

Assessment of White-Box Adversarial Attacks on Convolutional Neural Networks: Impact on Performance Metrics

Rakesh Podder

Department of Computer Science
Colorado State University

January 20, 2024

Abstract

The paper presents a comprehensive investigation into the susceptibility of Convolutional Neural Networks (CNNs) to white-box adversarial attacks. This research meticulously evaluates how various sophisticated attacks, like DeepFool, Carlini & Wagner (C&W), Projected Gradient Descent (PGD), Jacobian-based Saliency Map Attack (JSMA), and Basic Iterative Method (BIM), compromise key performance metrics of CNNs, such as accuracy, precision, and loss. By simulating these attacks, the paper offers deep insights into the vulnerabilities of CNNs in critical applications, underscoring the urgent need for robust defense strategies. The research is particularly relevant for fields where the reliability and security of machine learning models are crucial, such as autonomous vehicle navigation, financial systems, and healthcare diagnostics. The findings not only highlight the challenges in ensuring CNN security against these adversarial threats but also pave the way for developing more resilient machine learning models.

Key words: white-box adversarial attacks, CNN, test generation, accuracy & loss, minst.

1 Introduction

In the escalating landscape of cyber warfare, adversarial attacks on machine learning (ML) models have emerged as a sophisticated vector for undermining AI-driven systems. The inherent susceptibility of ML algorithms to specially crafted inputs that can lead to incorrect outputs, known as adversarial examples, has introduced a pressing challenge for the field of cybersecurity. The utilization of ML models in critical applications such as autonomous vehicles [5], financial systems [6], and healthcare diagnostics [4] has become prevalent. The trust and reliability placed on these ML systems by industries and end-users are fundamental to their widespread adoption.

Adversarial attacks have rapidly evolved from theoretical constructs to practical threats. These attacks leverage knowledge of the model's structure and data processing to introduce subtle perturbations, often imperceptible to humans but catastrophic for the model's decision-making accuracy. The consequences of successful adversarial attacks can range from trivial misclassifications to life-threatening situations. Therefore, understanding and mitigating these attacks are not just academic exercises; they are urgent requirements for the safe deployment of ML in real-world scenarios.

Despite the burgeoning research on adversarial ML, the gap between attack sophistication and defense mechanisms is widening. The challenge addressed in this research is twofold:

Assessment: There is a lack of comprehensive evaluation methodologies that can simulate an array of adversarial attack vectors effectively and measure their impact on diverse ML models.

Defense Enhancement: Current defense strategies often fail to generalize across different types of attacks and can degrade the performance of ML models in non-adversarial contexts.

The goal of this project is to conduct a systematic and thorough evaluation of various white-box adversarial attacks^[2] on datasets processed by Convolutional Neural Networks (CNNs). Through detailed experimentation in a range of test scenarios, the study aims to discern how different metrics are affected, specifically focusing on the accuracy, precision, and loss incurred by the model under adversarial conditions.

This research will delve into the intrinsic vulnerabilities of CNNs when exposed to white-box attacks—a scenario in which the attacker has complete visibility into the model’s architecture, parameters, and training data. By utilizing a comprehensive set of white-box attack strategies, the study will quantify the degradation of performance in CNNs and investigate the robustness of these networks against such exploits.

The paper’s expanded objective is to rigorously evaluate the impact of white-box adversarial attacks on the image quality and integrity as processed by Convolutional Neural Networks (CNNs), with a specific focus on the assessment of widely recognized image analysis metrics such as ERGAS, PSNR, SSIM, and SAM. By integrating these metrics into our test scenarios, we aim to quantify the degradation of image quality and the effectiveness of adversarial attacks in altering the image characteristics that CNNs rely on for accurate classification and analysis.

2 Related Work

Past research in the field of adversarial machine learning has made significant strides. The paper [1] provides a linearity-based theory for adversarial examples, proposes fast adversarial training, and refutes some alternative hypotheses. The view of adversarial examples as a fundamental property of linear models in high dimensions sparked significant subsequent research into understanding and improving model robustness. The paper [9] provides a comprehensive analysis of contemporary threats to machine learning systems and defenses across the system lifecycle. It highlights open challenges like physical attacks and efficient privacy preservation. The review of evaluations and future directions makes this a broad reference for machine learning security.

Xu et al. [8] provides a comprehensive review of adversarial attacks and defenses across multiple modalities including images, graphs, and text. It categorizes the main threat models, attack methods, defense strategies, and explanations for adversarial examples. The review covers key ideas like transferability of attacks and robustness vs accuracy tradeoffs. It summarizes the state-of-the-art in building robust deep learning models against adversarial threats. This relates to our work by providing background and motivation on improving model security through adversarial defenses. The review of attack and defense techniques across domains gives context for our proposed defense method in the image modality. This paper [3] outlines a set of principles for evaluating the robustness of machine learning defenses against adversarial examples, emphasizing the importance of a well-defined threat model and skepticism towards one’s own results. It advocates for rigorous testing using adaptive attacks, caution against security through obscurity, and the necessity of public code and model release for reproducibility. Additionally, the paper provides a checklist to avoid common pitfalls in such evaluations, encouraging comprehensive testing and comparison with existing work.

The paper [7] provides useful insights into factors influencing adversarial transferability and proposes a simple but effective smoothed gradient attack to enhance it. The attack has implications for evaluating model robustness.

This work aims to expand upon these foundations by specifically focusing on the impact of white-box adversarial attacks on CNN performance metrics. Unlike previous studies that broadly addressed adversarial threats in machine learning, our research delves into the detailed analysis of how these attacks affect CNNs, providing a more focused understanding of their vulnerabilities and potential defenses. This specificity in studying the direct effects of attacks on CNNs sets our work apart and underscores its importance in the broader context of machine

learning security.

3 Approach

Our study evaluates the resilience of Convolutional Neural Networks (CNNs) against a suite of sophisticated white-box adversarial attacks, each chosen for its relevance and prevalence in current literature and real-world applicability. The key approaches are:

- **Fast Gradient Sign Method (FGSM):** This method leverages the gradients of the loss with respect to the input image to create new images that are classified incorrectly. It's a straightforward yet powerful approach that demonstrates the vulnerability of neural networks to slight, often imperceptible, changes in the input data.
- **Jacobian-based Saliency Map Attack (JSMA):** A more refined technique that uses the model's Jacobian matrix to determine which pixels in the input image to alter to change the classification outcome. This method focuses on changing the least number of pixels, thus making the alterations less detectable.
- **Carlini & Wagner (C&W) Attack:** Recognized for its effectiveness, the C&W attack formulates adversarial example creation as an optimization problem. It aims to find the smallest perturbation that can mislead the model, ensuring that the adversarial examples remain as close as possible to the original images.
- **Projected Gradient Descent (PGD):** Often regarded as a more powerful version of FGSM, PGD applies the attack iteratively with small steps, making it more effective at finding adversarial examples that are misclassified by the network.
- **DeepFool:** This algorithm iteratively perturbs the input image in a way that is intended to cross the decision boundary of the classifier. It aims to be as efficient as possible, resulting in minimal perturbation.
- **Basic Iterative Method (BIM):** An extension of FGSM, BIM applies the gradient sign attack iteratively with small steps, allowing for finer control over the perturbation process and often resulting in more effective adversarial examples.

Each method will be systematically applied to a curated dataset of images processed by the CNN. The evaluation will be based on a comparison of key performance metrics pre- and post-attack, including accuracy, precision, loss, and image quality metrics such as ERGAS, PSNR, SSIM, and SAM. This comprehensive analysis will allow us to understand not just the effectiveness of each attack in degrading the CNN's performance but also the subtleties of how different types of perturbations can impact various aspects of the model's functionality.

By comparing these approaches under controlled conditions, the study aims to provide insights into the relative strengths and weaknesses of each method and establish a groundwork for developing more sophisticated defense mechanisms.

4 Evaluation Goals, Questions, and Metrics

The overarching goal of this evaluation is to measure and understand the impact of white-box adversarial attacks on the performance and reliability of Convolutional Neural Networks (CNNs) in the context of image processing. We aim to establish a robust framework for detecting vulnerabilities within CNNs and to quantify the effectiveness of adversarial attacks in degrading model performance. The following key objectives will guide our evaluation:

Goals:

1. To assess the impact of these attacks on the accuracy and integrity of the image classification process.
2. To identify the attack methodologies that result in the most significant degradation of performance metrics.
3. To provide insights into the development of more robust CNN architectures and training processes.

Questions:

1. How do various white-box adversarial attacks affect the classification accuracy of CNNs?
2. Which adversarial attack is most effective in inducing the highest error rates?
3. What is the correlation between the perceived image quality metrics (ERGAS, PSNR, SSIM, SAM) and the classification performance of CNNs under attack?
4. How does the iterative nature of certain attacks (e.g., BIM, PGD) compare to single-step attacks (e.g., FGSM) in terms of effectiveness?

Metrics:

To answer the above questions and achieve our goals, we will utilize a combination of traditional performance metrics and specialized image quality assessments:

- **Loss:** This is a measure of how well the model is performing from an error perspective. Specifically, it represents the “cost” incurred for inaccurate predictions. In your code, the loss is calculated using `'sparse_categorical_crossentropy'`, which is a common loss function for classification tasks. It compares the predicted probability distribution (output of the softmax function in the last layer) with the true distribution, where the true distribution is the label of the class that the input image belongs to. A lower loss indicates better performance of the model, as it means the model’s predictions are closer to the true labels.
- **Accuracy:** This metric measures the proportion of correctly predicted instances out of all predictions made. In a classification task like MNIST (which involves classifying images of handwritten digits into 10 classes, from 0 to 9), accuracy is calculated by the number of images correctly classified divided by the total number of images classified. Higher accuracy means the model has better predictive performance.
- **ERGAS (Relative Dimensionless Global Error in Synthesis):** A global measure of image fidelity, with lower values indicating better synthesis quality.
- **PSNR (Peak Signal-to-Noise Ratio):** A measure of peak error, with higher values indicating smaller differences between original and perturbed images.
- **SSIM (Structural Similarity Index):** A perception-based model that considers changes in texture, contrast, and luminance.
- **SAM (Spectral Angle Mapper):** A measure of the spectral similarity between two images, with lower values indicating higher similarity.

By analyzing these metrics before and after the application of adversarial attacks, we will be able to paint a comprehensive picture of CNN robustness, identifying both strengths and vulnerabilities. This information will be critical for advancing the field of machine learning security and for developing systems that can trustfully be deployed in the real world.

5 Study Design

The study design is meticulously crafted to ensure a systematic and thorough evaluation of the susceptibility of CNNs to white-box adversarial attacks. Our design incorporates specific tasks, tools for generating and testing adversarial examples, and a detailed schedule of measurements.

5.1 Tasks

- **Model Selection and Preparation:** A set of CNNs will be chosen for their prevalence and relevance in the field. These will be prepared by training on benchmark datasets.
- **Dataset Curation:** A comprehensive dataset will be curated, taking into account the variety and complexity required to challenge the CNNs under test.
- **Adversarial Example Generation:** Leveraging state-of-the-art adversarial techniques, we will systematically generate examples that aim to mislead the CNNs while preserving image quality.
- **Model Retraining (Optional):** Depending on the initial results, an iterative process of model retraining with adversarial examples may be conducted to assess any changes in model robustness.

5.2 Tools

- **Attack Algorithms:** Utilization of established libraries to implement FGSM, JSMA, C&W, PGD, DeepFool, and BIM attacks.
- **Evaluation Framework:** A framework will be either developed or adopted from existing solutions like CleverHans, ART, or Foolbox for automating the evaluation process. Here we will be using ART library mostly.
- **Statistical Analysis Software:** Python packages such as SciPy, Seaborn will be employed for in-depth data analysis.

5.3 Measurements

- **Pre-attack Performance:** Baseline measurements of accuracy, precision, recall, F1 score, ERGAS, PSNR, SSIM, and SAM will be established.
- **Post-attack Performance:** The same metrics will be re-assessed post-adversarial attack to evaluate the impact.
- **Adversarial Success Rate:** The rate at which adversarial inputs successfully deceive the CNN will be meticulously recorded.
- **Robustness Threshold:** Identification of the minimal perturbation magnitude necessary to compromise the model.

5.4 Running Example:

FSGM Attack

The given code is a Python script using the TensorFlow and Adversarial Robustness Toolbox (ART) libraries to train a neural network on the MNIST dataset, then generate adversarial examples using the Fast Gradient Sign Method (FGSM) attack, and finally evaluate the model's performance on the adversarial examples. Here's a step-by-step explanation:

1. Import necessary libraries: The necessary libraries and modules are imported, including TensorFlow, NumPy, ART, and matplotlib.
2. Load MNIST dataset: The MNIST dataset, which is a dataset of 60,000 28x28 grayscale images of handwritten digits, along with a test set of 10,000 images, is loaded and pre-processed (normalized).
3. Create the model: A simple neural network model is created using TensorFlow’s Keras API. The model consists of a Flatten layer that converts each 28x28 image into a 784 element vector, followed by a Dense layer with 128 nodes, a Dropout layer that randomly sets 20% of the input units to 0 during training, and a final Dense layer with 10 nodes corresponding to the 10 possible digits (0-9).
4. Compile and train the model: The model is compiled using the ‘adam’ optimizer and the ‘sparse_categorical_crossentropy’ loss function, and then trained on the training images and labels for 5 epochs.
5. Evaluate the model: The trained model is evaluated on the test images and labels to get the baseline loss and accuracy.
6. Create a TensorFlowV2Classifier: ART’s TensorFlowV2Classifier is created using the trained model. This classifier will be used as the input to the FGSM attack.
7. Create a FastGradientMethod attack: ART’s FastGradientMethod attack is created using the classifier and an epsilon value of 0.1.
8. Generate adversarial examples: Adversarial examples are generated from the test images using the FGSM attack.
9. Evaluate the model on adversarial examples: The model is evaluated on the adversarial examples to see how well it performs in the presence of an attack.
10. Display an original and adversarial image: The original and adversarial images for the first test example are displayed side by side using matplotlib.

The script trains a simple neural network on the MNIST dataset, generates adversarial examples using the FGSM attack, evaluates the model’s performance on the adversarial examples, and displays an original and adversarial image. This demonstrates how easy it is for an attacker to fool a neural network using adversarial examples and highlights the importance of developing robust models that can withstand adversarial attacks.

Metric	FGSM_ART
Accuracy	0.10
loss	3.08
ERGAS	27.083559
PSNR	22.267371
SSIM	(0.882, 0.945)
SAM	0.280872

Table 1: Metric Evaluation table for FGSM Adversarial Attacks

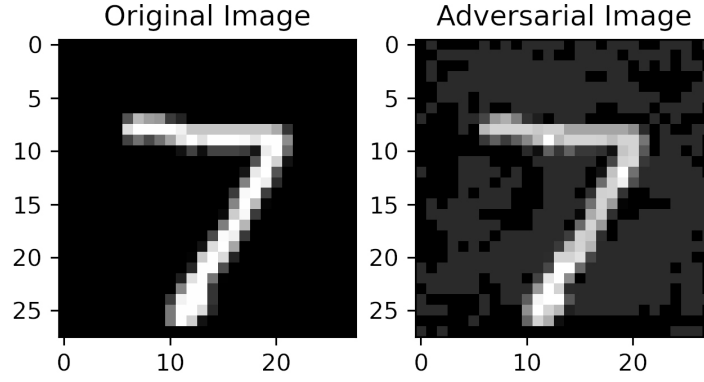


Figure 1: Comparison b/w Original image and modified image by FSGM Attack

6 Empirical Evaluation

We ran our experiment on a 3.1 GHz Dual-Core Intel Core i5 processor, with 8 GB 2133 MHz, and LPDDR3 memory. The empirical evaluation section of the paper encompasses several parts, each focusing on different adversarial attacks and their impacts. Here are short summaries for each part:

6.1 Metric Evaluation

This part evaluates the impact of the attack on CNNs. Metrics like accuracy, loss, ERGAS, PSNR, SSIM, and SAM are used to assess the effect on image quality and model performance.

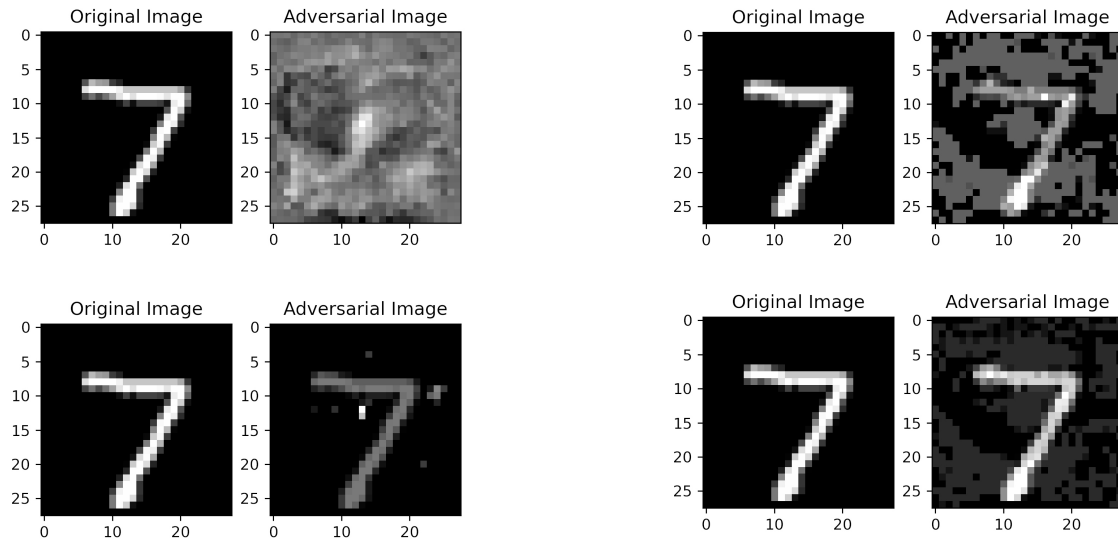


Figure 2: DeepFool, PGD, JSMA, BIM original and attack image.

6.2 ISP & BCC

Input and State Variable The input variable are:

- No. of Neurons: Numeric value.

Metric	FGSM	DF	C&W	PGD	JSMA	BIM
Accuracy	0.10	0.00977	0.977	0.0135	0.075	0.02
loss	3.08	2782.88	0.074	80.89	0.975	16.5
ERGAS	27.08	3254.48	78.336	79.73	29.72	26.56
PSNR	22.5	2.617	13.407	13.52	23.56	22.55
SSIM	(0.89, 0.94)	(-0.13, -0.24)	(0.632, 0.71)	(0.62, 0.71)	(0.93, 0.934)	(0.896, 0.95)
SAM	0.28	1.365663	0.748	0.749	0.236	0.27

Table 2: Metric Evaluation table for all Adversarial Attacks

- Dropout Rate: Numeric percentage.
- NB_Classes: Numeric value.
- Optimizer: Alphanumeric value.

The state variable is the dataset (MINST).

6.2.1 Input Space Partitioning

Table 3 displays the input state Partitioning table.

Variables	Characteristics	Partitions	Values
No. of Neuron (N)	Numeric: $0 < N < finitevalue$	a1 = true a2 = false	N = 128 N = NULL
Dropout rate (R)	Numeric: $0 < R \leq 1$	b1 = true b2 = false	R = 0.2 R = - 0.3
NB_classes (nb)	Numeric: $0 < N < finitevalue$	c1 = true c2 = false	$nb \geq 2$ $nb < 2$
Optimizer (O)	Alphanumeric	d1 = true d2 = false	O = "Adadelta " O = NULL
Dataset Size (len)	Length: Integer	e1 = nonEmpty e2 = Empty	$len > 0$ $len = 0$

Table 3: ISP for the Evaluation

6.2.2 Base choice coverage (BCC)

Test	Partition 1	Partition 2	Partition 3	Partition 4	Partition 5	Oracle
T_1 (base)	a1	b1	c1	d1	e1	Pass
T_2	a2	b1	c1	d1	e1	Fail
T_3	a1	b2	c1	d1	e1	Fail
T_4	a1	b1	c2	d1	e1	Fail
T_5	a1	b1	c1	d2	e1	Fail
T_6	a1	b1	c1	d1	e2	Fail
T_7	a2	b2	c1	d1	e1	Fail
T_8	a2	b1	c2	d1	e2	Fail
T_9	a1	b2	c1	d2	e2	Fail
T_{10}	a1	b2	c2	d2	e1	Fail

Table 4: BCC Table for Evaluation

6.3 Testing

6.3.1 Testing the ISP on FGSM

in this section various test suite created to validate and evaluate the various type attacks on cnns. A test coverage had be shown below for FGSM

Test Case	Partition 1	Partition 2	Partition 3	Partition 4	Partition 5	loss	accuracy
TC_1	a1 = 128	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	3.08	0.105
TC_2	a1 = 100	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	3.10	0.09
TC_3	a1= 150	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	3.04	0.108
TC_4	a1 = 500	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	2.9	0.12
TC_5	a1=1000	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	2.5	0.13

Table 5: Testing the ISP on FGSM varying No. of Neurons.

Test Case	Partition 1	Partition 2	Partition 3	Partition 4	Partition 5	loss	accuracy
TC_1	a1 = 128	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	3.08	0.105
TC_2	a1 = 128	b1 = 0.02	c2 = 10	d1 = Adadelta	e1 = MINST	3.19	0.099
TC_3	a1 = 128	b1 = 0.001	c2 = 10	d1 = Adadelta	e1 = MINST	3.8	0.06
TC_4	a1 = 128	b1 = 0.5	c2 = 10	d1 = Adadelta	e1 = MINST	2.88	0.11
TC_5	a1 = 128	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	2.75	1.24

Table 6: Testing the ISP on FGSM varying dropout rate.

Test Case	Partition 1	Partition 2	Partition 3	Partition 4	Partition 5	loss	accuracy
TC_1	a1 = 128	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	3.08	0.105
TC_2	a1 = 128	b1 = 0.2	c2 = 2	d1 = Adadelta	e1 = MINST	3.10	0.09
TC_3	a1 = 128	b1 = 0.2	c2 = 50	d1 = Adadelta	e1 = MINST	3.04	0.108
TC_4	a1 = 128	b1 = 0.2	c2 = 100	d1 = Adadelta	e1 = MINST	2.9	0.12
TC_5	a1 = 128	b1 = 0.2	c2 = 200	d1 = Adadelta	e1 = MINST	2.5	0.13

Table 7: Testing the ISP on FGSM varying nb_classes.

7 Discussion

Based on the empirical evaluation data, we observe that various adversarial attacks on Convolutional Neural Networks (CNNs) demonstrate different impacts when applied to the MNIST dataset. This variety in impact is shown by significant fluctuations in key performance indicators such as accuracy, loss, and image quality metrics including ERGAS, PSNR, SSIM, and SAM. These findings underscore the susceptibility of CNNs to sophisticated adversarial methods. It's evident that the development of more robust defense mechanisms is crucial to ensure the reliability and security of CNN applications across different domains.

From Table 2's comparison, we can see that the DeepFool attack results in the lowest accuracy among all the attacks, indicating a substantial amount of data loss. In contrast, the FGSM attack shows a comparatively better synthesis quality as it records the lowest ERGAS value among the tested attacks. The JSMA attack stands out with the highest peak error. However, it is interesting to note that in terms of SSIM values, which reflect changes in texture, contrast, and luminance, the impacts are quite similar across all attacks on the MNIST dataset. A particularly noteworthy observation is that despite JSMA's high peak-to-peak error, it has the lowest SAM value, suggesting that it maintains higher similarities between the attacked

Test Case	Partition 1	Partition 2	Partition 3	Partition 4	Partition 5	loss	accuracy
TC_1	a1 = 128	b1 = 0.2	c2 = 10	d1 = Adadelta	e1 = MINST	3.08	0.105
TC_2	a1 = 128	b1 = 0.2	c2 = 10	d1 = adam	e1 = MINST	8.6	0.074
TC_3	a1 = 128	b1 = 0.2	c2 = 10	d1 = SGD	e1 = MINST	3.85	0.104
TC_4	a1 = 128	b1 = 0.2	c2 = 10	d1 = Adagrad	e1 = MINST	3.29	0.092
TC_5	a1 = 128	b1 = 0.2	c2 = 10	d1 = RMSProp	e1 = MINST	8.48	0.073

Table 8: Testing the ISP on FGSM varying Optimizer.

image and the original image. This characteristic of JSMA could be critical for understanding and countering adversarial attacks on CNN models.

Further, the evaluation extends to different test cases for the FGSM attack as detailed in Tables 5, 6, 7, and 8. These tables illustrate how the accuracy and loss metrics of a CNN are influenced by varying factors such as the number of neurons, the dropout rate, the choice of optimizer, and the number of classes. This analysis is pivotal in understanding how CNNs react to adversarial attacks under different configurations and settings. Such insights are vital for developing CNN models that are not only effective in their predictive accuracy but also resilient against sophisticated adversarial attacks, ensuring their reliability in practical applications.

From the evaluation we can achieve the goals by getting the answers of the questions from study design section. For example, Different adversarial attacks affect CNN classification accuracy variably. For instance, the DeepFool attack significantly reduces accuracy due to substantial data loss. The most effective attack in inducing high error rates is the DeepFool attack, as evidenced by its low accuracy scores. There’s a notable correlation between image quality metrics (ERGAS, PSNR, SSIM, SAM) and CNN performance under attack. Attacks like FGSM, which have lower ERGAS values, indicate better image synthesis quality despite the adversarial modification. Iterative attacks like BIM and PGD are more effective compared to single-step attacks like FGSM. This is due to their ability to apply perturbations iteratively, allowing for a more refined and effective attack process.

8 Threats to validity

External Validity: The generalizability of the results is a concern. While the study effectively demonstrates the impact of adversarial attacks on CNNs using the MNIST dataset, the same results might not hold for more complex or larger datasets. Although this limitation exists, the study still provides valuable insights into the vulnerabilities of CNNs, which can inform further research in more varied contexts.

Internal Validity: The causal relationship between the treatment (adversarial attacks) and the observed effects (changes in performance metrics) needs careful examination. Other factors, such as the specific architecture of the CNN or the nature of the dataset, might also influence the outcomes. Ensuring that the observed effects are solely due to the adversarial attacks is crucial for accurate conclusions.

Construct Validity: This concerns whether the study accurately measures what it intends to. In this context, it’s about ensuring that the performance metrics like accuracy, loss, ERGAS, PSNR, SSIM, and SAM truly reflect the impact of adversarial attacks on the CNNs. Misinterpretation or inaccuracies in these measurements could lead to incorrect conclusions about the vulnerability and robustness of the CNNs to adversarial attacks.

9 Conclusion

Convolutional Neural Networks (CNNs) are highly vulnerable to various sophisticated adversarial attacks. These attacks significantly degrade the performance of CNNs, as evidenced by

changes in key metrics such as accuracy, precision, loss, and image quality assessments. The results highlight the importance of developing more robust defense mechanisms for CNNs, particularly in applications where security and reliability are crucial. This research contributes to a deeper understanding of CNN vulnerabilities and paves the way for future advancements in creating more resilient neural network models.

The main lessons from the study are that Convolutional Neural Networks (CNNs) exhibit significant vulnerability to a range of adversarial attacks, with these attacks leading to notable degradation in performance metrics like accuracy, loss, and image quality. The research underscores the importance of developing more resilient CNN architectures and defense mechanisms to counteract these vulnerabilities, particularly in critical applications where CNN reliability is paramount. The findings provide valuable insights for future research aimed at enhancing the security and robustness of CNNs against sophisticated adversarial threats.

References

- [1] N. Carlini, A. Athalye, N. Papernot, W. Brendel, J. Rauber, D. Tsipras, I. Goodfellow, A. Madry, and A. Kurakin, "On evaluating adversarial robustness," *arXiv preprint arXiv:1902.06705*, 2019.
- [2] J. Ebrahimi, A. Rao, D. Lowd, and D. Dou, "Hotflip: White-box adversarial examples for text classification," *arXiv preprint arXiv:1712.06751*, 2017.
- [3] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [4] K. Kalaiselvi and M. Deepika, "Machine learning for healthcare diagnostics," *Machine Learning with Health Care Perspective: Machine Learning and Healthcare*, pp. 91–105, 2020.
- [5] A. Qayyum, M. Usama, J. Qadir, and A. Al-Fuqaha, "Securing connected & autonomous vehicles: Challenges posed by adversarial machine learning and the way forward," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 2, pp. 998–1026, 2020.
- [6] K. Singh, M. Hasan, R. P. Rajendran *et al.*, "Opportunities and challenges of ai/ml in finance," *The Impact of AI Innovation on Financial Sectors in the Era of Industry 5.0*, pp. 238–260, 2023.
- [7] L. Wu and Z. Zhu, "Towards understanding and improving the transferability of adversarial examples in deep neural networks," in *Asian Conference on Machine Learning*. PMLR, 2020, pp. 837–850.
- [8] H. Xu, Y. Ma, H.-C. Liu, D. Deb, H. Liu, J.-L. Tang, and A. K. Jain, "Adversarial attacks and defenses in images, graphs and text: A review," *International Journal of Automation and Computing*, vol. 17, pp. 151–178, 2020.
- [9] M. Xue, C. Yuan, H. Wu, Y. Zhang, and W. Liu, "Machine learning security: Threats, countermeasures, and evaluations," *IEEE Access*, vol. 8, pp. 74 720–74 742, 2020.