=1}^C$ and
model predictions $M(x_k) = \hat{y}_k$. The activation
values of layer l for sample k in model M is de-
noted as $a_{k,l}^{(M)}$. Define statistic $S = S_C$ or $S = S_W$.
Let Q be number of Monte Carlo simulations. Let
 $F^{i,l}$ and $T^{i,l}$ be the set of correctly and incorrect
predictions for class i and layer l .

```

for  $l$  in  $L$  do
  for  $i$  in  $1, \dots, C$  do
     $P_i = \{k \mid \hat{y}(x_k) = i\}$ : Extract samples
    predicted to belong to class  $i$ .
    for  $k \in P_i$  do
       $x^{k,l} = a_{k,l}^{(M)}$ .
      if  $\hat{y}(x_k) = y_k$  then
         $T^{i,l} \leftarrow T^{i,l} \cup \{x^{k,l}\}$ : Add activations
        to set  $T^{i,l}$ .
      else
         $F^{i,l} \leftarrow F^{i,l} \cup \{x^{k,l}\}$ : Add activations
        to set  $F^{i,l}$ .
      end if
     $\hat{\beta}^{i,l} = SPARDA(T^{i,l}, F^{i,l}, d_w)$ : Use
    SPARDA to compute projection  $\hat{\beta}^{i,l}$  that
    maximizes the Wasserstein distance,  $d_w$ ,
    between  $T^{i,l}$  and  $F^{i,l}$ .
    Estimate  $f_{i,l}(\rho)$  using GMM from the
    projections along  $\hat{\beta}_{i,l}$  in  $T^{i,l}$ .
    Estimate  $g_{i,l}(\rho)$  using GMM from from
    the projections along  $\hat{\beta}_{i,l}$  in  $F^{i,l}$ .
    if  $S = S_C$  then
       $\rho_1, \dots, \rho_Q \sim f_{i,l}(\rho)$ : Monte-Carlo gen-
      erations
      Compute  $Q$  realizations from the statis-
      tic in (1).
      Estimate CDF of  $S_C$  via ECDF of the
       $Q$  realizations and collect.
    else if  $S = S_W$  then
       $\rho_1, \dots, \rho_Q \sim g_{i,l}(\rho)$ : Monte-Carlo gen-
      erations
      Compute  $Q$  realizations from the statis-
      tic in (1).
      Estimate CDF of  $S_W$  via ECDF from
      the  $Q$  realizations and collect.
    end if
  end for
end for

```

Algorithm B.2 Inference phase of SPARDACUS

Given input x' with model prediction \hat{y}' , and
threshold $\alpha \in [0, 1]$.

```

for  $l$  in  $L$  do
  Compute  $p$ -value based on estimated CDF
  (ECDF) of  $S^{\hat{y}',l}$ .
end for
Compute Cauchy combination test statistic
 $p_{cauchy}$  from observed  $p$ -values from investigated
layers [8].
Declare prediction  $\hat{y}'$  as false if  $p_{cauchy} < \alpha$ 

```

555 **C All SPARDACUS results**

556 Table C.1 shows the maximum MCC values achiev-
 557 able for SPARDACUS method on the test data for
 558 several different parameter sets including which lay-
 559 ers are used (output, hidden, penultimate, all) and
 560 which statistic (S_C or S_W).

Data	Method	S	L	Thresh.	Prec.	Recall	Spec.	NPV	MCC
MNIST	SPARDACUS	S_C	out	1.46E-02	0.420	0.700	0.970	0.993	0.529 ± 0.002
MNIST	SPARDACUS	S_C	all	3.25E-02	0.332	0.721	0.967	0.993	0.473 ± 0.001
MNIST	SPARDACUS	S_W	out	2.43E-06	0.391	0.551	0.980	0.990	0.450 ± 0.014
MNIST	SPARDACUS	S_C	pen	4.08E-02	0.207	0.406	0.965	0.986	0.267 ± 0.001
CIFAR	SPARDACUS	S_W	all	6.80E-06	0.645	0.790	0.543	0.712	0.345 ± 0.001
CIFAR	SPARDACUS	S_W	out	3.28E-06	0.650	0.773	0.563	0.703	0.344 ± 0.001
CIFAR	SPARDACUS	S_C	out	3.80E-01	0.650	0.690	0.609	0.652	0.301 ± 0.000
CIFAR	SPARDACUS	S_C	all	5.37E-01	0.604	0.774	0.486	0.664	0.255 ± 0.001
CIFAR	SPARDACUS	S_W	pen	1.24E-01	0.593	0.692	0.500	0.608	0.196 ± 0.000

Table C.1. Results for the SPARDACUS method when applying different sets of layers for both the MNIST and CIFAR-10 datasets.