

International Journal of Intelligent Engineering Informatics

ISSN online: 1758-8723 - ISSN print: 1758-8715

<https://www.inderscience.com/ijiei>

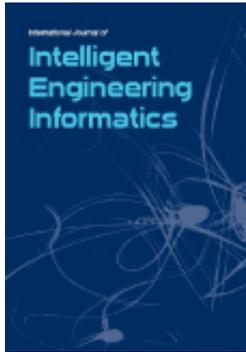
Similarity-based optimised and adaptive adversarial attack on image classification using neural network

Balika J. Chelliah, Mohammad Mustafa Malik, Ashwin Kumar, Nitin Singh, R. Regin

DOI: [10.1504/IJIEI.2023.10055562](https://doi.org/10.1504/IJIEI.2023.10055562)

Article History:

Received:	28 June 2022
Last revised:	23 January 2023
Accepted:	30 January 2023
Published online:	03 May 2023



International Journal of Intelligent Engineering Informatics

ISSN online: 1758-8723 - ISSN print: 1758-8715

<https://www.inderscience.com/ijiei>

Similarity-based optimised and adaptive adversarial attack on image classification using neural network

Balika J. Chelliah, Mohammad Mustafa Malik, Ashwin Kumar, Nitin Singh, R. Regin

DOI: [10.1504/IJIEI.2023.10055562](https://doi.org/10.1504/IJIEI.2023.10055562)

Article History:

Received:	28 June 2022
Last revised:	23 January 2023
Accepted:	30 January 2023
Published online:	03 May 2023

Similarity-based optimised and adaptive adversarial attack on image classification using neural network

Balika J. Chelliah*,
Mohammad Mustafa Malik, Ashwin Kumar,
Nitin Singh and R. Regin

Department of Computer Science and Engineering,
SRM Institute of Science and Technology,
Chennai, Tamil Nadu-600 089, India

Email: ballikaj@srmist.edu.in

Email: mm4947@srmist.edu.in

Email: ak1054@srmist.edu.in

Email: ns5861@srmist.edu.in

Email: regin12006@yahoo.co.in

*Corresponding author

Abstract: Image classification, natural language processing (NLP), and speech recognition have embraced deep learning (DL) techniques. Unrealistic adversarial samples dominate model security research. True hostile attacks are worryingly understudied. These attacks compromise real-world applications. This technique helps comprehend adversarial resistance in real-world challenges. We use real-world cases and data to test whether unreal hostile samples can protect models from genuine samples. Nodal dropouts from the first convolutional layer reveal weak and steady deep-learning neurons. Adversarial targeting links neurons to network adversaries. Neural network adversarial resilience is popular. Its DL network fails to skilfully manipulate input photographs. Our results show that unrealistic examples are as successful as realistic ones or give small enhancements. Second, we investigate the hidden representation of adversarial instances with realistic and unrealistic attacks to explain these results. We showed examples of unrealistic samples used for similar purposes and helped future studies bridge realistic and unrealistic adversarial approaches, and we released the code, datasets, models, and findings.

Keywords: deep neural network; DNN; interactive gradient shielding; generative adversarial networks; adversarial samples.

Reference to this paper should be made as follows: Chelliah, B.J., Malik, M.M., Kumar, A., Singh, N. and Regin, R. (2023) 'Similarity-based optimised and adaptive adversarial attack on image classification using neural network', *Int. J. Intelligent Engineering Informatics*, Vol. 11, No. 1, pp.71–95.

Biographical notes: Balika J. Chelliah is an Associate Professor in Computer Science and Engineering at SRM Institute of Science and Technology, Ramapuram, Chennai, India. He received his Master's and PhD in Computer Science and Engineering from SRM Institute of Technology. He has authored more than 50 papers in journals and conferences.

Mohammad Mustafa Malik works as a Business Technology Solutions Associate at ZS Associates. He earned his BTech in Computer Science and Engineering from SRM Institute of Science and Technology in 2022. He is passionate about Kubernetes, cloud computing, machine learning (ML), and DEVOPS and brings innovative thinking to his work. He was actively involved in student communities and volunteer initiatives during his studies. His strong work ethic and involvement outside of academics demonstrate his well-roundedness and drive.

Ashwin Kumar was a student at the SRM Institute of Science and Technology, Chennai. He has completed his Bachelor's degree in Computer Science and Engineering. He is currently working in an MNC as a System Engineer. He is interested in leading-edge technologies such as ML, neural networks.

Nitin Singh has completed his Engineering in Computer Science with a specialisation in Big Data Analytics from SRM University Ramapuram campus, Chennai. He gradually became interested in Artificial intelligence while working as a computer vision engineer in the LandT-NxT department. He has also done college projects like 'sentimental analysis of Twitter data' and 'AI chatbot for healthcare'. He is currently employed with Bank of America as a Software Engineer. He thanks his professors and mentors in college for lighting and guiding his path for such a beautiful and bright journey.

R. Regin is an Assistant Professor at the Department of Computer Science and Engineering at SRM Institute of Science and Technology in Ramapuram, Chennai, India. He has been awarded PhD in Computer Science and Engineering at Anna University. He has been a Software Developer and researcher for ten years and has published 40 international journals, 20 international conferences, 15 national conferences, and 10 Springer Book Chapters. He is on Ilmiah Teunuleh's advisory board and editor of the *International Journal of Technology Information and Computer*, Growing Scholar USA. He edits 'ICT-based Framework for Data Science and ML' He explores VANET, cloud computing, and infosec. He reviews for Springer and others.

1 Introduction

Deep self-learning systems have made significant and quick advances in solving several issues (Yao et al., 2019) involving complicated data processing. Major advancements in voice recognition and natural language processing (NLP), financial market research, and many more have shown high accuracy and efficiency (Shukla et al., 2019). As a result, deep neural networks (DNNs) are increasingly being used in areas including face payment, face unlocking, virtual assistants, fraud detection, and self-driven cars. Recent years have seen a rapid expansion of the deep learning (DL) field's applicability across various conventional applications. It is also crucial to note that it is a better approach than machine learning (ML) in many fields, including cybersecurity, NLP, bioinformatics, robotics, and control (Yao et al., 2019). Despite the wide range of applications given by DNNs, they are subject to unrealistic and adversarial attacks, as evidenced by many studies. Introducing a small distortion to the input sample (Feutrill et al., 2018), the model misclassifies the Adversarial Examples by complete inaccuracy, yet the naked eye cannot identify the difference.

Deep self-learning systems adversary assaults are widely recognised as one of the world's most pressing dangers to artificial intelligence security ML. Using minor adjustments to certain original instances, these assaults can construct adversarial examples specifically (Koswara and Asnar, 2019) designed to trick the model's decision-making processes. While much of the research on adversarial assaults has been done in computer vision, the techniques have also been used in various other fields, such as consumer credit and cybercrime, malware systems, and computational linguistics. For ML models to be deployed safely in the real world, it is necessary to conduct adequate assessments of adversarial assaults on them and develop strategies for making models more resistant to such attacks (Bibi et al., 2022).

In addition to assessing a system's resiliency, it is common to determine the system's correctness for hostile instances created by an assault from a set of initial cases. Additionally, adversarial hardening is a well-established approach to hardening ML models, which is defined as the application of training strategies that lead models to learn to generate correct predictions in adversarial scenarios. Even though different models' architectures and training data are unique, similar models may be attacked with the same group of adversarial examples (Jabeen and Ping, 2019). They are seen to be very susceptible to DNNs (Wu et al., 2018). As a result, scientists have created copyright adversary threats that either modify solid objects through a sequence of issue transformations or generate function instabilities that satisfy domain restrictions (also referred to as feature perturbations) to circumvent these limitations (i.e., constrained feature space attacks). In contrast to standard attacks, the instances generated by these assaults are designed to seem genuine, albeit at the risk of a larger computing price, as contrasted to the usual methods of assault (Yuan et al., 2019). On the other hand, the extra samples may be so costly in certain cases that it limits the number of samples that algorithms specialists may use to analyse and enhance resilience.

1.1 Research gaps and drivers

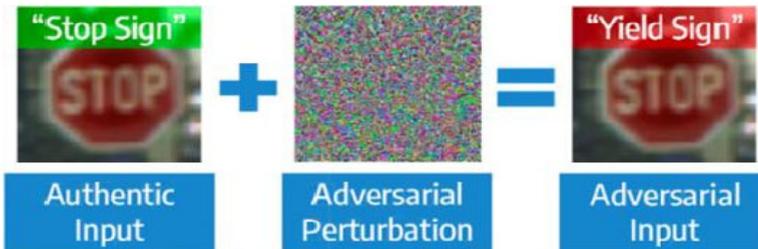
We investigated the research gap and found whether adversarial hardening on non-realistic conditions can increase model robustness when confronted with genuine scenarios as part of our effort to resolve the trade-off between realism and processing cost. Successful model hardening would be possible without creating specialised attacks that are efficient in the particular field under consideration and at an acceptable computational cost. The attackers may create these hostile examples and use them in real life, which will have many negative repercussions (Jabeen and Ping, 2019). For instance, the attacker may apply minor alterations to various traffic signs to confuse self-driving cars. Alternatively, the attacker could stage actual attacks or disturbances to confuse vehicles operating without a driver. Attackers can use different speech recognition technologies to issue menacing orders.

In order to get around these restrictions, scientists have developed copyright adversary threats that either transform solid objects through some kind of sequence of issue transformations or produce function instabilities that satisfy domain restrictions (also known as feature perturbations) (i.e., constrained feature space attacks). Most of the study on adversarial instances and their risks concentrates on fanciful hostile samples. Various methods for creating adversarial samples can successfully attack many DNNs; however, the generated images still lack a theoretical foundation.

In Huster et al. (2018), the current system recommended iterative gradient shielding (IGS), which enables gradient-based attack and region selection. The introduction of adjustable gradient shielding is also made. We can determine the vulnerability of DNNS with the help of the automatic disregard of insensitive gradient information during each iteration of an assaulting procedure. However, this method cannot determine the influence of adversarial samples in actual situations. The method requires a high level of specialised knowledge and is challenging to implement in real-time. An important concern concerning the absence of research into realistic hostile samples and their effects on the security of actual systems has been brought up by a recent study. Testing ML models' resilience to adversarial attacks and formulating plans to make them more resilient are required before they can be used in real-world circumstances.

Traditional adversarial attacks can not be utilised to measure robustness since they provide examples that aren't realistic (i.e., do not map to real-world entities). As a result, domain-specific adversarial attacks have been studied in which real objects are modified via a series of entity-space transformations (Figure 1).

Figure 1 Effect of distortions on input (see online version for colours)



1.2 Main contribution

In order to assess the stability of the networks of deep self-learning systems in real-world application settings, we will create domain-specific adversarial assaults that will impact physical objects in this project. Feature space attacks, which are meant to look realistic, can be used for this. In order to exploit a text classification model, the attacker could, for instance, replace terms with synonyms. Although this tactic has a hefty processing cost, it works. One way to get around this restriction is to utilise adversarial hardening on fictitious cases (Xie et al., 2020), which involves training techniques that let models make accurate predictions on hostile examples. For ML models to be securely used in the real world, it is necessary to create strategies for increasing model resistance and analyse adversarial attacks on them appropriately. A system's robustness can be assessed by evaluating the prediction performance on hostile instances produced by an assault from a collection of source cases. Similar to this, adversarial hardening, a tried-and-true methodology for strengthening ML models, entails training methods that instruct models to make accurate predictions in adversarial scenarios (Pahadiya et al., 2021).

According to recent discoveries (Lin et al., 2018), adversarial instances depend not on the model but on the dataset itself. As a result, they are an inherent feature of the dataset and cannot be altered. Using the same dataset across different architectures may be virtually flawlessly translated to the other by harmful samples found in one architecture (Qin et al., 2020). To save time and money, it is more common to use pre-existing data

than to create brand-new ones. An important security flaw is created even when the production process is kept under wraps. Various forms of defence have risen to meet this challenge due to the ongoing arms race between attackers and defenders.

It is necessary to employ the boundary defence approach to safeguard the model from black-box assaults. These threats can be classified as either soft or hard labels and can also be classified as targeted or untargeted. With other methods, black box assaults can be mitigated in intensity. The technique utilised is the boundary defence algorithm, which protects the DNN from BlackBox attacks (both soft and hard label, both targeted and untargeted), which include both soft and hard label attacks.

With other methods, black box assaults can be mitigated in intensity. During the adversarial attack's optimisation phase, the adversarial samples reach the DNN's classification border, indicating that the adversarial assault has been successful (Zhao and Zeng, 2021). The border defence approach recognises these boundary samples and modifies their logits by adding white Gaussian noise. In addition, it prevents attackers from upgrading their malicious material and ensures minimal DNN performance deterioration throughout attacks. In addition to being simple to construct, the approach is easy to integrate into DNN models because it requires little or no code. Another advantage is that it is adaptable and reliable in its operation. During the optimisation phase of the adversarial attack, the holistic samples are located with DNN, indicating that they are adversarial samples. The boundary defence approach, which adds white Gaussian noise to the logits of these border samples, detects the presence of these border samples (Gu et al., 2021). Therefore, attackers will not be able to optimise their adversarial samples, resulting in only minor degradations in DNN performance for the DNN itself.

The algorithm has the following advantages:

- It can be implemented quickly and efficiently.
- It is simple to include in DNN models and requires little programming.
- It is adaptable and dependable in its operation.

During the adversarial attack's optimisation phase, the adversarial samples reach the DNN's classification border, indicating that the adversarial assault has been successful. Our technique identifies boundary samples with a classification confidence score less than the threshold and adds disturbances. In addition, it will prevent attackers from improving their adversarial samples and ensure a minimal level of DNN performance deterioration throughout attacks. There are different attack techniques: The white box method and the black box method. The major difference between the white box attack method and the black box attack method is that the attacker can get complete information about the DNN in the case of the white box attack method. In the case of the black box attack method, the attacker does not possess information regarding the DNN.

According to the data produced by the black-box neural networks, black-box attacks create adversarial samples. The research contributes significantly to a deeper comprehension of adversary durability against realistic attacks. The suggested approach efficiently generates adversarial instances, which are then included in the model's training process, mixing the adversary data samples with the source data across the training set. The most common technique to reduce the classification error is to produce the worst-case adversarial instances (those with the greatest loss) each time and then adjust the model parameters. This differs from adversarial training, in which fresh

adversary instances are continually created and blended with the original training set. Both classical ML algorithms and neural networks benefit from adversarial retraining.

To sum up, the following are discussed:

- The work on realistic adversarial attacks is discussed.
- We present a boundary defence approach that reduces black-box threats by taking advantage of adversarial optimisation, frequently needing samples on the classification boundaries.
- We provide broad insights into future research that might help bridge the gap between unrealistic and actual adversarial samples.

The organisation of the paper is as follows. Section 2 is based on a review of the literature. Section 3 gives preliminary knowledge of the technique, along with examples. In Section 4, the proposed work, along with details of image pre-processing, is given. Section 5 is based on the proposed algorithms. Results and analysis of the study are given in Section 6. Section 7 is followed by a conclusion.

2 Related work

In this portion, we have described the vulnerabilities in the DNNs regarding the generation and the effect of Adversarial examples in real-life use cases and the challenges regarding adversarial examples (Gu et al., 2021). Because of their extensive use and the fact that they deal with sensitive information, web browsers are considered a security risk and are prime targets for hackers. However, it is limited because many solutions are ineffectual, fall into the local optimum, and have a lengthy training period. A VDM is offered to evaluate the overall number of faults in the web browser while considering vulnerability intensity. It may be used to estimate the number of vulnerabilities, the pace at which vulnerabilities are discovered, the likelihood of vulnerabilities occurring in the future, risk evaluation, etc. It contains a new explanation for the seismic vulnerability index and a new expression.

Previous research by the authors examined the transferability of features among neural networks, while only a handful demonstrated the possibility of adversarial data being incorrectly categorised across models. According to their findings, once an adverse sample is produced for a specific perception, it is also likely to be misidentified in neural networks with alternative designs, which accounts for the attack's success. The quality and size of the substitute dataset that the adversary has gathered and the suitability of the adversarial network utilised to create adversarial samples determine how successful this type of assault is (Abbas et al., 2021).

The new formula is only associated with the vulnerability matrix, which shows the contribution of various damage grades with a few drawbacks, such as the problem of diminishing feature reuse, shortening the processing time at the expense of reducing detection accuracy, and difficulties in achieving better performance. It mainly talks about AJAX. AJAX applications are online applications that use this technology and have started a new trend in web applications. But with some limitations, as the computational complexity will increase as the number of hidden layers increases, the approach is a bit time-consuming, computationally intensive, and requires a relatively large size.

The modern network administrator must watch for publicly disclosed security flaws and respond appropriately (Singh and Jindal, 2019) using patching, configuration, and other techniques. Due to inadequate representations of spatial feature space, models are expensive to train. Due to the rapid development of corporate information management, numerous hackers illegally access firm data by taking advantage of information system flaws, causing the company to sustain considerable financial losses. As a result, governance mechanisms and vulnerability screening have emerged as crucial elements of commercial data security. In order to find and handle high-risk vulnerabilities in the information systems of the fabrication control and management control areas, this work discusses a new method for security technology based on publicly available info and designs a quick cross-regional vulnerability management platform. People's lifestyles are becoming increasingly dependent on various mobile phone applications (Apps), including those used for shopping, budgeting, and browsing the internet. These apps now necessitate extensive okay of the test dataset criteria, a lengthy training timeframe, and a fall into the local optimum (Herrera et al., 2020).

In the software area, trust is becoming increasingly vital. People encounter significant hurdles due to its complicated composite notion (Gu et al., 2021), particularly in today's dynamic and continuously evolving internet technology. The paper summarises the strategies for creating adversarial examples and highlights current discoveries on adversarial examples for DNNs. For hostile examples, further information on countermeasures is provided. In addition, three main issues in adversarial instances are highlighted, and proposed solutions. Many principles must be thoroughly described and answered, and issues like transferability and strong effectiveness evaluation. This study looks at deep convolutional neural networks and shows that attackers may readily create adversarial instances even if they don't know anything about the target network. Our native-based algorithm does not provide adversarial images that indicate the source data distribution when the pixels are reduced.

Taherdoost (2019) explains that DNNS are extremely effective in various applications. We propose a scalable mathematical approach that leads to limitations on the influence of these input disturbances on the network output, given that modest perturbations create adversarial instances to the input. This provided a mathematical approach for estimating a DNN model's adversarial susceptibility. We researched ways to make DNN models less sensitive to adversarial data manipulation assaults by considering constrained adversarial manipulation. This study examines five adversarial threats and four defence mechanisms in-depth on three driving models. (Singh and Jindal, 2019) experiments reveal that these models are very susceptible to adversary assaults, like classification models. This can lead to a significant risk to the security of automated cars and must be considered in practice, as none of them can adequately guard against all five attacks. This work examines adversarial attacks and countermeasures against automated driving models in depth. To that goal, we used three CNN-based driving model attacks to construct five adversarial attacks and four defensive strategies.

The danger of adversarial samples on deep self-learning models for remote sensing scene categorisation is extensively examined in this paper. Targeted and untargeted assaults create minor adversarial perturbations undetectable to the naked eye but can readily trick DL networks. Because many remote sensing jobs are closely tied to national defence security, stability and dependability are critical considerations (Sharma et al., 2020). We systematically examine the danger of adversarial instances on DNNS for distant sensing. Investigates (Singh and Jindal, 2019) hostile texts that might deceive

typical ML networks' sentiment analysis. The ensemble word addition (EWA) method is presented, which takes out and filters a limited amount of words with high attack potential and adds them after the original text. This paper presents a query and perturbation distribution-based improved black-box attack (IBBA) technique. This approach requires the attacked/disturbed models' top-1 label to produce adversarial instances. We optimise the algorithm's performance from two perspectives: query distribution and perturbation distribution, based on current black-box assaults. This paper proposes a query and distortion distribution-based black-box assault. Non-targeted and targeted assaults are both subjected to different query and perturbation distributions.

3 Preliminary knowledge

Black-box adversarial attacks use iterative optimisation and repeated searches to produce adversarial samples. It is proved difficult to defend DNNs against such attacks. We present a boundary defence approach that mitigates black-box threats by using the fact that adversarial optimisations frequently need samples on the classification boundary. Our technique identifies the boundary samples as having poor classification confidence, which adds white Gaussian distortion to their logits. The influence of the strategies on the classification accuracy of deep networks is studied theoretically. Rigorous testing reveals that the boundary defence approach can consistently defend against hard-label black-box threats. The suggested approach efficiently generates adversarial instances, which are included in the model's training phase, blending the adversarial data samples with the source data across the training set.

3.1 *Generation of adversarial examples*

The suggested approach efficiently generates adversarial instances, which are then included in the model's training process, mixing the adversarial data samples with the source data across the training set (Sharma and Kumar, 2022). The most common technique to reduce the classification error is to produce the worst-case adversarial instances (those with the greatest loss) each time and then adjust the model parameters. We try to develop the attacked samples that looked to be produced by adding information to the original dataset by looking at the local region of the attack image. In order to develop the mechanism to eliminate noise and fight against this type of assault, we typically require enough data. However, creating the adversarial picture for each image is very difficult and time-consuming. We must prepare the mechanism to reconstruct noise to supplement the dataset with just a few image pairs. The adversarial and clean images can be subtracted to produce noise. We train a mechanism to predict the noise and learn the distributions of noise based on the created real noise dataset. The clean image is combined with the noise data produced by the trained process to create new adversarial images.

This differs from adversarial training, in which fresh adversary instances are continually created and blended with the original training set. Both classical ML algorithms and neural networks benefit from adversarial retraining. The model's accuracy, achieved on these adversarial samples, an attack created through a collection of original samples, is a typical technique to test the resilience (Figure 2).

Figure 2 Flow to denote the adversarial attack

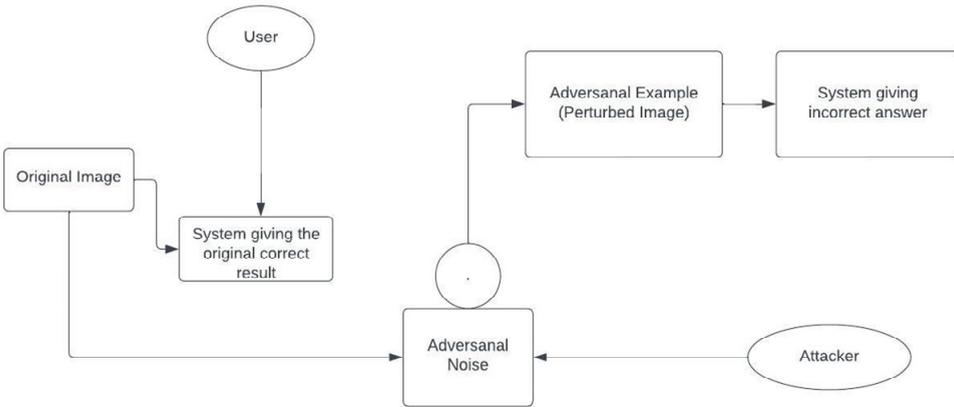
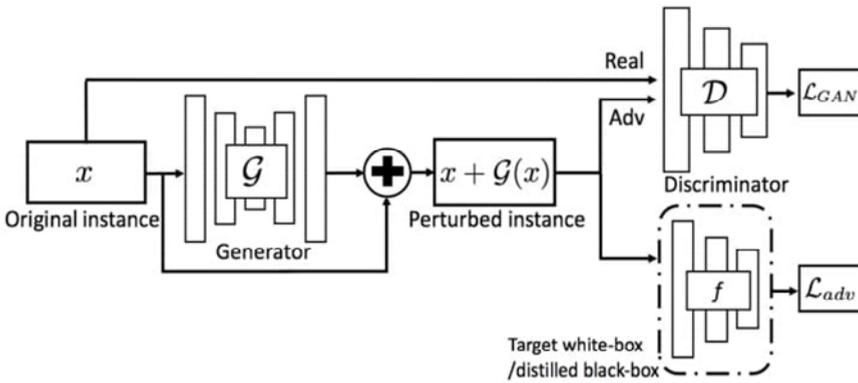


Figure 3 AdvGAN



The accuracy obtained from the model through the adversarial samples that the attack created from the source data samples is a typical technique to test resilience. Similarly, adversarial hardening or training techniques that teach models to make correct predictions on adversarial instances is a well-established method for hardening ML models. The attacker creates the adversarial samples using the fast gradient sign algorithm. The fast gradient sign approach forms an adversarial sample by leveraging the gradients present in the neural network. The method takes a new sample image that will take the maximum loss for an input sample by considering the input image’s loss gradients. This new picture is referred to as the adversarial sample. The following expression helps to summarise this:

$$adv_x + \epsilon * sign(\nabla_x J(\theta, x, y))$$

The generation of adversarial examples also takes place with the help of advance generative adversarial networks (AdvGAN). A generator G, a discriminator D, and a target neural network f are the three essential components of the AdvGAN. The generator G generates the perturbation G, which accepts the initial instance x as its input (x). The discriminator D, which will be utilised to differentiate the data produced from the original

version x , will then be supplied $x + G(x)$. D's goal is to make the produced instance unrecognisable from its source class in terms of data. We initially employ the white-box strategy to achieve our intention of deceiving a learning model, with the target model, in this case, being h . The gap between the forecasted and the target class (targeted attack), or the inverse of the distance between the predictions and the bottom class, is the input and output for h . The architecture of AdvGAN is depicted in the diagram (Figure 3).

3.2 Existence of adversarial examples

When constructing attacks and countermeasures in adversarial ML, the first step is generally to develop a knowledge of the presence and attributes of adversarial instances by reasoning why they impact the projection of ML methods. The problems faced by the models before the attack by the adversarial examples can be explained by the Clever-Hans effect. The phrase was popularised with the release of the CleverHans library. Hans is the name given to this phenomenon after a German-origin horse. His owner used to claim that Hans had the intellectual ability by having it answer mathematical questions by tapping its hoof the number of times that corresponded to the right answer. However, after repeated trials on Hans, scientists found that the horse was not solving mathematical problems but had evolved the capacity to recognise behavioural cues from the audience through claps and yells, which made the horse think to beat his hoof. In other words, Hans has created a technique of observing and analysing its environment in order to answer the questions, rather than an adaptive intelligence properly. Learning models, like Hans, can typically offer correct solutions to complicated issues like image recognition and classification, but they don't learn from training data, making them vulnerable to adversarial assaults.

4 Proposed work

4.1 Image pre-processing

Processing can help you improve the quality of your image or extract valuable information from it. Normally, downloaded image dataset includes complexity, inaccuracy, and inadequacy. Before building the model, we will pre-process the image dataset (cleaned and processed to the desired format) to achieve the desired results. We can eliminate undesired abnormalities and enhance certain features crucial for the program we are developing through pre-processing. Those qualities could alter according to the application. An image must be pre-processed for the software to work properly and deliver the required results. The major goal of processing the source data images is to improve image data (features) by suppressing undesired deformities and/or enhancing some critical image attributes so that ML and DL models can operate with better data. As a fundamental and crucial component of DNN models, data gathering is an extremely important and critical component. Because of the extensive usage of ML models, simply having a large dataset on a domain-specific task does not imply superior performance in that domain area (Figure 4).

Figure 4 Image pre-processing

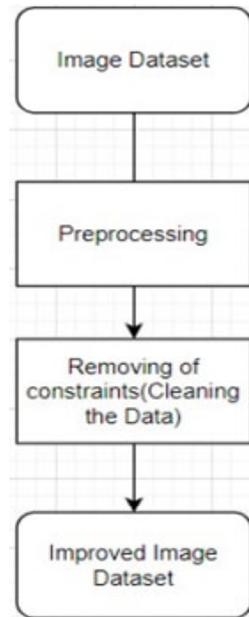
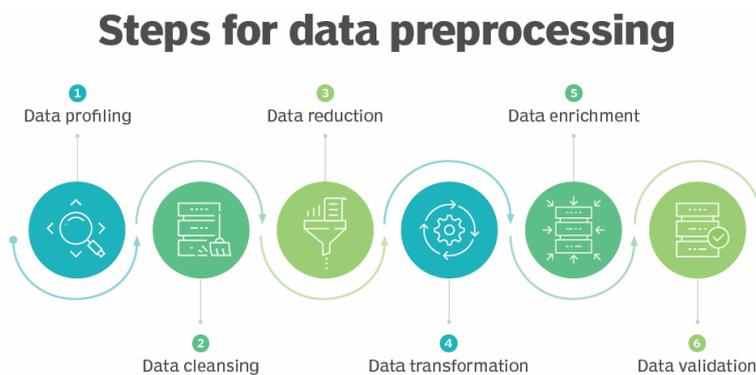


Figure 5 Data cleaning



Figure 6 Data pre-processing (see online version for colours)



Data cleansing can only obtain correct, consistent, and usable data. To avoid recurrence, errors and corruptions can be identified, corrected, deleted, or manually processed. The

goal of data cleaning is to fix any data that is erroneous, misleading, imperfect, improperly organised, reproduced, or even unimportant to the dataset’s purpose (Figure 5).

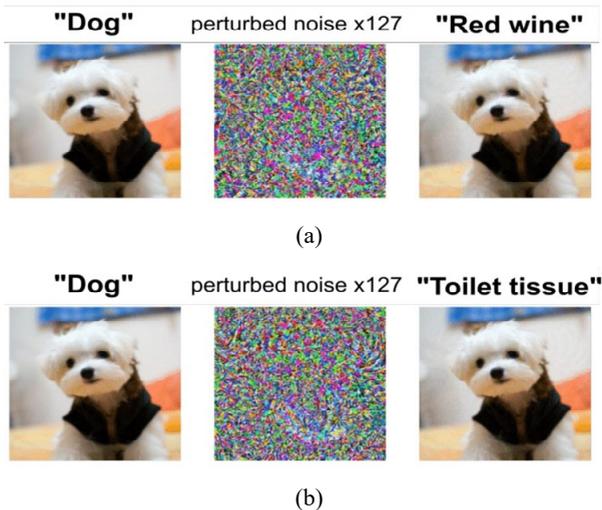
Processing raw data in an understandable format is known as pre-processing. There are several reasons we need to do this stage in data mining. To properly train a ML or data mining algorithm, thoroughly inspect the input data (Figure 6).

4.2 *Feature space attack*

The limitations imposed by domain attributes restrict many datasets used for classification or how features are designed (Yao. et al., 2019). The taken example cannot change individual features, and different relationships and correlations between different features occur. These limitations impose extra restrictions on the validity of adversarial instances (in addition to the distortion size). As a result, restricted assaults cause the original example to behave consistently with the limitations. The adversary can modify any attribute while conducting an unrestricted attack as long as the original and modified samples closely resemble each other. The traditional feature space attacks are known as unconstrained feature space attacks, whereas the domain-specific feature space attacks are called constrained feature space attacks (Figure 7).

When executing an unconstrained attack, the adversarial sample[4] can change any characteristic as far as the created adversarial sample is near enough to the original when several classification datasets are subject to limitations (for example, those imposed by intrinsic domain traits or the way features are created, constrained feature-space threats to arise.

Figure 7 Effect of distortions (see online version for colours)



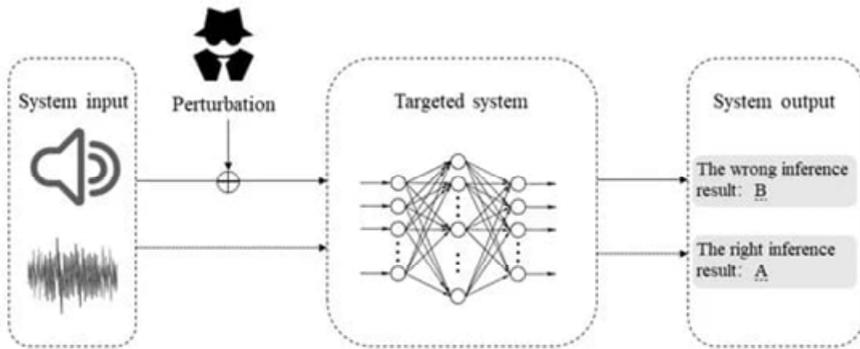
4.3 *Adversarial model attack*

The attacker is assumed to be aware of the model’s test distribution and have access to specific examples correctly classified by the method. We can examine either a black-box

or a white-box attack for, eg. The attacker in a traffic sign detection method can use the source data of traffic signs to initiate the attack. In a black-box exploit, the attacker is uninformed regarding the system’s exact specifics like source data, architecture, train dataset, etc. We suppose the attacker has the computing power to generate some actual adversarial sample that can trick the system. This indicates that a hacker has created a viable assault regarding the methodology discussed. As a result, starting with a correctly classified original example.

DL has taken centre stage in the evolution of ML and artificial intelligence. Current applications of neural networks include rebuilding the brain circuit (Yao et al., 2019), analysing ANN mutations, and evaluating data from particle accelerators. It has emerged as the driving force behind driverless cars, monitoring, and security applications in computer vision. Deep networks have been shown effective at solving complicated problems, but new research indicates that they are susceptible to small input disruptions. This minor perturbation can result in the classifier in the class label producing an incorrect result. While picture disturbances are frequently too minor for people to notice, they fully trick the DL model. DL applications face several dangers as a result of adversarial attacks. For instance, with face detection, the attacker is identified as a regular individual in order to steal the user’s personal information. In automated driving, a mistake in the roadside icon identification allows the control strategy to make a bad decision and produce bad behaviour consequently. It is essential to research the defence against hostile samples.

Figure 8 Generating adversarial attack



The accurate analysis of ML models’ resilience against adversarial attacks and the development of strategies for making models more resilient are required to allow their secure deployment in the real world. In order to determine the resilience of a system, When an attack generates adverse instances from a collection of initial samples, it is common practice to calculate the validity of the model. As an additional means of strengthening ML models, adversarial hardening uses training strategies that require the model to learn to make correct predictions for hostile scenarios, a well-established method of hardening ML models in general. The attacker can make adversaries conceptually similar to the source example but categorised incorrectly by the model. Although it is very expensive to carry out these adversarial attacks (particularly when used for improvising the system/model), a successful example can typically be obtained

by using them a few times. An example of the same is shown in the diagram below Figure 8.

5 Algorithm used

5.1 Boundary defence algorithm

To defend against black-box attacks, the model is protected using the boundary defence technique. These dangers might be classified as either mild or hard, targeted or untargeted. The severity of black-box attacks is lessened by using this technique. By taking advantage of adversarial optimisation typically requiring samples on the classification boundary, the boundary defence strategy reduces BlackBox threats. The boundary samples are identified as having low classification confidence by our method, which causes their logits to undergo white Gaussian distortion. The diagram below shows how the boundary defence system works (Figure 9).

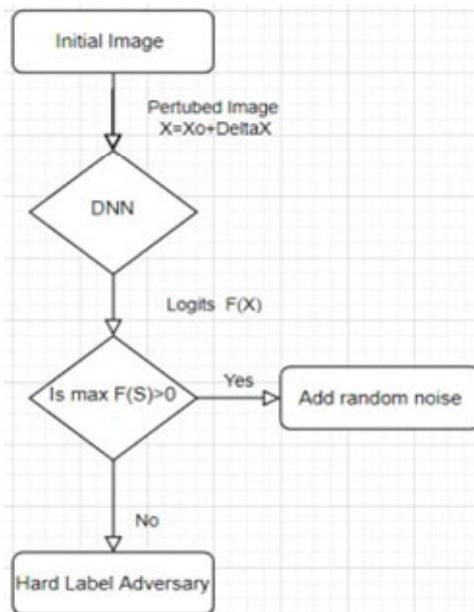
$$FBD(X) = F(S), \text{ if } \max F(S) > \theta$$

Else

$$F(S) + V$$

The adversarial samples are on the DNN’s classification boundary during the optimisation process of the adversarial attack. These border samples are detected by the boundary defence technique, which adds white Gaussian noise to their logits. As a result, attackers will be unable to optimise their adversarial samples, resulting in modest DNN performance deterioration.

Figure 9 Representation of black box attack and boundary defence



We first examine these algorithms from the standpoint of run time overhead, demonstrating that they cannot succeed within constrained time budgets, in addition to the standard criteria of attacker success and noise budget. The Actual Unencrypted Model, which we suggest, integrates the two adversarial attack types and quickly produces input-specific attacks. In particular, the model generates regions that act as a toasty for the internet connection during an unavailable generation stage. On the other hand, the online mode specialises in the content patch to the present input.

The advantages of the algorithm are:

- It can be implemented efficiently
- It can be easily added to DNN models with little code.
- It is flexible and works reliably.

5.2 Fast gradient sign algorithm

The attacker creates the adversarial samples using the fast gradient sign algorithm. The fast gradient sign approach forms an Adversarial sample by leveraging the gradients present in the neural network. The method takes a new sample image that will take the maximum loss for an input sample by considering the input image's loss gradients. This new picture is referred to as the adversarial sample. The following expression helps to summarise this:

$$adv_x = x + \epsilon * sign(\nabla_x J(\theta, x, y))$$

where

Adv(Real input) Adversarial image.

X Real input.

Y Real input label.

θ System parameters

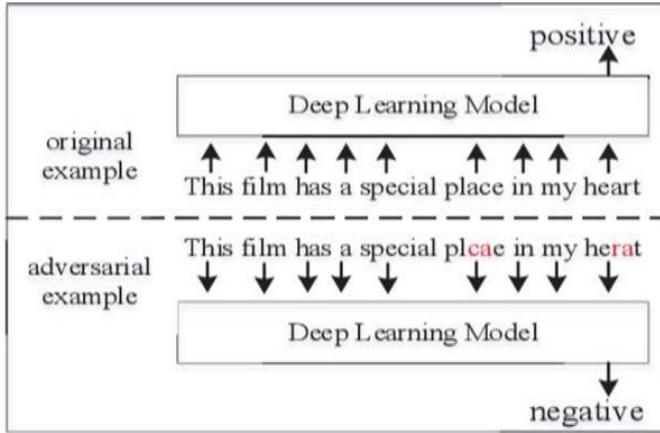
ϵ Multiplier variable to guarantee that the perturbations are minimal.

J Loss

5.3 Adversarial hardening

Development of domain-specific adversarial attacks that modify real-world objects to test the robustness of DNNs in real-world use cases. This may be accomplished by employing feature space assaults designed to be realistic. E.g., to attack a text classification model, the attacker can replace words with synonyms. This approach, however, has a considerable computational cost. To overcome this, adversarial hardening on unrealistic examples is used, which involves training techniques that teach models to generate accurate predictions on adversarial examples? (Figure 10).

Figure 10 Difference between the two samples



6 Result analysis

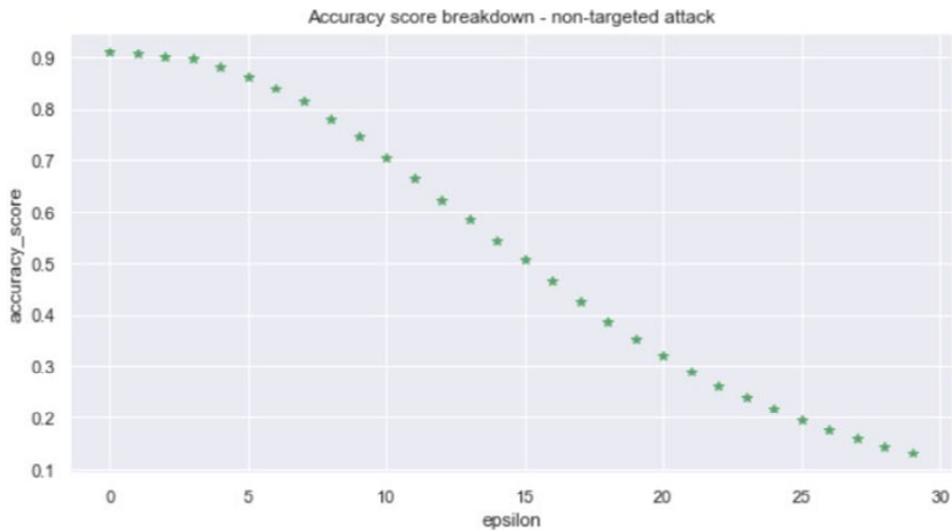
6.1 Hardware and software

Software, on the other hand, is hardware-executable instructions. Hardware is hard to change, while software is soft. Hardware usually follows software commands. Some computer systems use simple hardware, but others use hardware and software. i3 Dual-Core processor, Ethernet or Wi-Fi, 100 GB minimum hard disc, 200 GB recommended. Python, Anaconda, Jupyter Notebook, and TensorFlow require 8 GB RAM.

6.2 Results

The model is implemented using python: The first will be a tool for loading and parsing class labels from the dataset, including different pictures/images. Our next Python script will do basic picture classification (demonstrating ‘standard’ image classification) on the given dataset, showing correct results for each original input image. The dataset is later split into both the training and test sets. Even though the two pictures appear similar to the human eye, the final Python script will undertake an adversarial attack and produce an image that will purposefully mislead our model. The model receives an input image and classifies it, and we can generate an adversarial example that fools the model using the script’s output. The training and test sets are added with the distortions through small epsilon values and then train both the test and the training data. Different values of epsilon are tested, and the results are noted. The general trend shows that the higher the value of epsilon, the more distortion. The values are different once the distortion is added, which shows that the model is affected (Figures 11 and 12).

Figure 11 Higher the epsilon value, the lower the accuracy



As the epsilon’s value keeps increasing, the samples and their results become more and more inaccurate (Figure 13).

Figure 12 Difference between the actual values and the values after confusing the model (see online version for colours)

Out[45]:

	y_true	y_foiled	y_predicted	id
0	3	8	3	0
1	6	6	6	1
2	9	9	9	2
3	5	8	5	3
4	6	6	6	4

We compare the conventional’s effectiveness to that of the competitors. We evaluate the prediction performance in both attack scenarios while adjusting the time budget to account for varying confidence levels. Take note that the target accuracy of the categorisation for an adversarial instance is the credibility factor for the model. Finding the minimal fluctuation required to result in misclassification with a particular target confidence interval differs from another technique that creates adversarial cases within a certain perturbation level. Additionally, higher confidence takes longer to reach since it regulates the distance between the choice limit and the created adverse attack. This also leads to these affected samples becoming more and more visible to the naked eye, which contradicts the fact that they should be invisible to the human eye. As the number of distortions being added increases, the result of the outputs becomes more and more inaccurate, which could prove very efficient in confusing the system and producing incorrect results (Figures 14 and 15).

Figure 13 Distinct results for different epsilon values

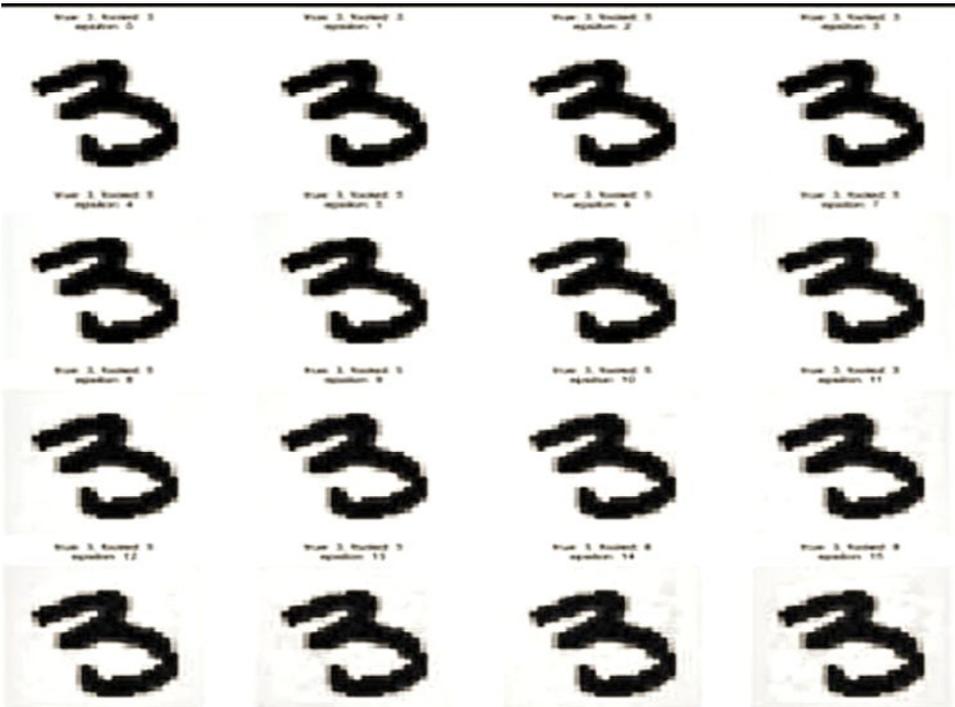


Figure 14 Different results on increasing distortions (see online version for colours)

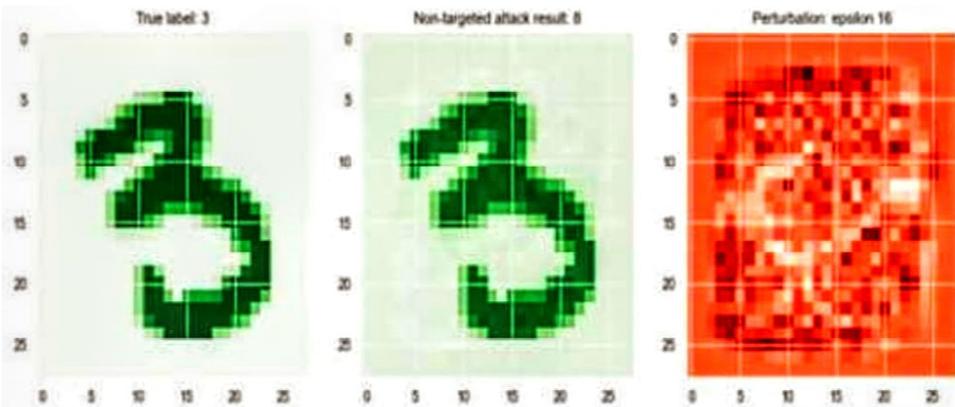


Figure 15 shows how often the result was wrong for a particular input.

Figure 16 shows the heat map for how often y_{true} was predicted as some y_{fooled} digit in percentage. The label holds the true digit and the other columns, all 784 pixels of an image with 28 times 28 pixels. Let's split our data into train and test. This way, we can measure our model performance on the test set and see how this score breaks down during the attack. The script moves on to perform the targeted attacks, producing the results for fooling the model for the intent class. An example of this is depicted in

Figure 14, where the intended targeted distortion is added to get the prediction for the desired result (Figure 17).

Figure 15 Plot predicted vs. count (see online version for colours)

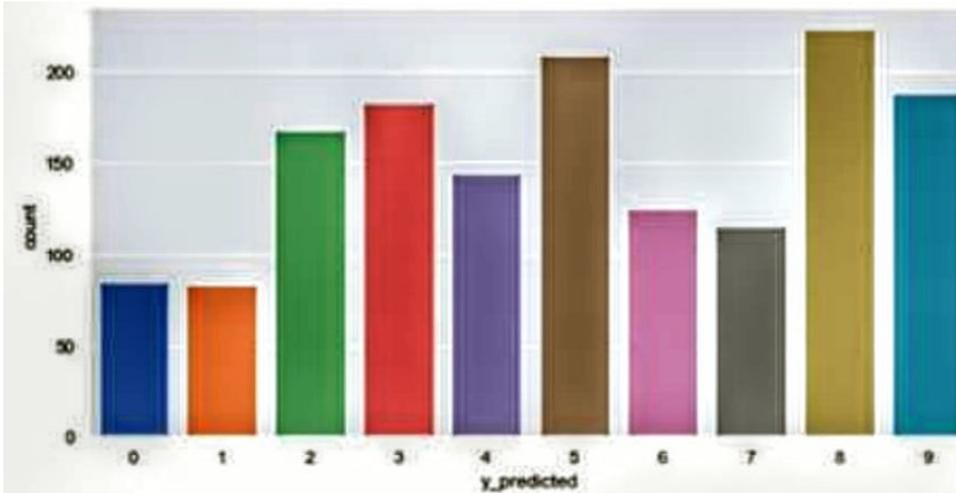
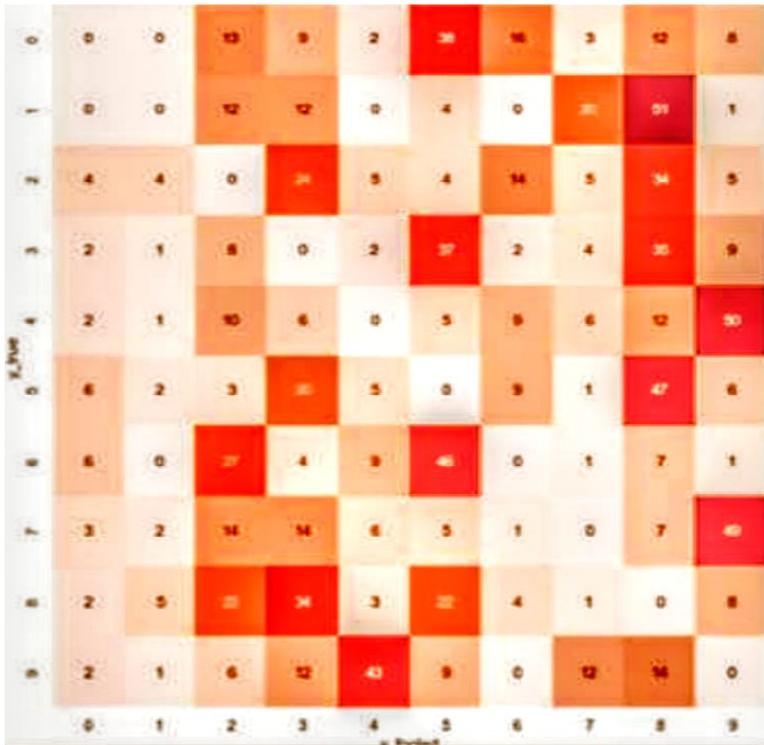


Figure 16 Heat map (see online version for colours)



We also try to evaluate the accuracy of both the intended and the unintended attacks and plot how the epsilon values affect their accuracy. Figure 18 helps in denoting the same.

Figure 17 Different values for the same lab (see online version for colours)

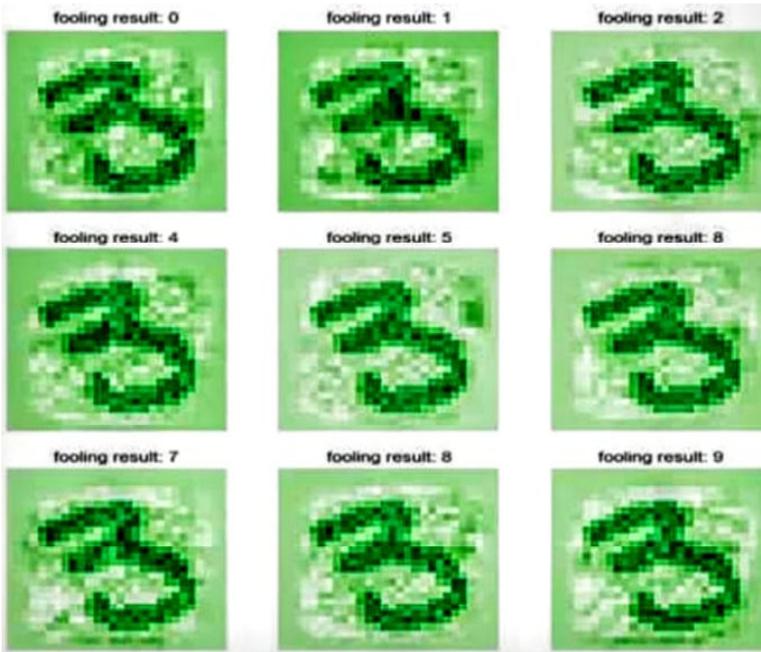
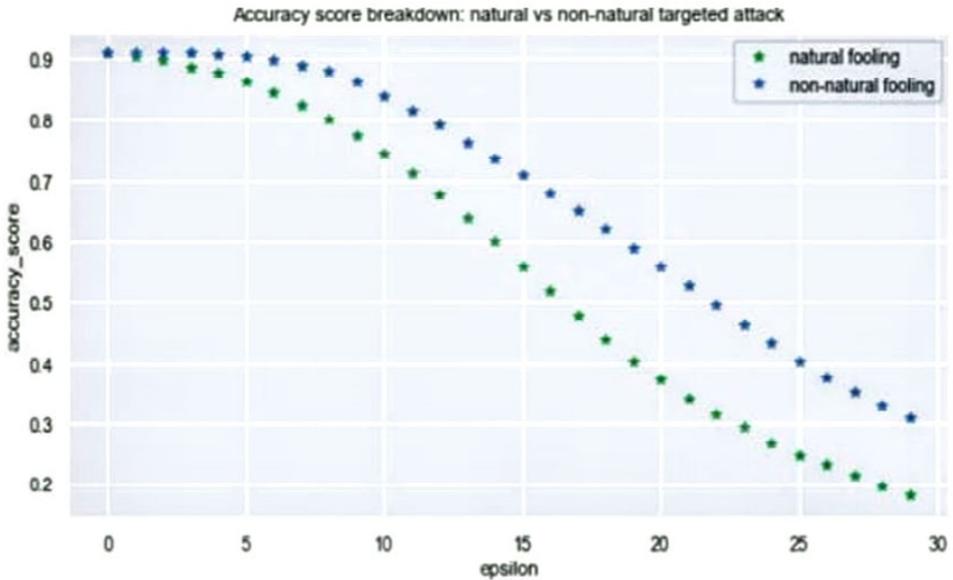


Figure 18 Plot epsilon vs. accuracy score (see online version for colours)



The plot above shows that as the value of epsilon increases, the model's accuracy decreases because it becomes much easier to fool the model; however, the accuracy decreases because the distortions become much more visible to the image. Adversarial Examples were developed and distributed onto a test network to validate the system. The primary objective of these expertly created examples is to utterly trick and confuse the model, allowing it to be tested for the susceptibility of deep self-learning systems. This technique paves the way for more in-depth knowledge of adversarial resistance in the face of real-world threats. Working with real-world instances and data, we seek to determine whether or not unreal adversarial samples that are very real could be used to safeguard frameworks from them. We use nodal dropouts from the first convolutional layer in deep self-learning systems topologies to identify weak and steady neurons. We build a connection between the neurons and the adversarial attacks in the network using an adversarial targeting method. Our findings reveal differences in use cases, with implausible examples succeeding as well as believable ones or offering marginal advantages over believable ones. To explain these results, the hidden representation of adversarial scenarios created with realistic and improbable attacks is studied. We have provided examples illustrating irrational samples' usage for comparable purposes.

6.2.1 Designing and evaluating defences

6.2.1.1 Step 1: Establish a threat model

The defence should always specify that it is resistant to adversarial attacks in its threat model. The threat model must be detailed in full, preferably under the taxonomy defined, so that reviewers and attackers may limit their assessments to the needs the defence declares to be secure.

6.2.1.2 Step 2: Use adaptive adversaries

A competent evaluation should put a defence by replicating adaptive adversaries that employ the threat model to create potent attacks. Without exception, all settings and assault situations that have the potential to overcome the defence should be considered. The outcomes of the tests do not lead to credible results and support the defence's arguments and resilience bounds. Therefore, an evaluation based purely on a non-adaptive adversary is of limited benefit. This can be described by the expression given below:

$$\min_{\theta} \max_{D(x, x') < \eta} J(\theta, x', y)$$

where x' is the adversarial input, the inner maximisation problem is handled using a well-designed adversarial approach like FGSM to locate the most effective adversarial samples.

6.2.1.3 Step 3: Perform sanity checks

Sanity checks are vital for detecting abnormalities and contradictory data that might lead to inaccurate conclusions by researchers. Some of the steps are:

- On authentic samples, report model accuracy: While protecting self-learning systems from adversarial instances is an important security concern, sacrificing a significant amount of legitimate data to improve the model’s robustness may be unreasonable in scenarios where there is a less chance of the adversarial attack, and there are fewer chances of getting the wrong result. For reactive defences, it is crucial to consider how the rejection of disturbing data affects the model’s accuracy on genuine samples.
- Sequential vs. iterative attack: Iterative attacks are more potent. Suppose adversarial samples generated by a sequential approach have a greater impact on classification models than those created using iterative algorithms. In that case, this might imply that the iterative attack’s execution is incorrect.
- Increase the distortion budget: Assaults that are permitted to cause more distortion in source samples are more likely to trick classifiers than attacks with lower distortions or disturbance budgets. As a result, if the attack’s success rate falls as the distortion budget rises, the attack method is probably faulty.

6.2.1.4 Step 4: Making the source code available

It is critical that all study materials and codes used to conduct the experimental studies and algorithms referred to in the paper be made publicly available online so that interested reviewers can replicate the initial work’s results and ensure their accuracy (Figure 19).

In the overall study, the following architecture of CN is used, whose e layered example is shown (Figure 20).

Figure 19 Steps in testing process (see online version for colours)

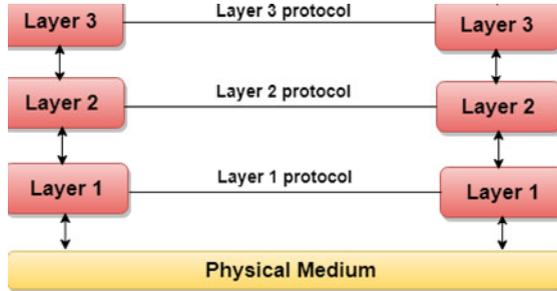


6.2.2 Merits and Demerits

After training it, the authors noticed that the model’s correctly categorised inputs were close to neighbouring adversarial inputs that had been misclassified. This indicates that learning models are inherently vulnerable to adversarial examples, regardless of the training method utilised. This has also refuted that stated hostile cases lie in a distinct range when opposed to legitimate databased on empirical findings. So, this study is

related to relate to previous research. It is clear that in order to obtain robustness, claim ml algorithms must generalise strongly, i.e., with the aid of robust optimisation. In essence, the authors found that the presence of adversarial cases is an inevitable byproduct of working in a statistical environment rather than always a flaw of particular classification techniques. The authors concluded that no viable methods achieve adversarial durability, mainly because the available datasets are too small to effectively train powerful models.

Figure 20 Architecture of CN (see online version for colours)



7 Conclusions

To defend against black-box assaults, we present an efficient and effective border defence mechanism in this study. By assessing classification confidence scores, this approach discovers boundary samples and adds random noise to the query results of these boundary samples. This study has proven that reducing the attack success rate classification accuracy degradation for image models reduces the attack success rate classification accuracy degradation. This simple and practical defence mechanism was examined and tested to show that it could effectively defend neural network models against state-of-the-art black-box assaults. To defend against black-box assaults, we present an efficient and effective border defence mechanism in this study. By assessing classification confidence scores, this approach discovers boundary samples and adds random noise to the query results of these boundary samples. This study has proven that reducing the attack success rate classification accuracy degradation for image models reduces the attack success rate classification accuracy degradation. This simple and practical defence mechanism was examined and tested to show that it could effectively defend neural network models against state-of-the-art black-box assaults. Finally, we hope our findings will spur more research into GAN-based adversarial examples. We believe our preliminary work on this innovative approach might call attention to the vast potential of image-to-image.

7.1 Future work

In plans, detecting hostile samples will still be a challenge in subsequent research. It is interesting to note that according to the authors' high predictive theorem, one learning algorithm is sufficient to truly portray any function. Thus, it makes intuitive sense that

strengthening the training phase is essential for fending off hostile samples. We intend to discuss cyclical auto-encoder and the drawbacks of an unsupervised approach. Additionally, as the majority of the models in our taxonomy have not previously been studied, this gives room for additional research in other adversarial contexts. In the future, this strategy will eventually serve as a common defence architecture. Our technique is effective against most generally thought-of-attack strategies, according to the experimental findings we obtained after testing it on various datasets and target models. The suggested strategy also offers many advantages over the most cutting-edge defence strategies. It is important to note that, despite being a realistic and straightforward defence mechanism, our system still presents certain practical challenges in terms of implementation and application. For instance, for complicated datasets, our experimental performance will suffer. Also, we'll concentrate on modifying the defence framework's network topology to enhance its performance in complex circumstances.

References

- Abbas, M., Dwivedy, S.K. and Narayan, J. (2021) 'Adaptive iterative learning-based gait tracking control for paediatric exoskeleton during passive-assist rehabilitation', *International Journal of Intelligent Engineering Informatics*, Vol. 9, No. 6, p.507, <https://doi.org/10.1504/ijiei.2021.10046403>.
- Bibi, A., Attique Khan, M., Younus Javed, M., Tariq, U., Kang, B-G., Nam, Y.H. and Sakr, R. (2022) 'Skin lesion segmentation and classification using conventional and deep learning based framework', *Computers, Materials and Continua*, Vol. 71, No. 2, pp.2477–2495, doi:10.32604/cmc.2022.018917.
- Feutrill, A., Ranathunga, D., Yarom, Y. and Roughan, M. (2018) 'The effect of common vulnerability scoring system metrics on vulnerability exploit delay', in *2018 Sixth International Symposium on Computing and Networking (CANDAR)*, pp.1–10, IEEE, <https://doi.org/10.1109/CANDAR.2018.00009>.
- Gu, Z., Hu, W., Zhang, C., Lu, H., Yin, L. and Wang, L. (2021) 'Gradient shielding: Towards understanding vulnerability of deep neural networks', *IEEE Transactions on Network Science and Engineering*, Vol. 8, No. 2, pp.921–932, <https://doi.org/10.1109/tNSE.2020.2996738>.
- Herrera, M.M., Morales, A.M. and Quiroga, J.M. (2020) 'Exploring the linkages between the patent applications and energy transitions: a system dynamics perspective', *International Journal of Intelligent Engineering Informatics*, Vol. 8, Nos. 5–6, p.526, <https://doi.org/10.1504/ijiei.2020.115739>.
- Huster, T.P., Chiang, C-Y.J., Chadha, R. and Swami, A. (2018) 'Towards the development of robust deep neural networks in adversarial settings', in *MILCOM 2018 – 2018 IEEE Military Communications Conference (MILCOM)*, pp.419–424, IEEE, <https://doi.org/10.1109/MILCOM.2018.8599814>.
- Jabeen, G. and Ping, L. (2019) 'A unified measurable software trustworthy model based on vulnerability loss speed index', in *2019 18th IEEE International Conference on Trust, Security and Privacy in Computing and Communications/13th IEEE International Conference on Big Data Science And Engineering (TrustCom/BigDataSE)*, pp.18–25, IEEE, <https://doi.org/10.1109/TrustCom/BigDataSE.2019.00013>.
- Koswara, K.J. and Dwi Wardhana Asnar, Y. (2019) 'Improving vulnerability scanner performance in detecting AJAX application vulnerabilities', in *2019 International Conference on Data and Software Engineering (ICoDSE)*, IEEE, pp.1–5, <https://doi.org/10.1109/ICoDSE48700.2019.9092613>.
- Lin, G., Zhang, J., Luo, W., Pan, L., Xiang, Y., De Vel, O. and Montague, P. (2018) 'Cross-project transfer representation learning for vulnerable function discovery', *IEEE Transactions on Industrial Informatics*, Vol. 14, No. 7, pp.3289–3297, <https://doi.org/10.1109/tii.2018.2821768>.

- Pahadiya, P., Saxena, S. and Vijay, R. (2021) 'Optimisation of thresholding techniques in de-noising of ECG signals', *International Journal of Intelligent Engineering Informatics*, Vol. 9, No. 5, p.487, <https://doi.org/10.1504/ijiei.2021.10044780>.
- Qin, J., Zhang, H., Guo, J., Wang, S., Wen, Q. and Shi, Y. (2020) 'Vulnerability detection on android apps—inspired by case study on vulnerability related with web functions', *IEEE Access: Practical Innovations, Open Solutions*, Vol. 8, pp.106437–106451, <https://doi.org/10.1109/access.2020.2998043>.
- Sharma, I. and Kumar, V. (2022) 'Multi-objective tunicate search optimisation algorithm for numerical problems', *International Journal of Intelligent Engineering Informatics*, Vol. 10, No. 2, p.119, <https://doi.org/10.1504/ijiei.2022.125859>.
- Sharma, N., Sharma, H., Sharma, A., and Bansal, J. C. (2020) 'Dung beetle inspired local search in artificial bee colony algorithm for unconstrained and constrained numerical optimisation,' *International Journal of Intelligent Engineering Informatics*, 8(4), p. 268. Available at: <https://doi.org/10.1504/ijiei.2020.10034275>.
- Shukla, A., Katt, B. and Nweke, L.O. (2019) 'Vulnerability discovery modelling with vulnerability severity', in *2019 IEEE Conference on Information and Communication Technology*, p.16, IEEE, <https://doi.org/10.1109/CICT48419.2019.9066187>.
- Singh, I. and Jindal, R. (2019) 'A survey on database intrusion detection: approaches, challenges and application', *International Journal of Intelligent Engineering Informatics*, Vol. 7, No. 6, p.559, <https://doi.org/10.1504/ijiei.2019.10026278>.
- Taherdoost, H. (2019) 'Electronic service quality measurement: development of a survey instrument to measure the quality of e-service', *International Journal of Intelligent Engineering Informatics*, Vol. 7, No. 6, p.491, <https://doi.org/10.1504/ijiei.2019.10026271>.
- Wu, Q., Liu, Y., Wen, B., Zhou, A. and Chen, Z. (2018) 'Research of cross-regional vulnerability governance for power grid enterprise information system', in *2018 11th International Symposium on Computational Intelligence and Design (ISCID)*, pp.181–184, IEEE, <https://doi.org/10.1109/ISCID.2018.10143>.
- Xie, Y., Gu, Z., Fu, X., Wang, L., Han, W. and Wang, Y. (2020) 'Misleading sentiment analysis: Generating adversarial texts by the ensemble word addition algorithm', in *2020 International Conferences on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData) and IEEE Congress on Cybermatics (Cybermatics)*, IEEE.
- Yao, X-Q., Sun, B-T., Yang, Z-L. and Cao, J-Q. (2019) 'A new method for vulnerability analysis and application in rural dwellings', in *2019 13th Symposium on Piezoelectricity, Acoustic Waves and Device Applications (SPAWDA)*, IEEE, <https://doi.org/10.1109/SPAWDA.2019.8681872>.
- Yuan, X., He, P., Zhu, Q. and Li, X. (2019) 'Adversarial examples: attacks and defenses for deep learning', *IEEE Transactions on Neural Networks and Learning Systems*, Vol. 30, No. 9, pp.2805–2824, <https://doi.org/10.1109/TNNLS.2018.2886017>.
- Zhao, W. and Zeng, Z. (2021) 'Improved black-box attack based on query and perturbation distribution', in *2021 13th International Conference on Advanced Computational Intelligence (ICACI)*, pp.117–125, IEEE, <https://doi.org/10.1109/ICACI52617.2021.9435907>.