

# Distributed Adversarial Attacks

Sander Prenen

Thesis voorgedragen tot het behalen  
van de graad van Master of Science  
in de ingenieurswetenschappen:  
computerwetenschappen, hoofdoptie  
Artificiële intelligentie

**Promotoren:**

Prof. dr. ir. W. Joosen  
Dr. ir. D. Preuveneers

**Assessoren:**

Ir. W. Eetveel  
W. Eetrest

**Begeleiders:**

V. Rimmer  
I. Tsingenopoulos

© Copyright KU Leuven

Without written permission of the thesis supervisors and the author it is forbidden to reproduce or adapt in any form or by any means any part of this publication. Requests for obtaining the right to reproduce or utilize parts of this publication should be addressed to the Departement Computerwetenschappen, Celestijnenlaan 200A bus 2402, B-3001 Heverlee, +32-16-327700 or by email [info@cs.kuleuven.be](mailto:info@cs.kuleuven.be).

A written permission of the thesis supervisors is also required to use the methods, products, schematics and programmes described in this work for industrial or commercial use, and for submitting this publication in scientific contests.

Zonder voorafgaande schriftelijke toestemming van zowel de promotoren als de auteur is overnemen, kopiëren, gebruiken of realiseren van deze uitgave of gedeelten ervan verboden. Voor aanvragen tot of informatie i.v.m. het overnemen en/of gebruik en/of realisatie van gedeelten uit deze publicatie, wend u tot het Departement Computerwetenschappen, Celestijnenlaan 200A bus 2402, B-3001 Heverlee, +32-16-327700 of via e-mail [info@cs.kuleuven.be](mailto:info@cs.kuleuven.be).

Voorafgaande schriftelijke toestemming van de promotoren is eveneens vereist voor het aanwenden van de in deze masterproef beschreven (originele) methoden, producten, schakelingen en programma's voor industrieel of commercieel nut en voor de inzending van deze publicatie ter deelname aan wetenschappelijke prijzen of wedstrijden.

# Preface

*Sander Prenen*

# Contents

<b>Preface</b>	<b>i</b>
<b>Contents</b>	<b>ii</b>
<b>List of Figures</b>	<b>iv</b>
<b>List of Tables</b>	<b>v</b>
<b>List of Abbreviations</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Background</b>	<b>3</b>
2.1 Neural networks . . . . .	3
2.2 Adversarial attacks . . . . .	6
2.3 Particle swarm optimization . . . . .	9
<b>3 Related work</b>	<b>13</b>
3.1 Boundary attack . . . . .	13
3.2 HopSkipJumpAttack . . . . .	15
3.3 SurFree attack . . . . .	17
3.4 Stateful defense . . . . .	18
3.5 PSO and distributed attacks . . . . .	20
<b>4 Approach</b>	<b>21</b>
4.1 Distribution . . . . .	21
4.2 Optimization . . . . .	22
4.3 Approach . . . . .	22
4.4 Threat model . . . . .	22
<b>5 Evaluation</b>	<b>25</b>
5.1 Evaluation protocol . . . . .	25
5.2 Determining the baseline . . . . .	26
5.3 Applying PSO to BBA . . . . .	27
5.4 Towards distribution of the attack . . . . .	29
5.5 Throwing the defense off the scent . . . . .	32
5.6 Optimizing the attack . . . . .	32
<b>6 Conclusion</b>	<b>35</b>
<b>A SurFree calculations</b>	<b>39</b>

<b>Bibliography</b>	<b>41</b>
---------------------	-----------

# List of Figures

2.1	Linear separability . . . . .	5
2.2	Decision boundaries and decision regions . . . . .	5
2.3	Activation functions . . . . .	5
2.4	Example of a convolution layer . . . . .	6
2.5	Difference targeted and untargeted attack . . . . .	8
2.6	Adversarial training . . . . .	10
2.7	Particle swarm optimization . . . . .	12
3.1	Difference between noise patterns . . . . .	14
3.2	Influence of frequency on Perlin noise . . . . .	15
3.3	Intuition of the Boundary Attack . . . . .	16
3.4	Intuition of the HopSkipJumpAttack . . . . .	17
3.5	Geometric configuration of SurFree . . . . .	19
5.1	Some examples of the MNIST dataset . . . . .	26
5.2	Some examples of the CIFAR-10 dataset . . . . .	26
5.3	Intuition of multiple starting points . . . . .	30
5.4	Detections of the different attacks . . . . .	30
5.5	Schematic overview of the query submission distribution . . . . .	32
5.6	Overview of the round-robin based distribution scheme . . . . .	33

# List of Tables

5.1	Model architectures for the MNIST and CIFAR model . . . . .	26
5.2	Baseline results . . . . .	27
5.3	PSO-BBA results . . . . .	29
5.4	Distributed PSO-BBA results . . . . .	33





# List of Abbreviations

- AI** Artificial Intelligence. 3
- ANN** Artificial Neural Network. 3, 4, 6
- API** Application Programming Interface. 7, 22, 23, 31
- BA** Boundary Attack. 13, 14, 15, 16
- BBA** Biased Boundary Attack. 14, 16, 20, 22, 26, 27, 28, vii, 29, 30, 33
- CNN** Convolutional Neural Network. 3, 4
- DB** Distance-Based. 32, 33
- DDoS** Distributed Denial of Service. 20
- DkNN** Deep k-Nearest Neighbors. 9
- DNN** Deep Neural Network. 3
- EA** Evolutionary Algorithm. 9, 20
- EDB** Embedded-Distance-Based. 32, 33
- FGSM** Fast Gradient Sign Method. 7
- HSJA** HopSkipJumpAttack. 15, 16, 17
- MRR** Modified Round-Robin. 31, 33
- PSO** Particle Swarm Optimization. 9, 10, 11, 20, 22, 27, 28, vii, 29, 30, 31, 33
- ReLU** Rectified Linear Unit. 3, 5, 26
- RNN** Recurrent Neural Network. 3
- RR** Round-Robin. 31, 32, 33



## Chapter 1

# Introduction



## Chapter 2

# Background

### 2.1 Neural networks

Ever since the invention of computer systems, it has always been a goal of scientists and engineers to create **Artificial Intelligence (AI)**. Current state of the art approaches are mimicking the human brain, more specifically the neurons inside the brain. Already in the fifties, Rosenblatt introduced his perceptron [1]. The perceptron is a single neuron able to learn linearly separable patterns. It does so by finding a hyperplane that separates the two classes. This hyperplane is called the decision surface or decision boundary and the perceptron itself is called a classifier. Geometric regions separated by a decision boundary are called decision regions. The concept of linear separability is explained in Figure 2.1 in two dimensions. In Figure 2.2, the decision boundaries and decision regions are explained visually.

Unfortunately not all patterns are linearly separable. To overcome this problem, the neurons can be layered, creating an **Artificial Neural Network (ANN)** in the process. Layering neurons sequentially is essentially a linear combination of neurons. This in itself does not create non-linear decision surfaces. Non-linear activation functions are added for the **ANN** to be able to learn more complex decision boundaries. Some commonly used activation functions are **Rectified Linear Unit (ReLU)** [2], Heaviside step function [3] and softmax [4] (or sigmoid when used on scalars). In Figure 2.3 the plots of some activation functions can be found.

The neurons can be combined in different ways to create different **ANN** architectures. Each architecture has its own strengths and weaknesses. **Convolutional Neural Networks (CNNs)** excel in classifying visual data [5, 6], whilst **Recurrent Neural Networks (RNNs)** are widely used when there exist dependencies inside the data, such as in speech recognition [7, 8] or time series prediction [9]. More recent research focuses on **Deep Neural Networks (DNNs)**, due to the ever increasing computational power available. **DNN** approaches are able to compare to and even surpass human performance on very specific tasks [10, 11].

## 2. BACKGROUND

---

Most **ANNs** are not built from individual neurons, but rather from layers of neurons. An **ANN** architecture describes the different layers of which the network consists. The most basic type of layer is the fully connected or dense layer. This layer, as the name suggests, consists of neurons that are connected to all neurons of the previous layer. Every neuron in this layer outputs a value based on the perceptron rule:

$$y = b + W \cdot X$$

Here  $y$  is the output of the neuron,  $b$  is the bias,  $W$  is the weight matrix and  $X$  is the vector consisting of the outputs of the previous layer. The values of the bias and weights are learned based on the training data.  $W \cdot X$  is the dot product  $\sum_{i=1}^n w_i x_i$ , where  $n$  is the number of neurons in the previous layer [12].

Activation layers apply activation functions to their inputs. They can be standalone layers, but more frequently they are integrated in other layers.

Convolution layers [13] are the main building blocks of **CNNs**. Convolution layers map inputs to so called feature maps. A filter  $W$  slides over the input  $X$  from left to right top to bottom. At every position, the dot product  $W \cdot X$  of the frame and the underlying values of the input is computed. This value is placed at the corresponding position of the feature map. Higher values of the dot product (for a fixed value of  $\|W\|_2$ ) mean that  $W$  and the underlying values of  $X$  are more similar. The feature map therefore indicates where in the input  $X$  a certain pattern  $W$  can be found. This process is translation-invariant, meaning that a certain pattern can be detected at any location of  $X$ . Scale and rotation invariance can be achieved by repeating the process for scaled and rotated filters. This is the main reason why **CNNs** are widely used for visual data. An example of a convolution layer is shown in Figure 2.4.

The use of convolution layers tends to cause an explosion in the dimensionality of the data. One filter will produce a feature map that is smaller than the original input. But a convolution layer typically consists of many filters, causing the explosion in the number of connections and weights. Pooling or subsampling layers [13] alleviate this problem by aggregating nearby points. The aggregation happens using a sliding window similar to a convolution layer. This approach to reduce the dimensionality has the added benefit that the output of the model is less sensitive to the exact location of a pattern. Some commonly used pooling layers are the MaxPooling and AveragePooling layers. They aggregate nearby points using the maximum and average value respectively. Convolution and pooling layers are often used in combinations with each other and are repeated multiple times. This allows the **CNN** to detect patterns at different scales.

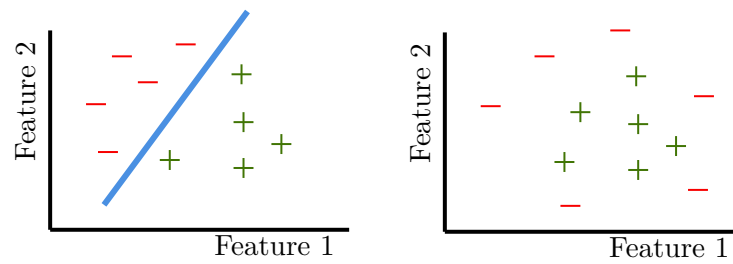


FIGURE 2.1: Linearly separable classes on the left and non-linearly separable classes on the right. Two classes are linearly separable if there exists a hyperplane for which all examples of one class are on the same side of this hyperplane, whilst all examples of the other class are on the other side of the hyperplane. In two dimensions, the hyperplane is a straight line.

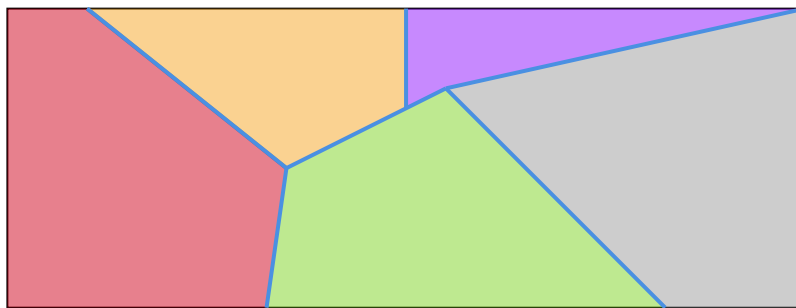


FIGURE 2.2: Decision boundaries (blue lines) separate different decision regions (colored regions).

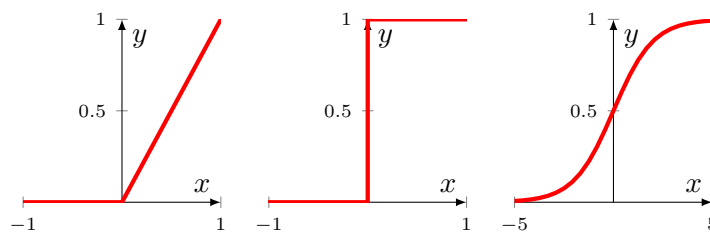


FIGURE 2.3: Plots of different activation functions. From left to right: ReLU, Heaviside step and sigmoid.

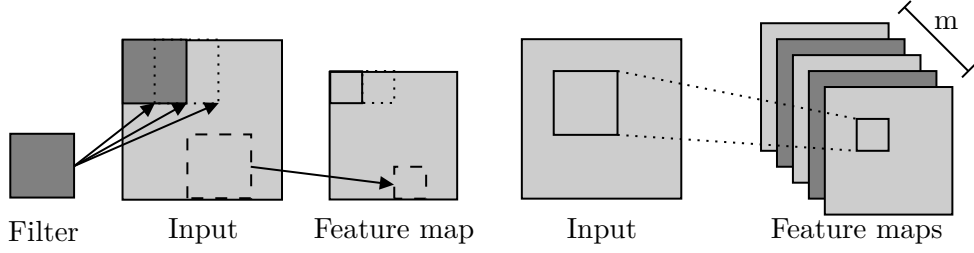


FIGURE 2.4: Example of a convolution layer. The filter slides over the input from left to right top to bottom. At every position, the dot product of the filter and input image is computed. This value is added to the corresponding position of the feature map. This process is repeated for  $m$  different filters leading to  $m$  feature maps.

## 2.2 Adversarial attacks

The expressiveness of ANNs is a double-edged sword. It is the cause for the near-human performance on some tasks, but also for counter-intuitive properties. As studied by Szegedy et al [14], one of these properties is the presence of discontinuous decision boundaries. This might cause seemingly identically images to be classified differently. They first defined adversarial examples as *imperceptibly small perturbations to a correctly classified input image, so that it is no longer classified correctly* [14]. This property of ANNs might not seem important at first glance, but it can be quite worrisome from a security point-of-view. Malicious users could craft images to bypass face recognition software [15] or attack the camera of a self-driving car to misclassify traffic signs [16]. Other fields where adversarial examples are of interest include malware detection [17], natural language processing [18] and industrial control systems [19]. Adversarial attacks are algorithms used to craft such adversarial examples.

All adversarial attacks are evaluated against a threat model. A threat model is *a structured representation of all the information that affects the security of an application* [20]. This information consists of the goals, knowledge and capabilities of the attacker, the accessibility of the model under attack and the costs of (un)successful attacks.

Most research on adversarial attacks is done using images. Researchers have the most freedom in this domain, since a slightly altered image is still an image with roughly the same contents. Slightly modifying an industrial control system however, might break the entire way the system works. Research in other domains is mostly conducted by altering existing image algorithms to the specific use case. For this reason this work only focuses on adversarial attacks on images.



### 2.2.1 Adversarial attacks terminology

Adversarial attacks are generally divided in two categories, white box attacks and black box attacks. In a white box attack, the attacker has complete knowledge of the classifier under attack. This knowledge consists of the architecture, parameters and thus their gradients and all output of the classifier. Examples of white box attacks are the **Fast Gradient Sign Method (FGSM)** [21] or the Carlini & Wagner attack [22].

In black box attacks, the only thing the attacker has access to is the output of the model. Depending on the literature, this output consists of class labels only (decision-based attack) or class labels and the corresponding confidence scores (score-based attacks). Black box attacks are more relevant in real-life scenarios, since most attacks are performed on a third-party **Application Programming Interface (API)**. These **APIs** generally do not reveal the underlying model.

Transfer attacks [23] try to overcome this hurdle by creating a surrogate model. This is an undefended model similar to the model under attack. This idea is based on the observation that adversarial perturbations often are transferable to other models [21]. Attacks can leverage information (such as gradients) from the surrogate model to breach the black box model. Due to the transferability of adversarial examples, it is also possible to perform so called zero-query attacks. Zero-query attacks are performed entirely on the surrogate model and the resulting adversarial example is forwarded to the black box model.

Both white box and black box attacks can be divided into targeted and untargeted attacks depending on their goal. In a targeted attack, the goal of the attacker is to create an adversarial example with a specific target class. In an untargeted attack the target class can be any class. Untargeted variants of attacks generally enjoy much more freedom and are therefore able to craft adversarial examples that are closer to the original. Figure 2.5 visually explains the difference between the two types.

What does it mean for images to be close to each other? This is easy to visualize in two dimensions as in Figure 2.5, but in higher dimensions, this is more difficult. Images reside in  $d$ -dimensional space, where  $d$  is the amount of pixels of the image. Two commonly used distances in higher dimensions are the  $L_2$ -distance and the  $L_\infty$ -distance. They are defined as follows:

$$L_2(X, Y) = \sqrt{\sum_{i=1}^d |x_i - y_i|^2}$$

$$L_\infty(X, Y) = \lim_{p \rightarrow \infty} \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{1/p}$$

$$= \max(|x_1 - y_1|, |x_2 - y_2|, \dots, |x_d - y_d|)$$

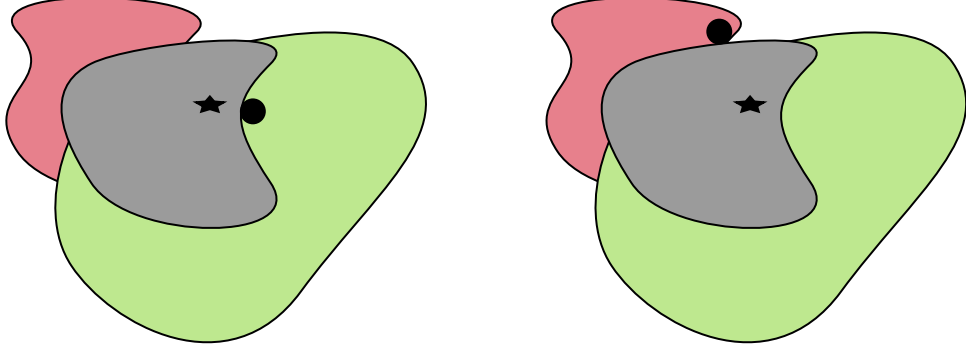


FIGURE 2.5: Different decision regions are shown in different colors. An adversarial example is being created starting from the source image (black star). On the left an untargeted attack is performed. The adversarial example is the image closest to the source image, that is classified differently (black circle). On the right a targeted attack is performed with the red decision region being the target class. The adversarial example is the image closest to the source image that is in the red decision region.

In both distances  $X$  and  $Y$  represent the images and  $(x_1, x_2, \dots, x_i, \dots, x_d)$  and  $(y_1, y_2, \dots, y_i, \dots, y_d)$  are the pixel values of  $X$  and  $Y$  respectively. The  $L_2$ -distance is also known as the Euclidean distance, which is a generalization of the Pythagorean theorem in more than two dimensions. It takes the pairwise distances between all pixels into account. The  $L_\infty$ -distance is also called the Chebyshev distance. This distance only depends on the maximal pairwise distance between the two images. By minimizing the  $L_\infty$ -distance, the maximal pixelwise difference is minimized [24]. Sometimes the norm notation  $\|\cdot\|$  is used instead of the distance notation. A distance is induced by a norm by  $d(X, Y) = \|X - Y\|$ . The other direction, a distance inducing a norm, is not as trivial. In order for a distance to induce a norm, the distance has to satisfy the following condition:  $d(\alpha X, \alpha Y) = |\alpha|d(X, Y)$ . This condition is satisfied in most commonly used distances, but not always. The discrete distance [25] is an example of such a distance and is defined as follows:

$$d(X, Y) = \begin{cases} 1, & \text{if } X \neq Y \\ 0, & \text{else} \end{cases}$$

### 2.2.2 Adversarial defenses

The existence of adversarial attacks naturally gave rise to adversarial defenses. These defenses can be categorized based on their objective. They can be either proactive or reactive. The goal of a proactive defense consists of making the models under attack more robust, while reactive defenses aim to identify attacks before they reach the model [26]. Different defensive countermeasures can be taken in each cat-

egory. The remainder of this section will discuss some commonly used techniques [27].

**Gradient masking** techniques hinder optimization-based attacks by having gradients "*that are not useful*" [28]. They are also sometimes referred to as obfuscated gradients [29]. Three types of obfuscated gradients can be identified. Shattered gradients introduce incorrect or non-existent gradients. Stochastic gradients are caused by random effects in the defense and exploding or vanishing gradients are primarily caused by chaining neural network evaluations. Besides intentionally introducing gradient masking in neural networks, they can also be introduced unintentionally due to the design of the network.

**Defensive distillation** [30] is a technique that can be classified as gradient masking. The goal of defensive distillation is to smooth the gradients of the model, making it more resilient to small input perturbations. The distillation procedure is done in two steps. First the model is trained on the original data and labels. This step produces a probabilistic output for each input due to the softmax activation function. Then the network is retrained using the original data and the probabilistic outputs as labels. The probabilistic labels contain additional knowledge that can be exploited to increase generalizability of the model.

**Adversarial training** [21] techniques inject adversarial examples in the training dataset and retrain the model in order to create a more resilient model. This is essentially a brute force method to correctly classify some adversarial examples. However this new model is still susceptible to new adversarial attacks, since the decision boundary has only been moved slightly. This can be seen in Figure 2.6.

**Preprocessing techniques** work on the inputs of a model. Different preprocessing techniques, such as denoising [31], dimensionality reduction [32] and image transformations [33] can be used to defend the model under attack. The goal of all techniques boils down to giving the attacker less control over the exact input that is being fed to the model.

Some defenses rely on **proximity measurements** between the input and the model. An example of this countermeasure is **Deep k-Nearest Neighbors (DkNN)** by Papernot and McDaniel [34]. **DkNN** computes support for a decision from a network based on a nearest neighbors search in the training data. Another example is region-based classification [35], where a prediction is made for a given input based on the proximity of training examples.

## 2.3 Particle swarm optimization

**Particle Swarm Optimization (PSO)** [36] is an optimization framework part of the **Evolutionary Algorithms (EAs)** family. In **EAs**, populations of candidate solutions

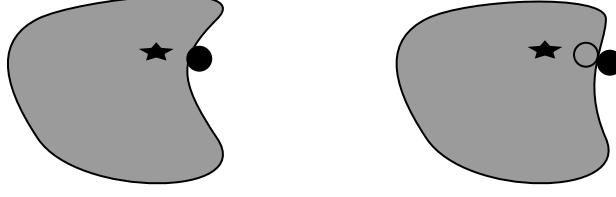


FIGURE 2.6: Decision boundaries before (left) and after (right) adversarial training. Before adversarial training, an adversarial example can be created (black circle). This example is added to the training set and the network is retrained. A new decision surface is created in the process. This surface classifies the previous adversarial example correctly (empty circle), but there is a new opportunity for an adversarial attack (black circle).

evolve based on mechanisms inspired by the field of biology, such as ant colonies [37], mutation and recombination [38]. The mechanism that inspired **PSO** is the behaviour of flocks of birds. The framework has been applied to numerous problems such as routing problems [39, 40], diagnosing diseases from imaging [41] and calculating heat transfer coefficients [42].

In **PSO**, different particles  $x$  move through the search space based on a set of rules. Their new position  $x_t$  is determined by their previous position  $x_{t-1}$  and a velocity  $v_t$ . The velocity depends on the distance to best position of the swarm  $g$  and the best known position of the particle  $p$ . The distances can be weighted by acceleration coefficients  $c$  to put more emphasis on exploration or exploitation. These values are multiplied by random parameters  $r$  uniformly distributed in  $[0, 1]$ . The velocity is also dependent on the previous velocity with a corresponding weight  $w$ . The best position is determined using a fitness function. This function states how ‘fit’ or good a certain position is with respect to the goal of the optimization problem. Equations 2.1 and 2.2 correspond to the update rules. Each particle  $i$  has its own position  $x_{t,i}$  and velocity  $v_{t,i}$  at time step  $t$ , but the  $i$  index has been omitted for readability. In Figure 2.7, the steps are graphically represented for a single particle.

$$v_t = \underbrace{wv_{t-1}}_{\text{Inertia}} + \underbrace{c_p r_p (p_{t-1} - x_{t-1})}_{\text{Individual best}} + \underbrace{c_g r_g (g_{t-1} - x_{t-1})}_{\text{Swarm best}} \quad (2.1)$$

$$x_t = x_{t-1} + v_t \quad (2.2)$$

Ever since the first mention of **PSO** in 1995, efforts have been made in order to improve the framework. The rest of this section will discuss some improvements.

The first version of **PSO** used a fixed value for the inertia weight  $w$ . Later it has empirically been shown that a linearly decaying weight improves performance [43]. This allows the algorithm to focus on exploration in the early iterations, while

shifting its focus to exploitation later on. The rate of decay depends on the high start value  $w_{start}$ , the lower end value  $w_{end}$  and the maximum number of iterations  $t_{max}$ . The weight in iteration  $t$  is calculated as follows:

$$w_t = w_{end} + (w_{start} - w_{end}) \left(1 - \frac{t}{t_{max}}\right) \quad (2.3)$$

In a large search space, particles can be far apart, which can cause the velocities to explode. By limiting the value of velocities to  $v_{max}$ , the swarm can be stabilized and the probability of finding an optimum is increased [36].

Clerc and Kennedy [44] studied PSO from a dynamic systems point of view. They provided a theoretically backed solution to the problem of the exploding velocities. Instead of limiting the value of the velocity, something which has only been shown empirically that it works, they proposed the constriction factor  $\chi$ . The constriction factor is able to prevent the velocity explosion, whilst avoiding premature convergence to local optima. The velocity update of equation 2.1 is altered as follows:

$$v_t = \chi(v_{t-1} + c_p r_p (p_{t-1} - x_{t-1}) + c_g r_g (g_{t-1} - x_{t-1}))$$

The inertia weight  $w$  is dropped and the entire sum is multiplied with  $\chi$ . The value of  $\chi$  is calculated as follows:

$$\chi = \begin{cases} \sqrt{\frac{2\kappa}{C-2+\sqrt{C^2-4C}}}, & C > 4 \\ \sqrt{\kappa}, & \text{else} \end{cases}$$

Here  $C$  is the sum of the acceleration coefficients  $c_p$  and  $c_g$  and  $\kappa$  is a user defined value to determine the rate of convergence.

TODO hybridization, multigroup

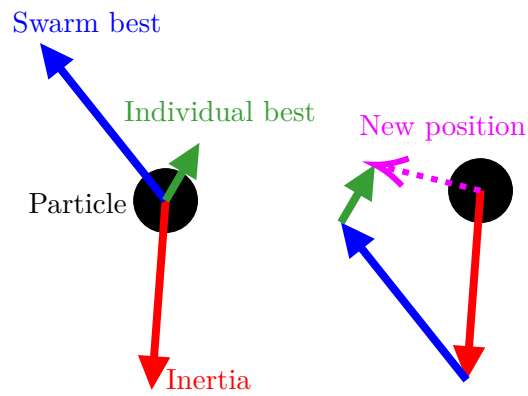


FIGURE 2.7: Particles move based on a step towards their best known position, the swarm's best known position and a step in the direction of the movement of the previous iteration. The different steps are combined to determine the new position of the particle.

## Chapter 3

# Related work

### 3.1 Boundary attack

**Boundary Attack (BA)** [45] is a decision-based adversarial attack. The basic intuition of **BA** differs from traditional adversarial attacks. Unlike these traditional adversarial attacks, where the original image is moved through search space in order to become adversarial, **BA** starts from an input that is already adversarial. This input is then moved closer to the original image, while staying adversarial.

The attack has to be initialized with an already adversarial input. Two different approaches can be taken depending on the attack setting. In the untargeted case, the input can be sampled from a maximum entropy distribution given the valid domain of this input. Samples that are not adversarial are rejected. An example of such a starting position can be seen in Figure 3.1a. In the case of a targeted attack, the input is a sample from the dataset that is classified as the target class by the model under attack.

**BA** iteratively updates the adversarial image by performing a step orthogonal to the original image and a step towards this image. In iteration  $k$ , a perturbation  $\eta_k$  is sampled from a uniform distribution. This perturbation is rescaled and added to the adversarial image. From this new position in search space, the step towards the original image is taken. This way the path of the attack follows the decision boundary, hence the name of the attack. The intuition of the **BA** is shown in Figure 3.3. The attack can only follow the boundary if the adversarial image is already near the boundary. The starting image is projected onto the boundary using binary search to ensure that the adversarial image is in the vicinity of the boundary.

The step sizes are adjusted according to local geometry of the boundary. The orthogonal step size  $\delta$  is adjusted so that approximately half of the orthogonal perturbations is still adversarial. This approach is based on trust region methods [46]. The step size towards the original image  $\epsilon$  is adjusted using the same principle, but here a user specified threshold is used. The decision boundary tends to become flatter, the closer to the original image the attack gets [47]. Therefore the algorithm converges

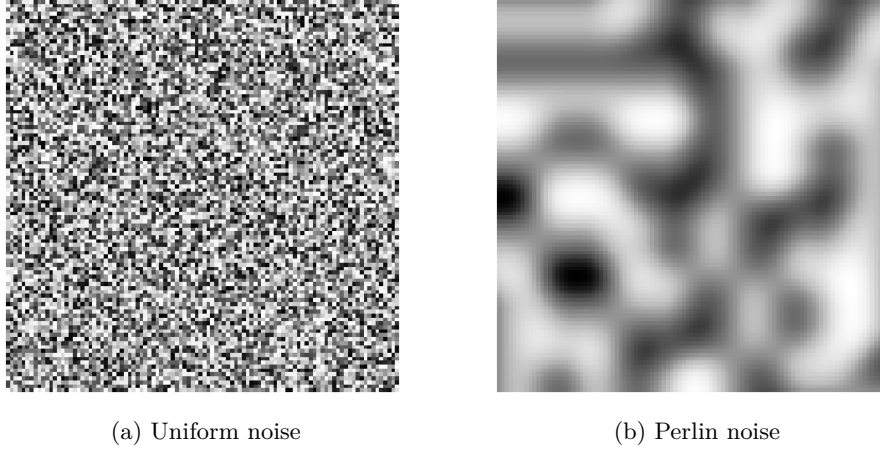


FIGURE 3.1: Difference between noise patterns.

when  $\epsilon$  converges to zero.

**Biased Boundary Attack (BBA)** [48] (previously known as *Boundary Attack++*) is an improvement on the original **BA** in three different ways. All three improvements will be discussed in order of the strength of their effect. The first improvement is a biased sampling technique. The key idea behind this is that most previous attacks yield adversarial examples with high frequencies in the image. By sampling the perturbations in the first step of the **BA** from a low frequency distribution, the frequency of the created adversarial example will be lowered as well. **BBA** does this by sampling from a Perlin noise [49] distribution instead of a uniform distribution. Lower frequency images yield more natural results and can more easily bypass simple preprocessing defense schemes. The difference between the two noise patterns can be seen in Figure 3.1. The noise patterns can be influenced by a frequency value. This value can be tuned depending on the size of the images at hand. Higher frequency values yield less smooth noise patterns. Figure 3.2 visually shows the influence of the frequency values.

The second improvement is to use a regional mask. The original **BA** applies a perturbation to the images as a whole. Every pixel will be perturbed with the same magnitude. This magnitude can be altered on a per-pixel basis when using a mask. Pixels that are further away from the target image will receive a larger perturbation than pixels that are already close to the corresponding pixel in the target. The mask  $m$  is constructed according to equation 3.1 based on the original image  $x_{orig}$  and the adversarial image  $x_{adv}$ . It is then pixel-wise applied to the sampled perturbation in equation 3.2. The masked perturbation is normalized afterwards.



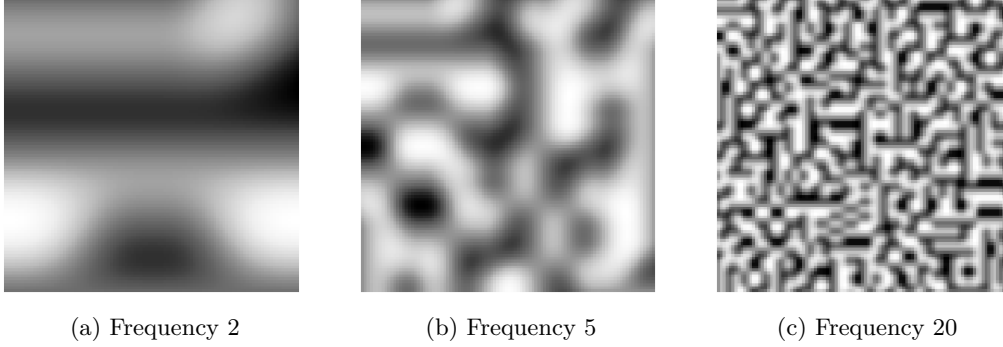


FIGURE 3.2: Influence of frequency on Perlin noise patterns.

$$m = |x_{adv} - x_{orig}| \quad (3.1)$$

$$\eta_k = m \odot \eta_k; \eta_k = \frac{\eta_k}{\|\eta_k\|} \quad (3.2)$$

This technique improves efficiency since the search space is significantly reduced. It is also possible to engineer masks for specific examples in order to incorporate other knowledge in the attack.

The final improvement is based on the idea of transfer attacks. A surrogate model is trained and will be used to calculate adversarial gradients. These gradients will then be used to bias the sampling direction for the orthogonal step. If the surrogate model does not closely resemble the defender, then the gradients will only hamper the speed of convergence of the attack instead of causing the attack to fail.

### 3.2 HopSkipJumpAttack

**HopSkipJumpAttack (HSJA)** [50], like **BA**, is a decision-based adversarial attack that starts from an adversarial input. The initial input is obtained in an identical manner as in **BA**. **HSJA** is an iterative algorithm that consists of three steps.

The first step is a projection onto the decision boundary of the model under attack. This projection is carried out using a binary search. The second step is to estimate the direction of the gradient at the boundary. Different directions are sampled from a uniform distribution over a  $d$ -dimensional sphere, where  $d$  is the input dimension. This random direction is added to the boundary point, generating a new query for the model. The results of these queries are combined to a gradient estimation  $\bar{\nabla}S$  using the Monte Carlo estimate of equation 3.3. In this equation  $u_b$  are the random directions and  $x_t$  is the boundary position.  $B$  is the number of random directions that needs to be sampled. This number increases based on the current iteration of

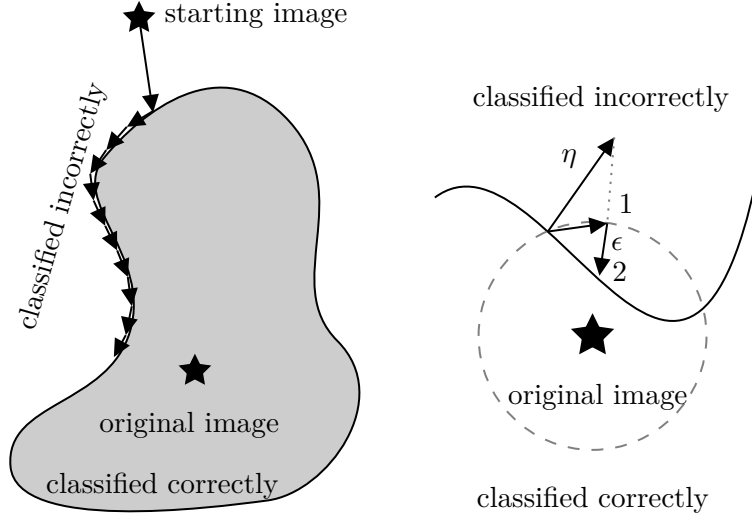


FIGURE 3.3: Intuition behind the Boundary Attack. On the left the path of the attack is shown. The first step is a projection onto the boundary, afterwards it follows the decision boundary of the class of the original image. Each arrow represents one iteration of the attack. On the right, the two different steps of each iteration can be seen. In the first step, a random direction is sampled and projected onto a sphere around the original image. The second step is to take a step towards the original image from this new position. Image inspired by [45].

the attack to reduce the variance of the estimate. The function  $\phi_{x^*}$  returns 1 if the new position is adversarial and -1 if it is not adversarial.  $\delta$  is a positive parameter determining the size of the  $d$ -dimensional sphere.

$$\widetilde{\nabla S}(x_t, \delta) := \frac{1}{B} \sum_{b=1}^B \phi_{x^*}(x_t + \delta u_b) u_b \quad (3.3)$$

Once the gradient has been estimated, the third and final operation is to take a step along this gradient. The step size is determined using a geometric progression scheme. These steps are iteratively repeated until the pre-set stopping criterion is met. Figure 3.4 represents the intuition behind **HSJA** in a graphical manner.

**HSJA** eclipses **BA** and **BBA** both on median distance against queries and attack success rates using a limited amount of queries. The untargeted version of **HSJA** is able to compete with white box attacks on the ImageNet dataset [51]. It also performs similar or superior to white box attacks such as the C&W attack [22] when evaluated against defensive mechanisms such as defensive distillation [30], region-based classification [35] and adversarial training [21].

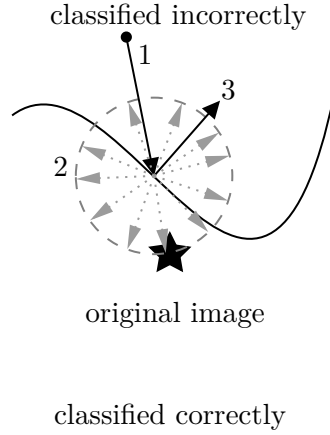


FIGURE 3.4: Intuition behind the HopSkipJumpAttack. Each iteration consists of three steps. The first step is a projection onto the boundary. The second step is the estimation of the gradient at this point. This is done by sampling directions from a uniform distribution and querying the model under attack from this new position (grey arrows). The results are combined via the Monte Carlo estimate. The third and final step is to take a step along the estimated gradient. Image inspired by [50].

### 3.3 SurFree attack

SurFree attack [52] is a decision-based attack, similar to HSJA. The authors of the paper observed that HSJA uses a significant part of the query budget to estimate the gradients of the decision boundary. These queries will not improve the distance of the best adversarial example to the original image and will therefore cause the progress to stall. In order for the attack to require less queries to be successful, SurFree works without any surrogate gradient estimation, hence the name. Instead directions are sampled randomly and a geometrical mechanism determines the biggest distortion in this direction.

#### 3.3.1 Geometrical mechanism

The geometrical mechanism starts from a point on the decision boundary:  $x_b \in \partial\mathcal{A}$ . This boundary point is located at distance  $d := \|x_b - x_o\|_2$  from the original image  $x_o$ . A unit vector  $\mathbf{u}$  is constructed starting in  $x_o$ , pointing towards  $x_b$ . The search space for a better adversarial point will be restricted to a two-dimensional random plane  $\mathcal{P}$ . The plane  $\mathcal{P}$  is spanned by vectors  $\mathbf{u}$  and  $\mathbf{v}$ , where  $\mathbf{v}$  is a random unit vector orthogonal to  $\mathbf{u}$ . The vector  $v$  will be sampled from a distribution  $\mathcal{T}$  as explained in section 3.3.2 and orthogonalized using the Gram-Schmidt process [53]. Both  $x_o$  and  $x_b$  are part of  $\mathcal{P}$ .

### 3. RELATED WORK

---

New adversarial candidates are points at a distance  $d(1 - \alpha)$  from  $x_o$  in  $\mathcal{P}$ . In polar coordinates, this results in the following equation:

$$\mathbf{z}(\alpha, \theta) = d(1 - \alpha)(\cos(\theta)\mathbf{u} + \sin(\theta)\mathbf{v}) + x_o \quad (3.4)$$

with  $\alpha \in [0, 1]$ ,  $\theta \in [-\pi, \pi]$  and with  $\theta$  being the angle with  $\mathbf{u}$ . If the new adversarial candidate yields to be adversarial, then the distance to the original image is decreased by  $100 \cdot \alpha\%$ .

Equation 3.4 has two parameters that can be chosen independently. However,  $\alpha$  can be set based on  $\theta$  in order to increase the probability of  $\mathbf{z}(\alpha, \theta)$  being adversarial. Lets assume that the intersection of the decision boundary and the random two-dimensional plane  $\partial\mathcal{A} \cap \mathcal{P}$  is a line through  $x_b$ . This line has a normal unit vector  $\mathbf{n}$  in  $\mathcal{P}$ . This vector makes an angle  $\psi$  with  $\mathbf{u}$ . The polar form of  $\mathbf{n}$  is the following:

$$\mathbf{n} = \cos(\psi)\mathbf{u} + \sin(\psi)\mathbf{v}. \quad (3.5)$$

The optimal adversarial point, minimizing the distance to  $x_o$ , would be the orthogonal projection of  $x_o$  on  $\partial\mathcal{A} \cap \mathcal{P}$ . This projection is achieved when  $\theta = \psi$  and  $\alpha = 1 - \cos(\psi)$ . However, it is not possible to create this projection in practice, since the attacker has no knowledge about the angle  $\psi$ . The calculation in appendix A shows that  $z(\alpha, \theta)$  is adversarial if the following equation holds:

$$\left| \frac{1 - (1 - \alpha)\cos(\theta)}{(1 - \alpha)\sin(\theta)} \right| \leq \tan(\psi) \operatorname{sgn}(\theta) \quad (3.6)$$

#### 3.3.2 Sampling random directions

### 3.4 Stateful defense

The defensive schemes discussed in section 2.2.2 all operate on the query level. They try to detect and flag possible attacks based on a single query without taking other context into account. The stateful detection mechanism by Chen, Carlini and Wagner [54] is different in this aspect. As the name suggests, it holds state of previously submitted queries. It is similar to the defenses that use proximity measurements, but the measurement is between queries instead of between the query and training data.

All queries submitted to the model equipped with a stateful detection mechanism are stored in a history buffer. Each user of the model has a distinct history buffer, where its queries are stored. These buffers can be bounded by time or number of queries depending on the resources available and the use case of the model. Each time a query is submitted to the model, the average distance to its  $k$  nearest neighbors is calculated and if this distance is lower than a certain threshold, then the user gets flagged by the mechanism. Appropriate actions such as banning the account can be taken.

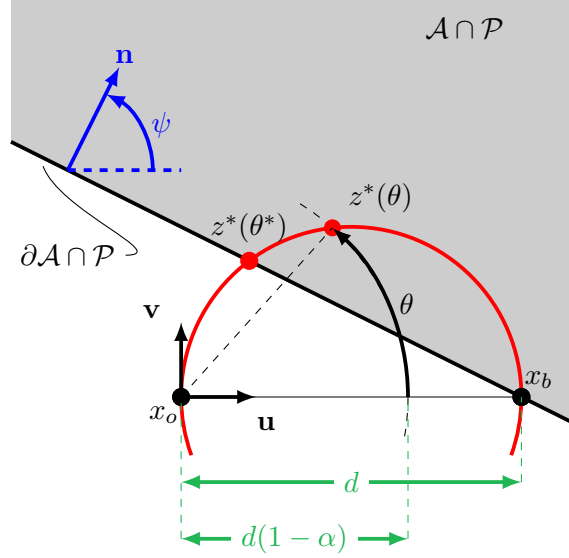


FIGURE 3.5: •

The distance metric is not calculated in input space. Each query is encoded by a deep similarity encoder [55] to an encoded space, typically of a lower dimension. In this encoded space, images which represent perceptually similar objects are clustered together. The advantage of the encoded space is twofold. Firstly, the dimension of the encoded space is smaller than the dimension of the input space. Therefore less space is needed to store the history buffers. Secondly, simpler distance metrics such as  $L_2$ -distance in input space can easily be evaded by an attacker. For example the  $L_2$ -distance can be significantly increased by simply rotating or shifting the input image.

The parameter  $k$ , the number of neighbors to consider is picked as follows. As the training data of the model consists of only benign queries, no attacks should be flagged when feeding the stateful detection mechanism with this data. To allow for some more leniency, a false positive rate of 0.1% is still acceptable. For each value of  $k$ , a different threshold will be required to maintain the selected false positive rate. Larger values have the benefit of larger thresholds causing the defense to be more resilient, since attackers' images need to be more diverse. But  $k$  is also the number of queries needed before an attack can be flagged. Therefore too large values for  $k$  are disadvantageous. Smaller values also reduce computational cost. Chen, Carlini and Wagner set the value of  $k$  to 50 for the CIFAR-10 dataset [56], since the thresholds increased sharply up to this value. Other datasets might require different values for  $k$ .

### 3.5 PSO and distributed attacks

There have been several attempts to craft adversarial examples using an EA. Previous attempts tried to reduce the number of queries needed to create a successful adversarial example by utilizing EAs [57, 58, 59, 60, 61, 62].

GenAttack [57] and the similar efficient attack by Dong et al. [58] use genetic algorithms in order to minimize the number of queries to the model. Both algorithms reduce the dimension of the search space to improve the efficiency of the attack. Once a promising perturbation is found in this lower dimensional space, it is upscaled using a bilinear transformation. By reducing the search space, the number of individuals in the genetic algorithm can be lowered, which in turn lowers the total amount of queries. GenAttack also uses annealing schemes to adaptively scale the parameters of the algorithm. This allows it to escape local optima and improve the adversarial example further.

AdversarialPSO [59] and the similar attack from [60] use PSO as optimization routine on images and audio fragments respectively. Each particle represents a possible adversarial example. Both attacks use the standard rules of PSO as specified by [36] improved with a linearly decaying inertia weight [63]. The former attack also uses a constriction factor to avoid premature convergence [44]. While the latter solves this problem by generating new particles using a genetic algorithm when premature convergence is detected. Both algorithms rely on confidence scores to assign fitness values to certain positions in the search space. PSO-BBA [61], is similar to AdversarialPSO, but only relies on distances to determine fitness values. This attack can therefore also be used in decision-based settings instead of solely in score-based settings.

The idea behind the multi-group PSO attack [62] is to use multiple PSO swarms to escape local optima. The intuition behind it is inspired by the Distributed Denial of Service (DDoS) attack [64]. The swarm is split into multiple smaller groups and each group is placed on a single node. The groups perform the standard PSO algorithm. The best position over all groups is communicated using a dedicated server. Each group submits its queries from its own node, tricking the defensive mechanism into thinking that multiple users are submitting queries. The authors state that this will ultimately result in less detections, but they have not evaluated this against a defensive scheme.

## Chapter 4

# Approach

The research concerning adversarial attacks and defenses is predominantly driven by a game of cat and mouse. Whenever a new attack is proposed, a defensive mechanism countering the novel attack is developed and vice versa. This work aims to create a new family of algorithms that can be used to perform both targeted and untargeted attacks. The goal of these algorithms is to craft adversarial examples comparable to state of the art approaches while remaining undetected by the stateful detection mechanism [54] described in section 3.4.

### 4.1 Distribution

The stateful detection mechanism [54] makes the assumption that queries can be traced back to their adversary and that there is no cooperation between different adversaries. This assumption can be problematic as  $N$  collaborating adversaries can theoretically reduce the number of submitted queries per adversary by a factor  $1/N$ . Even a single adversary could set up multiple accounts and submit queries on each account until it is banned. Due to the reduced number of submitted queries, less attacks will be detected since each buffer of the defense mechanism only holds a fraction of all queries.

This work will aim to evade the detection mechanism by distributing the query submissions over multiple nodes. Each node will represent a different user of the model under attack. The users could in theory be all different persons or they could be different accounts of the same person.

As described in section 3.5, several attempts have been made to distribute adversarial attacks [61, 62]. However, none of these attacks have been evaluated against the stateful detection mechanism, since the goal of the distribution was to make the attack more efficient. This work will distribute the query submission over multiple nodes in order to avoid detection.

## 4.2 Optimization

As previously mentioned, reducing the number of submitted queries per adversary by a factor  $1/N$ , where  $N$  is the number of collaborators, is straightforward. Adversaries can gain knowledge about the search space by cooperating with other adversaries. They can leverage this knowledge in order to reduce the number of submitted queries even more. This idea has been utilized by multiple algorithms that were mentioned in section 3.5. These algorithms used some form of **PSO** to optimize the final adversarial example. However all but the **PSO-BBA** algorithm by Xiang et al [61] rely on the confidence score of the model for the fitness value calculation. All attacks discussed use **PSO** as an attack in itself. This work will combine the benefits of state of the art black box attacks and **PSO**.

## 4.3 Approach

The remaining sections of this work will propose a new family of adversarial attacks. First a threat model is defined in section 4.4. All remaining experiments will be performed with this threat model in mind. Afterwards the novel adversarial algorithm is proposed in section 5.3. This algorithm is iteratively improved based on the results of the experiments. Finally, the final algorithm will be compared with state of the art decision-based attacks.

During the process of creating, improving and optimizing the attack, this work tries to answer the following research questions.

- What are the (dis)advantages of using **PSO** in relation to vanilla adversarial attacks?
- How can **PSO** be combined with state of the art adversarial attacks?
- What are the (dis)advantages of distributing an adversarial attack?

This work concludes with an answer on these questions in section 6.

## 4.4 Threat model

This section will describe the threat model that will be used for the remainder of this work. The first and most impactful assumption is that the model under attack will only output labels for the input. There will be no confidence scores or model parameters available. The proposed attack will have to be a decision-based attack. Most real world **APIs** will only expose the final decision to the user, causing decision-based attacks to be the only viable option. Decision-based attacks are the most restricted type of attack. Therefore they can also be applied to score-based and white-box models.



The proposed attack will be a targeted attack, as it is the most relevant type from a security point of view. A targeted attack implementation can also easily be ported to an untargeted one. This is done by running a targeted attack for every possible class and selecting the one closest to the original image.

The model will be defended by a stateful detection mechanism [54] as described in section 3.4. It will have a query bounded buffer for each account. Once an query has been flagged as potential attack, the account that submitted the query will be banned. The user will have to set up a fresh account in order to submit queries again.

There are two main fees associated with the model. The first fee is related to setting up an account. The assumption is made that setting up an account requires a valid credit card or phone number. Whenever an account is banned, the credit card or phone number is invalidated in the system, requiring the attacker to sign up for a new credit card or register for a new phone number. This cost is identical to the assumption made by Chen et al [54]. The second fee is related to the amount of queries submitted to the model. The more queries an attacker submits, the higher the total cost will be. Chen et al [54] did not incorporate a cost per query, but most vision APIs, such as Google Cloud Vision [65] and Amazon Rekognition [66] use this type of fee.

These fees ensure that attackers need to be both efficient and evasive in order to be successful. The evasiveness of the attack is measured as the number of detections by the stateful defense mechanism. Each detection essentially means that the attacker has to set up a new account and pay the associated cost. To compare the efficiency of different attacks, the following strategy will be used. Every run of an attack has a budget of queries. The attack that has created an adversarial example closest to the original image inside this budget of queries, is the most efficient. The evasiveness of an attack is the primary concern, since the cost of setting up a new account is significantly higher than the cost per query.



# Chapter 5

## Evaluation

This chapter will discuss the experiments performed, as well as the ideas and intuitions behind them. The adversarial algorithm will be refined throughout the different subsections and the results of the refinement will be discussed at the end of each subsection.

### 5.1 Evaluation protocol

This subsection described the evaluation protocol that will be followed for all experiments performed. Experiments will be performed with the MNIST [13] and CIFAR [56] datasets. Some examples of these datasets are visualised in Figures 5.1 and 5.2 respectively. A black box model is trained using the training data of the respective dataset. The architectures of the models are identical to the ones used in [22, 30]. A summary of the two architectures can be found in Table 5.1. The training parameters are also identical to the ones used in [22, 30]. These models remain unchanged for all experiments.

The trained models are then used to classify all instances in the test set of their respective dataset. The incorrectly classified examples are filtered out of this set as they are essentially already adversarial. The remaining examples can be used for experiments. A list of experiments is generated based on the remaining examples. An experiment consists of an original image, a target label and starting position(s). All future refinements will therefore perform the same set of experiments in order to make the comparison more fair. All random effects present in the algorithm are seeded for the same purpose.

The stateful defense mechanism [54] will use a query bounded detector buffer of 1000 queries per user. The value of  $k$  will be set to 50, as suggested by the authors of the paper. A threshold is determined in order to achieve a 0.1% false positive detection rate on the training data. The detector buffer will be cleared after each detection to simulate the creation of a new account.



FIGURE 5.1: Some examples of the MNIST [13] dataset.



FIGURE 5.2: Some examples of the CIFAR-10 [56] dataset.

TABLE 5.1: Model architectures for the MNIST and CIFAR models. The architectures are identical to [22, 30].

Layer type	MNIST Model	CIFAR Model
Convolution + ReLU	$3 \times 3 \times 32$	$3 \times 3 \times 64$
Convolution + ReLU	$3 \times 3 \times 32$	$3 \times 3 \times 64$
Max Pooling	$2 \times 2$	$2 \times 2$
Convolution + ReLU	$3 \times 3 \times 64$	$3 \times 3 \times 128$
Convolution + ReLU	$3 \times 3 \times 64$	$3 \times 3 \times 128$
Max Pooling	$2 \times 2$	$2 \times 2$
Fully Connected + ReLU	200	256
Fully Connected + ReLU	200	256
Softmax	10	10

The algorithm receives a budget of 25 000 queries to create an adversarial example. The final adversarial example will be evaluated on the  $L_2$ -distance to the original image and the total number of detections by the stateful defense mechanism.

## 5.2 Determining the baseline

The goal of this subsection is to determine a baseline to which all future algorithms will be compared. The baseline will be a vanilla BBA using the same hyperparameters as suggested in [48]. The orthogonal step will be set to 0.05 and the source step to 0.002. Both the Perlin noise improvement and the regional masking improvement will be used. No gradients of surrogate models will be calculated as the authors of

TABLE 5.2: Results for the **BBA** baseline approach on MNIST and CIFAR-10 dataset.

Attack	MNIST		CIFAR	
	Distance	Detections	Distance	Detections
Baseline <b>BBA</b>	2.807	413	1.306	474

the paper already mentioned that the improvement was marginal.

Table 5.2 reports the average  $L_2$ -distance and the average number of detections for the different datasets.

### 5.3 Applying PSO to BBA

As mentioned in section 4.2, previous attempts to craft adversarial example using **PSO** used **PSO** to guide adversarial examples through the search space closer to the original image. This work aims to combine the benefits of **PSO** and an existing decision-based attack such as **BBA**.

Vanilla **BBA** has the disadvantage that, depending on the starting position, it might get trapped in a local optimum. By using **PSO** in combination with **BBA**, multiple starting positions can be explored and the probability of getting trapped is lowered. The intuition behind this idea is shown in Figure 5.3. The more starting positions there are, the higher the probability of finding the global optimum. There is however a clear trade-off in terms of efficiency. By having  $n$  different starting positions, the query budget is essentially reduced by a factor  $n$  for each starting position.

The efficiency reduction does not have to pose a problem due to the implicit communication in the swarm. Particles can move to promising regions in the search space based on the information of their peers. The promising regions are therefore more queried.

The proposed **PSO-BBA** algorithm works as follows. Particles will perform a more aggressive version of **BBA**. The initial source step  $\epsilon$  is set significantly higher than in the vanilla version. The increased source step might cause the particle to end up in a non-adversarial decision region. Once this happens, the standard **PSO** equations (2.1 and 2.2) are used to guide the particle back to the adversarial region. Later refinements of this algorithms will deal with attack iterations. An attack iteration is defined as a specified budget of the aggressive **BBA** (50 queries in all experiments) or a single query submission using the **PSO** rules.

The inertia weight of the **PSO** equations is determined by the linearly decaying scheme of equation 2.3, with the weight decaying from 1 to 0. The acceleration

coefficients  $c$  are set based on an idea of multi-group **PSO**. Two groups with opposite acceleration coefficients are created. This approach helps escape local optima [67]. The equations for both  $c_p$  and  $c_g$  are:

$$c_p = \begin{cases} \max(A1, A2), & \text{if } i \bmod 2 = 0 \\ \min(A1, A2), & \text{else} \end{cases}$$

$$c_g = \begin{cases} \min(A1, A2), & \text{if } i \bmod 2 = 0 \\ \max(A1, A2), & \text{else} \end{cases}$$

Here  $i$  is the index of the particle in its swarm. The values of  $A1$  and  $A2$  are 1 and 2 respectively. These values are suggested by the authors of [62].

The source step  $\epsilon$  will be changed after every iteration of the attack. Two separate multipliers are used to respectively increase and decrease the value of this parameter. The value of  $\epsilon$  is slightly increased if the new position is still adversarial. Likewise it is decreased if the position is no longer adversarial.

By using **PSO** in combination with **BBA**, the advantage of multiple starting points, as explained in Figure 5.3, can be exploited, without having a less efficient attack as a whole. Whenever particles end up in non-adversarial regions, they will move closer to the best known position in the swarm due to the **PSO** equations. At the end of the attack, most particles will be in the same area of the search space, allowing for more exploitation in this specific area.

The **PSO** framework requires a fitness function to quantify the fitness of a position for the problem at hand. The authors of AdversarialPSO [59] suggested the following fitness function  $f$ :

$$f(x) = |p_x - p_{x'}| - \frac{c}{n} \|x - x'\|_2 \quad (5.1)$$

Where  $x$  is the position of the particle,  $x'$  is the original image,  $p_x$  and  $p_{x'}$  are the confidence scores of the model in predicting the label of  $x$  and  $x'$  respectively and  $c$  is a constant to weight the penalty. However, as discussed in section 4.4, the confidence scores are not available for this specific attack. The fitness function in equation 5.1 also assumes that the position  $x$  will always be adversarial. This will not be the case in the proposed algorithm. The fitness function will therefore be altered to the following:

$$f(x) = \begin{cases} \|x - x'\|_2, & \text{if } x \text{ is adversarial} \\ +\infty, & \text{else} \end{cases} \quad (5.2)$$

TABLE 5.3: Results for the **PSO-BBA** approach on MNIST and CIFAR-10 dataset.

Attack	MNIST		CIFAR	
	Distance	Detections	Distance	Detections
Baseline <b>BBA</b>	2.807	413	1.306	474
<b>PSO-BBA</b> (5 particles)	2.712	138	1.239	257
<b>PSO-BBA</b> (10 particles)	3.157	44	2.290	184

The infinite value for the fitness function inside non-adversarial decision regions acts as a penalty, causing the particles to quickly diverge from these regions.

The same set of experiments as in section 5.2 has been performed. The experiments have been done using attacks with both five and ten particles. The initial step sizes have been set to 0.25 and 0.20 for MNIST and CIFAR respectively. The values for the increasing and decreasing multiplier have been set to 1.05 and 0.99 for both datasets. The results of the experiments can be found in Table 5.3.

The five particle **PSO-BBA** algorithm outperforms the baseline in terms of distance to the original image on both datasets. This is not the case for the ten particle version. The high number of particles requires sufficient queries in the beginning of the attack in order to discover promising regions in the high dimensional search space of CIFAR. This exploration requires more queries than the query budget allows. The number of detections is lower for all variants of **PSO-BBA**. The different starting points have the added advantage that initial queries are more spread out over the search space. These queries are therefore less similar and the detector will not flag as much attacks. This effect can be seen in Figure 5.4 for the MNIST experiments.

The number of detections drops as the number of particles increases. This was to be expected from the intuition behind **PSO-BBA**. However, the average  $L_2$ -distance to the original image is higher for ten particles compared to five. To reduce the number of different parameters in future refinements, only swarms with five particles will be considered from here on.

## 5.4 Towards distribution of the attack

The stateful defense mechanism [54] makes the assumption that there is no collaboration between adversaries. Based on this assumption, each account or user will have its own detector buffer. This subsection aims to exploit this assumption by distributing the query submission over multiple accounts.

By distributing the query submissions over multiple accounts, the number of total detection should be reduced compared to **BBA** due to two reasons. The first reason

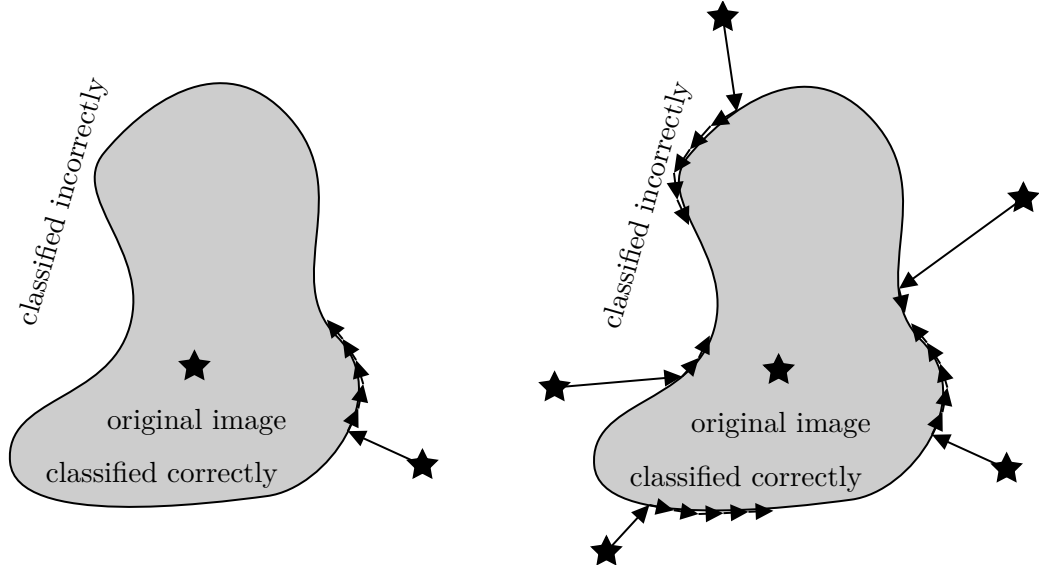


FIGURE 5.3: The vanilla version of **BBA** might get stuck in a local optimum depending on the starting point (left plot). By starting from multiple positions, the probability that **BBA** gets stuck in a local optimum is reduced (right plot). The multiple starting points are particles in a **PSO** swarm. Image inspired by [45].

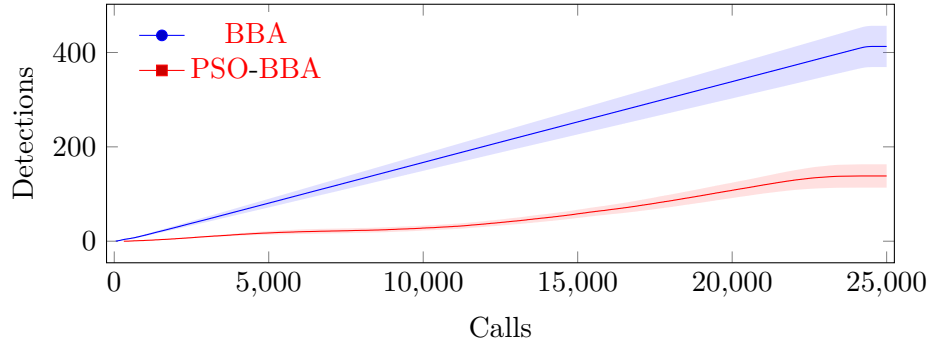


FIGURE 5.4: The cumulative number of detections for different attack algorithms for the MNIST experiments. The number of detections of the **BBA** algorithm steadily increases with the number of calls. The number of detections of the **PSO-BBA** algorithm increases more slowly in the beginning due to the dispersed nature of the attack. The increase is more sharp near the end of the attack when the particles converge. The plot shows the 95% confidence regions around the graphs.



is obvious. If a budget of 25 000 queries is distributed over  $N$  accounts, then every account will submit  $25\,000 / N$  queries. The less queries there are submitted, the less potential attacks will be detected. The second reason is due to **PSO**. As stated before, the different particles reside in different parts of the search space. This means that a detector buffer of a specific account will contain queries from all over the search space, causing the inter query distances to be larger. The latter reason was also present in the algorithm described in section 5.3.

It should be noted that there is only distribution at the query submission level. The algorithm itself will be executed on one machine without the need for distribution. Even the query submission can be done from one machine by constantly changing the **API** key or by logging in and out of a specific account whenever it is needed. This approach makes the assumption that the **API** under attack does not track the IP address of the machine that submitted the query. Multiple machines might therefore be more convenient in a real attack setting.

This work simulates the multiple machines setting by having separate node (or machine) objects on the same machine. Instead of having a detector buffer for each node at the location of the model, the responsibility is moved upstream to the nodes themselves. Each node will have a buffer to which the query will be added. All queries will be passed to the model afterwards. In Figure 5.5, the distribution is shown schematically.

The queries will be distributed over the nodes based on a distribution scheme. Three different distribution schemes will be used in the experiments. The first two schemes are inspired by the **Round-Robin (RR)** scheduling algorithm [68] and will be called the **RR** and **Modified Round-Robin (MRR)** distribution schemes. The third distribution scheme will be built on the assumption that the inter query distance in the detector buffers needs to be maximized.

The **RR** distribution scheme is fairly straightforward. The outgoing queries of the algorithm will be evenly distributed over the nodes in circular order. This scheme does not have a notion of the underlying attack. This can cause successive queries forwarded to the same node to be relatively close together by chance.

The **MRR** distribution scheme maps a particle to a node for the entire duration of an attack iteration. The queries submitted by a single particle during one attack iteration are relatively close together, but this approach ensures that all particles will submit queries to all nodes. This might ultimately cause less detections. The **MRR** distribution scheme maps a number of particles on a number of nodes. If the amount of particles is equal to the number of nodes then the mapping is straightforward. If the number of nodes is less than the number of particles then dummy nodes are introduced in the mapping. The dummy nodes forward their received queries to a random existing node. The random node changes every rotation. If the number of particles is less than the number of nodes then dummy particles that do not submit

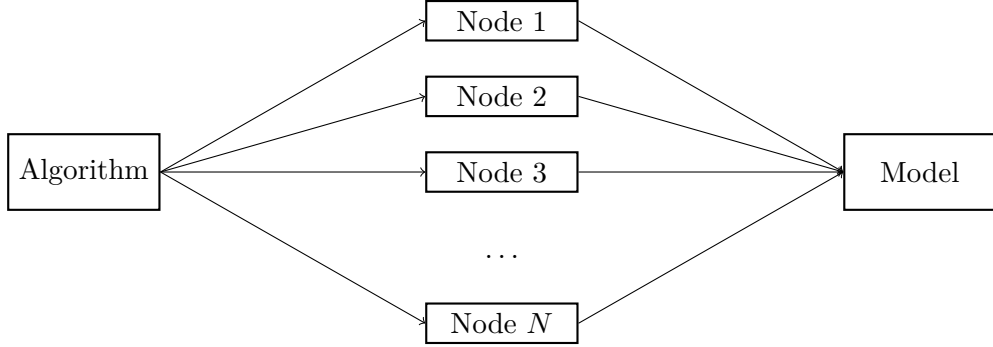


FIGURE 5.5: Schematic overview of the query submission distribution. Each arrow represents a query. The algorithm distributes the queries over the different nodes. Every node forwards its received queries to the model using the corresponding account. In a real setting, the model will have a detector buffer for every account.

In this work the buffers are located on the nodes themselves.

queries, are introduced. After every attack iteration, the mapping rotates. The mapping process is shown in Figure 5.6.

The ultimate goal of the attack is to evade the stateful detection mechanism that defends the model under attack. Potential attacks are flagged based on the inter query distance in a detector buffer. The distance-based distribution scheme tries to exploit the defense. Instead of distributing queries over nodes in a **RR**-fashion, queries are submitted to nodes based on the distance to the previously submitted queries to this node. The node with the highest average distance to the previously submitted queries is selected. For this purpose, every node will hold an internal buffer of queries. The size of this buffer can be chosen by the attacker, as can the distance metric. Standard  $L_2$ -distance can be chosen, but this has the same pitfalls as discussed in section 3.4. It is also possible to train another similarity encoder and use the  $L_2$ -distance in the embedded space in order to resemble the defense more closely. The distribution schemes will be called **Distance-Based (DB)** and **Embedded-Distance-Based (EDB)** distribution schemes respectively.

The experiments performed in this section will use a buffer size at each node of 20 for both the **DB** and **EDB** distribution schemes. The similarity encoder in the case of the **EDB** scheme will be trained on the test data of the respective dataset and has an output dimension of 128. The results of the experiments can be found in Table 5.4.

## 5.5 Throwing the defense off the scent

## 5.6 Optimizing the attack

TABLE 5.4: Results for the distributed **PSO-BBA** approaches on MNIST and CIFAR-10 dataset.

Attack	MNIST		CIFAR	
	Distance	Detections	Distance	Detections
Baseline <b>BBA</b>	2.807	413	1.306	474
<b>PSO-BBA</b>	2.712	138	1.239	257
<b>RR-PSO-BBA</b> (5 nodes)	2.712	124	1.239	224
<b>RR-PSO-BBA</b> (10 nodes)	2.712	104	1.239	202
<b>MRR-PSO-BBA</b> (5 nodes)	2.712	110	1.239	227
<b>MRR-PSO-BBA</b> (10 nodes)	2.712	93	1.239	206
<b>DB-PSO-BBA</b> (5 nodes)	2.712	107	1.239	230
<b>DB-PSO-BBA</b> (10 nodes)	2.712	87	1.239	207
<b>EDB-PSO-BBA</b> (5 nodes)	2.712	108	1.239	229
<b>EDB-PSO-BBA</b> (10 nodes)	2.712	88	1.239	205

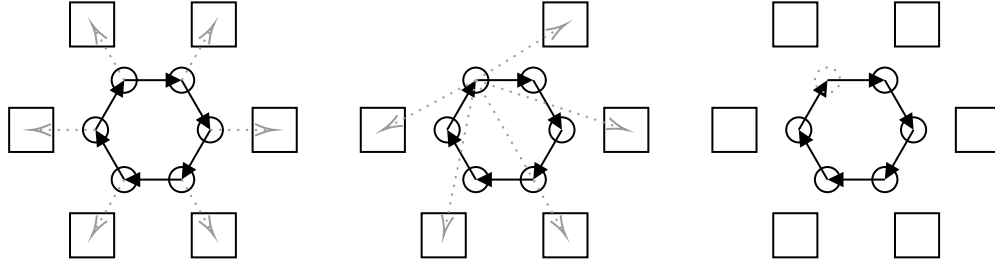


FIGURE 5.6: Overview of the round-robin based distribution scheme. The particles are represented by the small circles and the squares represent the different nodes. If the number of particles is equal to the number of nodes (left), then particles submit their queries to the corresponding node. If the number of nodes is less than the number of particles (middle), then the particles without nodes will submit their queries to a random node. This random node will change every attack iteration. If the number of particles is less than the number of nodes (right), then a dummy particles are introduced. These dummy particles will not submit any queries, therefore some nodes will not receive queries. After every attack iteration the particle mapping will rotate.



## Chapter 6

## Conclusion



# Appendices





# Appendix A

## SurFree calculations

The following calculations will show how equation 3.6 was derived. Assume that the normal vector  $\mathbf{n}$  is pointing outside, i.e. to the adversarial region  $\mathcal{A}$ . A candidate point  $\mathbf{z} \in \mathcal{P}$  is adversarial if  $(\mathbf{z} - x_b)^\top \mathbf{n} \geq 0$ . Under the assumption of a (local) flat decision boundary, this means that both vector  $n$  and vector  $(\mathbf{z} - x_b)$  are pointing outside. This assumption is justified according to [47]. Below follows the derivation of equation 3.6:

$$\begin{aligned}
& (\mathbf{z} - x_b)^\top \mathbf{n} \geq 0 \\
& \Leftrightarrow [d(1 - \alpha)(\cos(\theta)\mathbf{u} + \sin(\theta)\mathbf{v}) + x_o - x_b]^\top \mathbf{n} \geq 0 \quad \text{Equation 3.4} \\
& \Leftrightarrow [d(1 - \alpha)(\cos(\theta)\mathbf{u} + \sin(\theta)\mathbf{v}) - d\mathbf{u}]^\top \mathbf{n} \geq 0 \quad \mathbf{u} = (x_b - x_o)/d \\
& \Leftrightarrow [((1 - \alpha)\cos(\theta) - 1)\mathbf{u} + ((1 - \alpha)\sin(\theta))\mathbf{v}]^\top \mathbf{n} \geq 0 \\
& \Leftrightarrow ((1 - \alpha)\cos(\theta) - 1)\cos(\psi) + ((1 - \alpha)\sin(\theta))\sin(\psi) \geq 0 \quad \text{Equation 3.5} \\
& \Leftrightarrow ((1 - \alpha)\cos(\theta) - 1)\cos(\psi) \geq -((1 - \alpha)\sin(\theta))\sin(\psi) \\
& \Leftrightarrow (1 - (1 - \alpha)\cos(\theta))\cos(\psi) \leq ((1 - \alpha)\sin(\theta))\sin(\psi) \quad x \geq y \Leftrightarrow -x \leq -y \\
& \Leftrightarrow \left| \frac{1 - (1 - \alpha)\cos(\theta)}{(1 - \alpha)\sin(\theta)} \right| \leq \tan(\psi) \operatorname{sgn}(\theta) \quad \tan(x) = \frac{\sin(x)}{\cos(x)}
\end{aligned}$$

The absolute value and sign operator are present due to the fact that the inequality should reverse in case of a division by a negative number.



# Bibliography

- [1] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. 65(6):386–408.
- [2] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on International Conference on Machine Learning, ICML'10*, page 807–814, Madison, WI, USA, 2010. Omnipress.
- [3] Wikipedia. Heaviside step function. [https://en.wikipedia.org/wiki/Heaviside\\_step\\_function](https://en.wikipedia.org/wiki/Heaviside_step_function), 04 2022. Accessed: 2022-04-30.
- [4] Wikipedia. Softmax function. [https://en.wikipedia.org/wiki/Softmax\\_function](https://en.wikipedia.org/wiki/Softmax_function), 03 2022. Accessed: 2022-04-30.
- [5] M.V. Valueva, N.N. Nagornov, P.A. Lyakhov, G.V. Valuev, and N.I. Chervyakov. Application of the residue number system to reduce hardware costs of the convolutional neural network implementation. *Mathematics and Computers in Simulation*, 177:232–243, 2020.
- [6] Steve Lawrence, C. Lee Giles, Ah Chung Tsoi, and Andrew D. Back. Face recognition: A convolutional neural network approach. *IEEE TRANSACTIONS ON NEURAL NETWORKS*, 8(1):98–113, 1997.
- [7] Xiangang Li and Xihong Wu. Constructing long short-term memory based deep recurrent neural networks for large vocabulary speech recognition. *CoRR*, abs/1410.4281, 2014.
- [8] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. *CoRR*, abs/1409.3215, 2014.
- [9] Hansika Hewamalage, Christoph Bergmeir, and Kasun Bandara. Recurrent neural networks for time series forecasting: Current status and future directions. *CoRR*, abs/1909.00590, 2019.
- [10] Jon Russell. Google’s alphago ai wins three-match series against the world’s best go player. <https://techcrunch.com/2017/05/24/alphago-beats-planets-best-human-go-player-ke-jie/amp/?guccounter=1>, 05 2017. Accessed: 2021-12-08.

- [11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- [12] Wikipedia. Perceptron. <https://en.wikipedia.org/wiki/Perceptron>, 2022 03. Accessed: 2022-04-16.
- [13] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [14] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks, 2014.
- [15] Fatemeh Vakhshiteh, Ahmad Nickabadi, and Raghavendra Ramachandra. Adversarial attacks against face recognition: A comprehensive study, 2021.
- [16] Abhiram Gnanasambandam, Alex M. Sherman, and Stanley H. Chan. Optical adversarial attack, 2021.
- [17] Octavian Suciu, Scott E. Coull, and Jeffrey Johns. Exploring adversarial examples in malware detection, 2019.
- [18] Nicholas Carlini and David Wagner. Audio adversarial examples: Targeted attacks on speech-to-text, 2018.
- [19] Giulio Zizzo, Chris Hankin, Sergio Maffei, and Kevin Jones. Invited: Adversarial machine learning beyond the image domain. In *2019 56th ACM/IEEE Design Automation Conference (DAC)*, pages 1–4, 2019.
- [20] Victoria Drake. Threat modeling. [https://owasp.org/www-community/Threat\\_Modeling#:~:text=A%20threat%20model%20is%20a,through%20the%20lens%20of%20security](https://owasp.org/www-community/Threat_Modeling#:~:text=A%20threat%20model%20is%20a,through%20the%20lens%20of%20security), 09 2021. Accessed: 2022-03-31.
- [21] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2015.
- [22] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks, 2017.
- [23] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defenses. In *International Conference on Learning Representations*, 2018.
- [24] Wikipedia. Distance. <https://en.wikipedia.org/wiki/Distance>, 11 2021. Accessed: 2021-12-09.
- [25] Wikipedia. Discrete space. [https://en.wikipedia.org/wiki/Discrete\\_space](https://en.wikipedia.org/wiki/Discrete_space), 04 2022. Accessed: 2022-04-30.

- 
- [26] Gabriel Resende Machado, Eugênio Silva, and Ronaldo Ribeiro Goldschmidt. Adversarial machine learning in image classification: A survey toward the defender’s perspective. *ACM Comput. Surv.*, 55(1), nov 2021.
  - [27] Gabriel Resende Machado, Eugênio Silva, and Ronaldo Ribeiro Goldschmidt. Adversarial machine learning in image classification: A survey towards the defender’s perspective. *CoRR*, abs/2009.03728, 2020.
  - [28] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, ASIA CCS ’17, page 506–519, New York, NY, USA, 2017. Association for Computing Machinery.
  - [29] Anish Athalye, Nicholas Carlini, and David A. Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *CoRR*, abs/1802.00420, 2018.
  - [30] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. *CoRR*, abs/1511.04508, 2015.
  - [31] Shixiang Gu and Luca Rigazio. Towards deep neural network architectures robust to adversarial examples, 2015.
  - [32] Dan Hendrycks and Kevin Gimpel. Visible progress on adversarial images and a new saliency map. *CoRR*, abs/1608.00530, 2016.
  - [33] Chuan Guo, Mayank Rana, Moustapha Cissé, and Laurens van der Maaten. Countering adversarial images using input transformations. *CoRR*, abs/1711.00117, 2017.
  - [34] Nicolas Papernot and Patrick D. McDaniel. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. *CoRR*, abs/1803.04765, 2018.
  - [35] Xiaoyu Cao and Neil Zhenqiang Gong. Mitigating evasion attacks to deep neural networks via region-based classification. *CoRR*, abs/1709.05583, 2017.
  - [36] James Kennedy and Russell Eberhart. Particle swarm optimization. In *Proceedings of ICNN’95-international conference on neural networks*, volume 4, pages 1942–1948. IEEE, 1995.
  - [37] Marco Dorigo, Vittorio Maniezzo, and Alberto Coloni. Ant system: Optimization by a colony of cooperating agents. *IEEE transactions on systems, man, and cybernetics - part b. IEEE transactions on systems, man, and cybernetics. Part B, Cybernetics : a publication of the IEEE Systems, Man, and Cybernetics Society*, 26:29–41, 02 1996.

- [38] John H. Holland. *Genetic Algorithms and Adaptation*, pages 317–333. Springer US, Boston, MA, 1984.
- [39] Rami Abousleiman and Osamah Rawashdeh. Electric vehicle modelling and energy-efficient routing using particle swarm optimisation. *IET Intelligent Transport Systems*, 10(2):65–72, 2016.
- [40] Tadashi Yamada and Zukhruf Febri. Freight transport network design using particle swarm optimisation in supply chain–transport supernetwork equilibrium. *Transportation Research Part E: Logistics and Transportation Review*, 75:164–187, 2015.
- [41] Bruno Almeida and Victor Coppo Leite. *Particle Swarm Optimization: A Powerful Technique for Solving Engineering Problems*. 12 2019.
- [42] Romeela Mohee and Ackmez Mudhoo. Analysis of the physical properties of an in-vessel composting matrix. *Powder Technology*, 155(1):92–99, 2005.
- [43] Y. Shi and R.C. Eberhart. Empirical study of particle swarm optimization. In *Proceedings of the 1999 Congress on Evolutionary Computation-CEC99 (Cat. No. 99TH8406)*, volume 3, pages 1945–1950 Vol. 3, 1999.
- [44] M. Clerc and J. Kennedy. The particle swarm - explosion, stability, and convergence in a multidimensional complex space. *IEEE Transactions on Evolutionary Computation*, 6(1):58–73, 2002.
- [45] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models. *arXiv:1712.04248 [cs, stat]*, February 2018. arXiv: 1712.04248.
- [46] Ye Wenhe. Trust-region methods. [https://optimization.mccormick.northwestern.edu/index.php/Trust-region\\_methods](https://optimization.mccormick.northwestern.edu/index.php/Trust-region_methods), 06 2014. Accessed: 2021-12-09.
- [47] Alhussein Fawzi, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. Robustness of classifiers: from adversarial to random noise. *CoRR*, abs/1608.08967, 2016.
- [48] Thomas Brunner, Frederik Diehl, Michael Truong Le, and Alois Knoll. Guessing Smart: Biased Sampling for Efficient Black-Box Adversarial Attacks. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4957–4965, October 2019. arXiv: 1812.09803.
- [49] Ken Perlin. An image synthesizer. In *Proceedings of the 12th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH ’85, page 287–296, New York, NY, USA, 1985. Association for Computing Machinery.
- [50] Jianbo Chen, Michael I. Jordan, and Martin J. Wainwright. HopSkipJumpAttack: A Query-Efficient Decision-Based Attack. *arXiv:1904.02144 [cs, math, stat]*, April 2020. arXiv: 1904.02144.

- 
- [51] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
  - [52] Thibault Maho, Teddy Furon, and Erwan Le Merrer. Surf-free: a fast surrogate-free black-box attack. *CoRR*, abs/2011.12807, 2020.
  - [53] Wikipedia. Gram Schmidt process. [https://en.wikipedia.org/wiki/Gram-Schmidt\\_process](https://en.wikipedia.org/wiki/Gram-Schmidt_process), 04 2022. Accessed: 2022-04-30.
  - [54] Steven Chen, Nicholas Carlini, and David Wagner. Stateful Detection of Black-Box Adversarial Attacks. *arXiv:1907.05587 [cs]*, July 2019. arXiv: 1907.05587.
  - [55] Sean Bell and Kavita Bala. Learning visual similarity for product design with convolutional neural networks. *ACM Trans. Graph.*, 34(4), jul 2015.
  - [56] Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009.
  - [57] Moustafa Alzantot, Yash Sharma, Supriyo Chakraborty, Huan Zhang, Cho-Jui Hsieh, and Mani Srivastava. GenAttack: Practical Black-box Attacks with Gradient-Free Optimization. *arXiv:1805.11090 [cs]*, June 2019. arXiv: 1805.11090.
  - [58] Yinpeng Dong, Hang Su, Baoyuan Wu, Zhifeng Li, Wei Liu, Tong Zhang, and Jun Zhu. Efficient decision-based black-box adversarial attacks on face recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7714–7722, 2019.
  - [59] Rayan Mosli, Matthew Wright, Bo Yuan, and Yin Pan. They might not be giants: crafting black-box adversarial examples with fewer queries using particle swarm optimization. *arXiv preprint arXiv:1909.07490*, 2019.
  - [60] Hyunjun Mun, Sungwan Seo, Baehoon Son, and Joobeom Yun. Black-box audio adversarial attack using particle swarm optimization. *IEEE Access*, 10:23532–23544, 2022.
  - [61] Fengtao Xiang, Jiahui Xu, Wanpeng Zhang, and Weidong Wang. A distributed biased boundary attack method in black-box attack. *Applied Sciences*, 11(21), 2021.
  - [62] Naufal Suryanto, Hyoeun Kang, Yongsu Kim, Youngyeo Yun, Harashta Tatimma Larasati, and Howon Kim. A distributed black-box adversarial attack based on multi-group particle swarm optimization. *Sensors*, 20(24), 2020.
  - [63] Y. Shi and R.C. Eberhart. Empirical study of particle swarm optimization. In *Proceedings of the 1999 Congress on Evolutionary Computation-CEC99 (Cat. No. 99TH8406)*, volume 3, pages 1945–1950 Vol. 3, 1999.

- [64] Kaspersky. What is a ddos attack? - ddos meaning. [www.usa.kaspersky.com/resource-center/threats/ddos-attacks](http://www.usa.kaspersky.com/resource-center/threats/ddos-attacks), 01 2021. Accessed: 2022-03-18.
- [65] Google. Cloud Vision pricing. <https://cloud.google.com/vision/pricing#prices>, 04 2022. Accessed: 2022-04-08.
- [66] Amazon. Amazon Rekognition pricing. <https://aws.amazon.com/rekognition/pricing/>, 04 2022. Accessed: 2022-04-08.
- [67] Naufal Suryanto, Chihiro Ikuta, and Dadet Pramadihanto. Multi-group particle swarm optimization with random redistribution. In *2017 International Electronics Symposium on Knowledge Creation and Intelligent Computing (IES-KCIC)*, pages 1–5, 2017.
- [68] Wikipedia. Round-robin scheduling. [https://en.wikipedia.org/wiki/Round-robin\\_scheduling](https://en.wikipedia.org/wiki/Round-robin_scheduling), 04 2022. Accessed: 2022-04-18.