

# BAARD: Blocking Adversarial Examples by Testing for Applicability, Reliability and Decidability

Xinglong Chang<sup>1</sup>, Katharina Dost<sup>1</sup>, Kaiqi Zhao<sup>1</sup>, Ambra Demontis<sup>2</sup>, Fabio Roli<sup>3</sup>, Gillian Dobbie<sup>1</sup>, and Jörg Wicker<sup>1</sup>

<sup>1</sup> The University of Auckland, Auckland, New Zealand  
xcha011@aucklanduni.ac.nz, {katharina.dost, kaiqi.zhao, g.dobbie, j.wicker}@auckland.ac.nz

<sup>2</sup> University of Cagliari, Cagliari, Italy  
ambra.demontis@unica.it

<sup>3</sup> University of Genoa, Genoa, Italy  
fabio.roli@unige.it

**Abstract.** Adversarial defenses protect machine learning models from adversarial attacks, but are often tailored to one type of model or attack. The lack of information on unknown potential attacks makes detecting adversarial examples challenging. Additionally, attackers do not need to follow the rules made by the defender. To address this problem, we take inspiration from the concept of Applicability Domain in cheminformatics. Cheminformatics models struggle to make accurate predictions because only a limited number of compounds are known and available for training. Applicability Domain defines a domain based on the known compounds and rejects any unknown compound that falls outside the domain. Similarly, adversarial examples start as harmless inputs, but can be manipulated to evade reliable classification by moving outside the domain of the classifier. We are the first to identify the similarity between Applicability Domain and adversarial detection. Instead of focusing on unknown attacks, we focus on what is known, the training data. We propose a simple yet robust triple-stage data-driven framework that checks the input globally and locally, and confirms that they are coherent with the model’s output. This framework can be applied to any classification model and is not limited to specific attacks. We demonstrate these three stages work as one unit, effectively detecting various attacks, even for a white-box scenario.

**Keywords:** Adversarial Defense · Anomaly Detection · Applicability Domain · Evasion Attacks · White-box Adaptive Attacks

## 1 Introduction

Machine learning algorithms have shown promising results in many mission-critical fields, such as virtual drug screening [1] and autonomous driving [12].

Unfortunately, despite their high accuracy on benign examples, they are vulnerable to adversarial attacks, where malicious users exploit the classifiers’ weakness by manipulating the input data [7]. Starting from a benign data point, attackers craft a small perturbation that allows them to achieve the desired outcome: misclassification of the input example. For example, by adding a small artifact to a stop sign, a self-driving vehicle can be fooled into misclassifying the stop sign as a speed limit sign, with the risk of causing a car crash [12].

Adversarial detectors extract features from unlabeled examples and use them to identify adversarial examples based on certain thresholds [21]. Existing detectors often suffer from the following issues: First, many detectors focus on detecting adversarial examples with only minimal perturbations [20] and tend to fail to detect stronger ones. Second, many defenses are built on a single assumption or one attack, i.e., adversarial examples lead to overly confident predictions from the classifier [11]. However, attackers are not constrained by such assumptions, as they can easily bypass such a detector by altering their strategy. Third, most defenses are tailored to a specific machine learning architecture and do not generalize to other models [23]. There is a lack of flexible detectors that can detect unseen attacks on various classifiers.

In cheminformatics, models are trained on a finite number of compounds because the data-collecting process is expensive and time-consuming. However, the chemical space is vast and diverse in its properties, so models trained on one part of the space may not work on others. Hence, models typically struggle to generalize unseen compounds. To avoid false predictions, *Applicability Domain* (AD) is a concept that defines a domain in which a model can perform reliably. Compounds that are outside this domain are rejected, as the model cannot make reliable predictions on them [1]. Similar to cheminformatics, adversarial detectors only have the information for known attacks. However, new attacks come out so frequently that it is impossible to cover all attacks. In this paper, instead of defending against previously unseen attacks, we focus on what the classifier can reliably predict, the training data. Inspired by the idea of a triple-stage AD originally introduced by Hanser *et al.* [9], we propose the BAARD framework, **B**locking **A**dversarial examples by testing for **A**pplicability, **R**eliability and **D**ecidability.

To identify unknown attacks, BAARD investigates the example from three different perspectives, utilizing the training data in the following ways: 1. *Applicability Stage* uses the training data to validate the input globally; 2. *Reliability Stage* confirms that the example can be backed up by training data locally; and 3. *Decidability Stage* checks the model’s output to ensure it is coherent with the input. These three stages work as one unit to inspect the model’s interpretation of an unlabeled example. As shown in Fig. 1, BAARD rejects the example if there is an inconsistency between the input and the model’s prediction.

We summarize our contributions as follows:

- We are the first to demonstrate the effectiveness of linking two previously unlinked fields: the Applicability Domain in cheminformatics and adversarial detection in machine learning.

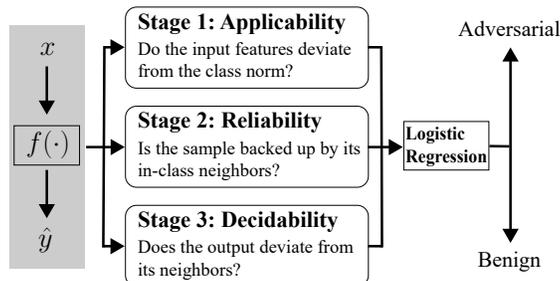


Fig. 1: An overview of BAARD. BAARD analyzes an example  $\mathbf{x}$ , the classifier  $f(\cdot)$ , and its prediction  $\hat{y}$  together by checking the Applicability, Reliability, and Decidability. Each stage outputs a score. The scores are used to train a logistic regression model to predict whether  $\mathbf{x}$  is benign or adversarial.

- Inspired by the Applicability Domain, we propose the BAARD framework (Blocking Adversarial examples by testing Applicability, Reliability, and Decidability), which utilizes training data to systematically detect adversarial examples from three different perspectives.
- By designing an adaptive white-box attack targeting BAARD, we show that it is difficult to penetrate all three stages, even under the worst scenario.
- We demonstrate BAARD is highly portable. This simple yet effective framework can detect adversarial examples with various constraints on a wide range of classifiers, including classifiers that have been neglected previously despite being vulnerable to attacks, such as support vector machines and decision trees.

We introduce the adversarial threat model, attacks and detectors relevant to this paper in Sec. 2. Sec. 3 and 4 present the BAARD framework and demonstrate its effectiveness, respectively. Sec. 5 concludes this paper.

## 2 Background

This paper focuses on detecting *evasion attacks*, where the attacker crafts malicious inputs by adding perturbations to existing examples which can deceive the classifier to make unexpected predictions [7]. Evasion attacks are the most common adversarial attacks since it is easier for a malicious user to interact with the model at inference time.

**Evasion Attacks.** One of the earliest attacks on *neural network* (NN) models is the *Fast Gradient Sign Method* (FGSM) [8], a single-step attack that forms the adversarial example as:  $\mathbf{x}' = \mathbf{x} - \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(\mathbf{x}, y))$  where  $\mathbf{x}$  is a benign input,  $y$  is the targeted label,  $\ell(\mathbf{x}, y)$  is the loss function used by the classifier, and hyperparameter  $\epsilon$  controls the amount of perturbation. *Auto Projected Gradient Descent* (APGD) [6] is the latest improved version of *Projected Gradient Descent*

(PGD) [15]. PGD is a multi-step variant of FGSM. It achieves a higher success rate by iteratively solving the optimization problem. Improving on PGD, APGD dynamically adjusts the number of iterations to ensure minimal perturbation while maintaining the success rate. Directly optimizing on the input space can be difficult, since NN models are highly non-linear. Instead of optimizing on the input space, the *Carlini and Wagner Attack* (CW) [5] transforms the image from the pixel space to the simpler tanh space. Not only NN models are vulnerable to adversarial attacks, the *Decision Tree Attack* (DTA) [18] exploits the data structure of a decision tree. The algorithm makes minimal changes at each node and keeps traversing from the leaf to the root until the prediction from the classifier deviates from the legitimate class.

**Detection.** Detecting adversarial examples with indistinguishable perturbations (hard to recognize by human) has been studied extensively [20]. One common assumption is that if the adversarial perturbation is small enough, the legitimate class can be restored by adding or removing noise. Detectors, such as *Feature Squeezing* (FS) [22] and the *Positive and Negative representation* (PN) detector [13] are motivated by image reconstruction techniques. FS is a defense motivated by using image filters to restore adversarial examples. He *et al.* [10] pointed out that strong adversaries can easily bypass FS. The PN detector assumes an adversary cannot simultaneously deceive a classifier trained on both the original and color-negative images. Such techniques have clear limitations, a detector that uses images’ properties cannot be generalized to other data types.

Another direction is to combine neighborhood relationship and noise generation. *Region-based Classification* (RC) [3] replaces the classifier with a region-based classifier by generating noisy samples centered at the example, and a decision is made via majority voting. Similar to RC, the *Odds are odd* (Odds) [19] detector assumes that adversarial examples are less robust to noise than benign examples. The assumption is that latent outputs significantly change when adding noise to an adversarial example. *Local Intrinsic Dimensionality* (LID) [14] is another neighbor-based algorithm that uses the intrinsic dimension metric by combining latent outputs from all hidden layers of a NN. The statistics are learned by comparing benign, noisy, and adversarial examples. ML-LOO [23] computes Leave-One-Out feature attribution maps on multiple hidden layers of a NN, and uses them to distinguish between benign and adversarial examples. Many detectors are based on certain assumptions of one type of attack. If the attacker’s goal is to bypass the system, such a constraint may not apply [4]. A detection that is tailored to one attack is not robust against white-box attacks, where the attacker knows a particular defense is placed [21].

### 3 BAARD: Blocking Adversarial Examples

This paper connects cheminformatics’ Applicability Domain with adversarial detection in machine learning. The goal of AD is to reject chemical compounds that the classifier cannot reliably predict. Therefore, AD analyzes the feature space

and the classifier together to define a tight region around the training instances but omits the rest of the space [17]. Adversarial examples are perturbations of legitimate example, and remain similar to the original example. However, adversarial examples are designed to cause misclassifications leading to inconsistencies between the predicted labels of the adversarial example and its legitimate neighbors. This observation leads us to believe that the idea used in AD can effectively detect adversarial examples. BAARD consists of three stages as shown in Fig. 1. The rest of this section explains the working of each stage and their effectiveness when combined together.

**Applicability Stage.** In chemistry, this stage checks the compound to confirm it is appropriate for the model to make a prediction [9]. Here, we know the model is trained on the training data, so we check the input feature space by comparing it with the training data globally. We conduct a Z-test by computing mean and standard deviation of input features for each class from the training data. Given an example  $\mathbf{x}$ , the Z-score is defined by  $\mathbf{z}_{\mathbf{x},\hat{y}} := (x - \mu_{X_{\text{train},\hat{y}}})/\sigma_{X_{\text{train},\hat{y}}}$ , where  $\mu_{X_{\text{train},\hat{y}}}$  and  $\sigma_{X_{\text{train},\hat{y}}}$  are the mean and standard deviation for examples in the training data that have the same label as the model’s prediction  $\hat{y}$ , and  $\mathbf{z}_{\mathbf{x},\hat{y}}$  has the same dimension as  $\mathbf{x}$ . Because we are only interested in the extrema and Z-test is two-tailed, we define the Applicability Score as: S1 score :=  $\max(|\mathbf{z}_{\mathbf{x},\hat{y}}|)$ . The Applicability Stage inspects each feature of the new, unlabeled example, individually. It outputs a high score if any feature is significantly different from the training samples that match the classifier’s predicted label.

**Reliability Stage.** Given a compound, this stage quantifies the relevance of information available to the model in chemistry. We implement this stage by examining the input locally using the compound’s neighbors in the training set. Unlike the previous stage, which considered each input feature independently, this stage accounts for all features together using the neighborhood relationship.

Adversarial examples aim to minimize the perturbation while forcing the model to make classification errors [8]. This moves the legitimate input closer to the decision boundary, causing the predicted label to change and potentially placing the example far away from its new in-class neighbors. The reliability test is based on the distances between adversarial examples and their neighbors. These distances are often higher than the distances between legitimate examples and their neighbors.

Choosing an appropriate distance metric is essential when measuring nearest neighbors. The Euclidean distance ( $L_2$ -norm) is well suited for low-dimensional space, but Cosine similarity has shown more robust results in high-dimensional sparse features [13]. Cosine similarity between two feature vectors  $A$  and  $B$  is defined as:  $S_C(A, B) := \sum_{i=1}^n A_i B_i / [(\sum_{i=1}^n A_i^2)^{\frac{1}{2}} (\sum_{i=1}^n B_i^2)^{\frac{1}{2}}]$ , where  $n$  is the dimension of the feature vector, and  $S_C \in [-1, 1]$ . If  $S_C$  is close to 1,  $A$  and  $B$  are positive co-linear vectors. If  $S_C = 0$ , they are independent vectors, and if  $S_C \approx -1$ , they are strong opposite vectors. This means neither minimal nor maximal indicates  $A$  and  $B$  are close. To properly present the distance between

---

**Algorithm 1** BAARD Stage 2 – Reliability Stage

---

**Input:**  $\mathbf{x}$ : unlabeled example,  $\hat{y}$ : its prediction,  $(X, Y)$ : training set,  $k_{S2}$ : number of neighbors, and  $m_{S2}$ : sample size.

**Output:**  $S2\_score \in [0, 2\pi]$

- 1:  $X_{\hat{y}} \leftarrow$  Random sampling  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{m_{S2}}, y_{m_{S2}})\}$ , where  $\mathbf{x}_i \in X$ ,  $y_i \in Y$ , and  $y_i = \hat{y}$
  - 2:  $D(\mathbf{x}, X_{\hat{y}}) \leftarrow$  Compute angular distances between example  $\mathbf{x}$  and subset  $X_{\hat{y}}$
  - 3:  $S2\_score \leftarrow \text{mean}(\text{top\_k}(D(\mathbf{x}, X_{\hat{y}}), k_{S2}))$   $\triangleright$  Compute the mean of top  $k_{S2}$  distances
  - 4: **return**  $S2\_score$
- 

two features using cosine similarity, we compute the angular distance  $D$ , which is defined as:  $D(A, B) := \arccos(S_C(A, B))/\pi$ .

Algorithm 1 provides the pseudocode for this stage. It takes two hyperparameters: the number of nearest neighbors  $k \in \mathbb{N}$ , and the sample size  $m$  that limits the computational expense. For an unlabeled example  $\mathbf{x}$ , the S2 score is the mean distance of the  $k$ -nearest neighbors of  $\mathbf{x}$  within a subset of training data where examples have the same label as the prediction  $\hat{y}$ . Because the angular distance is within  $[0, 2\pi]$ , the S2 score shares a similar scale as the S1 score. To reduce the computational cost, we randomly sample  $m$  instances from the training examples where the legitimate labels are the same as  $\hat{y}$ .

**Decidability Stage.** This stage confirms whether the model’s output is coherent with the evidence from previous stages. Machine learning models operate under the assumption that similar examples have similar labels. Hence, a trained model can generalize to new and previously unseen examples. However, this is often violated when the model tries to predict maliciously crafted adversarial examples. The prediction of an adversarial example often conflicts with the predictions of its neighbors. As shown in Fig. 1, we use the local neighborhood relationship to check adversarial examples based on this property.

Algorithm 2 uses the same distance metric as in previous stages. The critical difference is that the entire training data are used regardless of their labels. We apply the Softmax function so the model outputs probability estimates. Given an example, we run a Z-test on its probability estimates based on its  $k$  neighbors.

**Combining All Stages.** A single stage may be effective on a certain type of attack, but no stage alone can cover all attacks. The Applicability and Reliability Stages both check the feature space but from different perspectives. Once we collect enough evidence from the input space, the Decidability Stage checks the output to ensure the model’s output is coherent with the evidence. We fit a Logistic Regression model using the scores from BAARD on a hold-out training set to distinguish adversarial examples from legitimate inputs.

While being fast and memory-efficient, this approach has two issues when dealing with image data. 1. When the feature space is sparse, the S1 score becomes noise-sensitive. 2. The score varies under transformations, such as trans-

---

**Algorithm 2** BAARD Stage 3 – Decidability Stage

---

**Input:**  $\mathbf{x}$ : unlabeled example,  $\hat{y}$ : its prediction,  $(X, Y)$ : training set,  $k_{S3}$ : number of neighbors, and  $m_{S3}$ : sample size.

**Output:**  $S3\_score$

- 1:  $S \leftarrow$  Random sampling  $\{\mathbf{x}_1, \dots, \mathbf{x}_{m_{S3}}\}$  where  $\mathbf{x}_i \in X$
  - 2:  $D(\mathbf{x}, S) \leftarrow$  Compute angular distances between example  $\mathbf{x}$  and subset  $S$
  - 3:  $X' \leftarrow \text{top\_k}(D(\mathbf{x}, S), k_{S3})$   $\triangleright$  Find top  $k$ -nearest neighbors.
  - 4:  $P' \leftarrow \text{Softmax}(f(X'))$   $\triangleright$  Compute probability estimates for neighbors.
  - 5:  $\mu_{P'}, \sigma_{P'} \leftarrow \text{mean}(P'), \text{std}(P')$   $\triangleright$  Compute mean and standard deviation vectors.
  - 6:  $\mathbf{z} \leftarrow \left| \frac{\text{Softmax}(f(\mathbf{x})) - \mu_{P'}}{\sigma_{P'}} \right|$
  - 7: **return**  $\mathbf{z}_{\hat{y}}$   $\triangleright \mathbf{z}_{\hat{y}}$  is the value of  $\mathbf{z}$  index at  $\hat{y}$ .
- 

lation and rotation. Images are commonly modeled by convolutional neural networks, because the convolutional layers can learn internal representation in a two-dimensional space. Hence, these latent outputs represent the extracted feature space learned by the model. We overcome the above issues by using the latent outputs after the convolutional layers but before the fully connected layer. Note that tabular data does not suffer from the same issues. Moreover, anything related to the training data can be calculated beforehand to speed up the algorithm at inference time.

## 4 Experiments

We evaluate BAARD by analyzing its parameters, deconstructing it, and testing it against attacks in both white-box and gray-box settings. We repeat the experiments five times to ensure robustness. To ensure reproducibility, all data, pre-trained classifiers, hyperparameters, additional results, and code are available at <https://github.com/changx03/baard>.

### 4.1 Experimental Setup

**Data and Classifiers.** We test BAARD on both image and tabular data. We acquire MNIST and CIFAR10 with default train-test split from PyTorch for image datasets. We use the model from Carlini and Wagner [5] for MNIST and ResNet18 from PyTorch for CIFAR10. The pre-trained models are available in our repository. We remove the misclassified examples and sample 1000 images for generating adversarial examples and another 1000 for validating the detectors from the test set. We acquire all tabular data from the UCI ML repository <sup>4</sup>. All tabular data use a 60-20-20 split. The SVM and *Decision Trees* (DT) models for tabular data use the default parameters. Additional datasets are tested and included in our repository.

<sup>4</sup> Source: <https://archive.ics.uci.edu/ml>

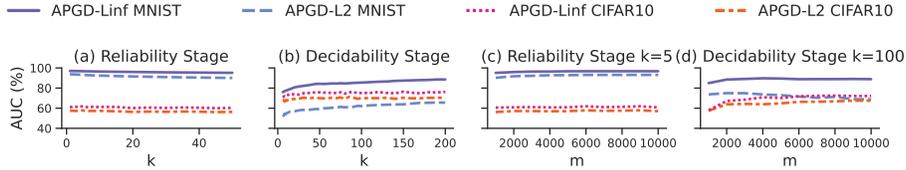


Fig. 2: Tuning hyperparameters for BAARD at the minimal adversarial perturbation. We first search for the optimal  $k$ , then tune the sample size  $m$ .

**Attack Algorithms.** We evaluated BAARD and other detectors under various attacks that are covered in Sec. 2, including PGD [15], APGD [6], CW- $L_2$  [5], and DTA [18]. We additionally include the results for FGSM [8], Boundary Attack [2], and DeepFool [16] in our repository. We define the adversary’s goal to have examples misclassified as any class except the true one so all attacks are untargeted. To test attack strengths, when there are multiple L-norm constraints, we test both  $L_\infty$  and  $L_2$  norm constraints. For each attack, we have considered a wide range of attack strengths. For instance, the parameter  $\epsilon$  in APGD controls the amount of perturbation allowed [6]. We set the minimal value to where the attack has at least 95% success rate. The minimal  $\epsilon$  for APGD is set to 0.22 and 4.0 for  $L_\infty$  and  $L_2$  on MNIST, 0.01 and 0.3 on CIFAR10, respectively. In Table 1, these values are used as the “Low”  $\epsilon$ , and the “High” is set to at least double the “Low” where there is a visible artifact on the example, but the legitimate label is still recognizable.

**Evaluation Metrics.** We report the *Area Under the Curve* (AUC) of the *Receiver Operating Characteristic* (ROC) curve as the performance metric. In practice, a single threshold may be selected based on the *False Positive Rate* (FPR). Hence, we also report TPRs when thresholds are chosen based on 5% FPRs (TPR@5FPR) when comparing different detectors.

## 4.2 Detection Results

**Parameter Analysis.** We treat each stage as an individual detector when tuning the hyperparameters. Since each stage’s performance directly links to  $k$  and the sample size  $m$  is for speeding up the algorithm, we first find the optimal  $k$  while using the entire training set, then use the optimal  $k$  to tune  $m$ .

The values of  $k_{S_2}$  and  $k_{S_3}$  are different. As shown in Fig. 2,  $k_{S_2}$  in the Reliability Stage becomes stable after the initial fluctuation. Reliability prefers a smaller  $k_{S_2}$  value, as it checks the closest representation of  $\mathbf{x}$  in the training samples with the same label as  $\hat{y}$ . Because Decidability finds neighbors from all training samples, a greater value of  $k_{S_3}$  is preferred. Once  $k_{S_2}$  and  $k_{S_3}$  are chosen, the optimal  $m_{S_2}$  and  $m_{S_3}$  should be the minimum value while maintaining the detector’s performance. Because the Reliability Stage uses the in-class training subset, the possible sample size is smaller than the Decidability Stage. Our

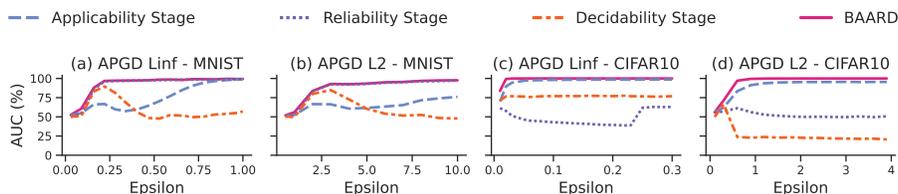


Fig. 3: BAARD’s performance under decomposition against adversarial attacks with a full range of perturbations.

results show that the detector’s performance is sturdy after initial turbulence, suggesting that the sub-sampling has minimal impact on the overall performance. The experiment concludes that BAARD requires minimal tuning. We set  $k$  to 5 and 100 and  $m$  to 1000 and 5000 for the Reliability and Decidability Stages respectively for all image datasets.

**Ablation Study.** We decompose BAARD to investigate how each stage contributes to the overall performance. Fig. 3 shows AUCs at various adversarial perturbations. Since attacks under a  $L_\infty$  constraint result in a significant deviation on the feature space [23], we find neither the S1 nor S2 score alone can detect such attacks. In Fig. 3d, the Decidability Stage’s AUC (orange dotted line) goes lower than 50% when  $\epsilon \geq 0.6$ , indicating that the correlation between the S3 score and the detector’s performance flip when  $\epsilon$  increases. It means the classifier becomes more confident with the misclassified predictions when  $\epsilon$  increases, leading to smaller S3 scores. Meanwhile, the S1 score becomes larger since the attack makes significant changes to the input. A low AUC on one stage indicates that stage alone is insufficient as a detector. However, by combining all stages, the results show BAARD is effective on a wide range of adversarial perturbations.

**White-Box Evaluation.** We address the robustness of BAARD against adaptive white-box attacks. To simultaneously attack the classifier and BAARD, the attacks’ loss function is  $\mathcal{L}^* := \mathcal{L} + \mathcal{L}_{S1} + \mathcal{L}_{S2} + \mathcal{L}_{S3}$ , where  $\mathcal{L}$  is the term for the evasion attack:  $\mathcal{L} := -\text{CrossEntropy}(f(\mathbf{x}'), y_{\text{target}})$ , and the rest of the terms are the losses for each stage. Because none of the stages are differentiable, a common approach is to apply gradient approximation [11]. Tramer *et al.* [21] pointed out that gradient approximation tends to fail when the loss function includes multiple indifferentiable terms, a more robust approach is



Fig. 4: Apply our Adaptive White-box Targeted  $L_2$  Attack to CIFAR10; When extreme parameters are used, it transforms a benign example into the target.

When extreme parameters are used, it transforms a benign example into the target.

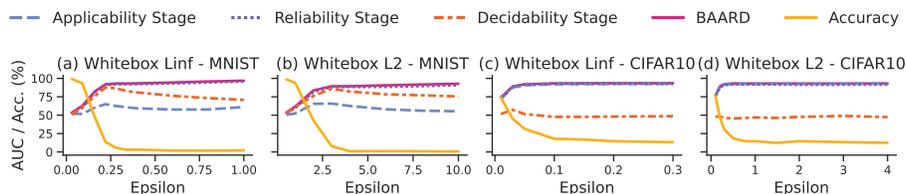


Fig. 5: BAARD’s performance against Adaptive White-box Targeted attacks. The accuracy indicates the classifier’s performance under such attacks.

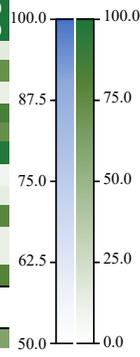
to find a target  $\mathbf{x}_{\text{target}}$  that can pass the detector and use it as a reference. Hence, we propose an *Adaptive White-box Targeted* (AWT) attack as follows: we find the nearest neighbor from the training data based on the same feature space BAARD uses, as  $\mathbf{x}_{\text{target}}$ , and then minimize the difference between  $\mathbf{x}$  and  $\mathbf{x}_{\text{target}}$  to bypass S1 and S2. To avoid  $f(\cdot)$  making over confident predictions, we use  $f(\mathbf{x}_{\text{target}})$  as a reference to bypass S3. The new loss function becomes  $\mathcal{L}^* := -\ell(f(\mathbf{x}), f(\mathbf{x}_{\text{target}})) - c\ell(\mathbf{x}, \mathbf{x}_{\text{target}})$ , where both terms use the *Mean Squared Error* (MSE) loss and the hyperparameter  $c$  controls the ratio on how much  $\mathbf{x}$  moves toward to  $\mathbf{x}_{\text{target}}$ . As shown in Fig. 4, if we relax the perturbation constraint  $\epsilon$  and dial  $c$  to an extreme, the adversarial example becomes indistinguishable from the target.

We present the evaluation of BAARD against our AWT attacks in Fig. 5 with  $c$  set to 1. The attack can successfully deceive the classifier and bypass S1 or S3, but not all stages. Previous works show similar algorithms are effective on detectors with multiple loss functions, such as the Odds detector [21]. We find such attacks are ineffective on BAARD, as three stages work together, which are robust against AWT attacks under both  $L_2$  and  $L_\infty$  constraints.

**Gray-Box Benchmark.** To benchmark the performance of BAARD against other detectors in Sec. 2, we use the same hold-out set to train logistic regression models for each dataset based on the features extracted from the detector. Table 1 presents both the AUC and TPR values obtained by varying the threshold of the regressors’ outputs. BAARD performs consistently well across different classifiers under attacks with various strengths, showing outstanding performance on attacks with high perturbations. One outlier is the APGD attack with an  $L_2$  constraint at a low  $\epsilon$  on CIFAR10, where most detectors are weak, except FS and Odds. However, FS and Odds are tuned explicitly for low perturbations and completely fail to detect attacks with high perturbations. RC can apply to any classifier in theory, but it only performs well in CW2. Meanwhile, the detectors tailored to images and neural networks cannot apply to SVM and DT classifiers. No detector performs reliably on the PGD attacks on the Breast Cancer dataset. However, BAARD is substantially faster than detectors with similar performance, such as LID, Odds, and ML-LOO. In conclusion, BAARD is the most versatile detector tested that can reliably detect adversarial examples with various constraints on a wide range of classifiers.

Table 1: Performance of detectors. The AUC scores (%) on the left are computed from logistic regression. The right side shows the corresponding TPR at 5% FPR. “Low” and “High” indicate perturbations allowed for the attack.

Image Data	Attack Perturbation Detector	AUC-ROC (%)						TPR@FPR5 (%)			
		APGDinf		APGD2		CW2	APGDinf		APGD2		CW2
		Low	High	Low	High		Low	High	Low	High	
MNIST (CNN)	RC	72.9	51.6	50.8	51.5	99.9	49.0	0.0	0.0	0.0	100.0
	FS	99.7	74.0	73.7	68.1	100.0	99.3	1.7	26.1	3.3	100.0
	LID	60.6	98.1	43.4	80.8	62.8	17.5	91.8	6.6	43.3	14.3
	Odds	98.9	99.7	96.5	96.4	95.7	97.8	100.0	81.4	79.5	76.5
	ML-LOO	99.8	100.0	93.2	100.0	60.0	99.2	100.0	70.9	100.0	10.8
	PN	89.7	62.1	55.3	54.3	97.1	64.9	7.0	9.7	4.9	89.0
	BAARD	97.0	98.4	92.8	96.8	96.0	84.4	92.8	61.2	82.6	77.0
CIFAR10 (ResNet18)	RC	49.6	54.8	55.9	54.7	99.4	4.7	0.0	12.1	0.0	98.5
	FS	95.7	70.7	95.1	82.5	90.7	75.7	24.7	78.8	42.5	6.9
	LID	82.2	99.2	63.4	98.7	40.7	44.0	96.5	22.3	94.2	13.0
	Odds	98.0	67.4	97.2	80.9	96.1	95.2	0.2	95.5	2.6	83.1
	ML-LOO	67.0	99.6	58.6	99.2	66.4	28.0	98.9	16.8	97.5	10.8
	PN	76.6	54.2	75.4	58.1	66.8	18.1	7.5	17.0	9.7	10.8
	BAARD	81.6	100.0	70.2	99.2	89.0	35.4	100.0	16.4	96.1	85.2
Tabular Data	Attack (Model)	PGDinf (SVM)			DTA (DT)		PGDinf (SVM)			DTA (DT)	
	Banknote	RC	79.7	99.0		86.4	46.0	100.0		63.4	
	FS, LID, etc.	-	-		-	-	-		-		
	BAARD	96.5	100.0		95.9	87.0	100.0		89.9		
Breast Cancer	RC	65.0	75.2		97.2	0.0	0.0		82.6		
	FS, LID, etc.	-	-		-	-	-		-		
	BAARD	77.7	52.0		96.8	21.8	7.6		85.8		



## 5 Conclusion and Future Work

In this paper, we connected two previously unlinked domains: the Applicability Domain (AD) in cheminformatics and adversarial detection in machine learning. By sharing solutions to similar problems, both areas can benefit. We proposed BAARD, a novel adversarial detection framework inspired by AD. Our experiments showed its robustness against various adversarial evasion attacks, including those with strong perturbations. BAARD is portable and versatile enough to work with any classifier, removing the need for redesigning a defense. Our framework overcomes challenging issues in the field while maintaining comparable performance. In future research, we will explore how the insights we have gained from adversarial detection can be transferred into cheminformatics.

**Acknowledgements.** The authors wish to acknowledge the use of New Zealand eScience Infrastructure (NeSI) national facilities - <https://www.nesi.org.nz>.

## References

1. Alvarsson, J., McShane, S.A., Norinder, U., Spjuth, O.: Predicting with confidence: using conformal prediction in drug discovery. *J. Pharm. Sci.* **110**(1), 42–49 (2021)
2. Brendel, W., Rauber, J., Bethge, M.: Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. In: *ICLR* (2018)

3. Cao, X., Gong, N.Z.: Mitigating evasion attacks to deep neural networks via region-based classification. In: ACSAC. pp. 278–287 (2017)
4. Carlini, N., Wagner, D.: Adversarial examples are not easily detected: Bypassing ten detection methods. In: AISec. pp. 3–14 (2017)
5. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: IEEE SSP. pp. 39–57 (2017)
6. Croce, F., Hein, M.: Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In: ICML. pp. 2206–2216. PMLR (2020)
7. Demontis, A., Melis, M., Pintor, M., Jagielski, M., Biggio, B., Oprea, A., Nita-Rotaru, C., Roli, F.: Why do adversarial attacks transfer? explaining transferability of evasion and poisoning attacks. In: USENIX Security. pp. 321–338 (2019)
8. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: ICLR (2015)
9. Hanser, T., Barber, C., Marchaland, J., Werner, S.: Applicability domain: towards a more formal definition. SAR and QSAR in Environmental Research **27**(11), 865–881 (2016)
10. He, W., Wei, J., Chen, X., Carlini, N., Song, D.: Adversarial example defenses: ensembles of weak defenses are not strong. In: USENIX WOOT. pp. 15–15 (2017)
11. Hu, S., Yu, T., Guo, C., Chao, W.L., Weinberger, K.Q.: A new defense against adversarial images: Turning a weakness into a strength. NIPS **32** (2019)
12. Kloukinotis, A., Papandreou, A., Lalos, A., Kapsalas, P., Nguyen, D.V., Moustakas, K.: Countering adversarial attacks on autonomous vehicles using denoising techniques: A review. IEEE OJ-ITS (2022)
13. Luo, W., Wu, C., Ni, L., Zhou, N., Zhang, Z.: Detecting adversarial examples by positive and negative representations. ASC **117**, 108383 (2022)
14. Ma, X., Li, B., Wang, Y., Erfani, S.M., Wijewickrema, S., Schoenebeck, G., Song, D., Houle, M.E., Bailey, J.: Characterizing adversarial subspaces using local intrinsic dimensionality. In: ICLR (2018)
15. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: ICLR (2018)
16. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: A simple and accurate method to fool deep neural networks. In: CVPR. pp. 2574–2582. IEEE (2016)
17. Netzeva, T.I., Worth, A.P., Aldenberg, T., Benigni, R., Cronin, M.T., Gramatica, P., Jaworska, J.S., Kahn, S., Klopman, G., Marchant, C.A., et al.: Current status of methods for defining the applicability domain of (quantitative) structure-activity relationships: The report and recommendations of ecvam workshop 52. ATLA **33**(2), 155–173 (2005)
18. Papernot, N., McDaniel, P., Goodfellow, I.: Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. arXiv preprint arXiv:1605.07277 (2016)
19. Roth, K., Kilcher, Y., Hofmann, T.: The odds are odd: A statistical test for detecting adversarial examples. In: ICML. pp. 5498–5507. PMLR (2019)
20. Tramer, F.: Detecting adversarial examples is (nearly) as hard as classifying them. In: ICML. pp. 21692–21702. PMLR (2022)
21. Tramer, F., Carlini, N., Brendel, W., Madry, A.: On adaptive attacks to adversarial example defenses. NIPS **33**, 1633–1645 (2020)
22. Xu, W., Evans, D., Qi, Y.: Feature squeezing: Detecting adversarial examples in deep neural networks. arXiv preprint arXiv:1704.01155 (2017)
23. Yang, P., Chen, J., Hsieh, C.J., Wang, J.L., Jordan, M.: MI-loo: Detecting adversarial examples with feature attribution. In: AAAI. vol. 34, pp. 6639–6647 (2020)