

Defense Implemented: Randomization

Randomization Approach: Random Cropping

Metric	Clean	No Defense Attack	w/ Defense Attack
Loss	0.00395	2.79	0.0284
Accuracy	100%	22%	98%

Example Misclassifications:

No misclassified images for stage: Clean
Attack: fgsm_bim_attack
Dataset: MNIST
Training Epochs: 10
Trained Clean Images: 54210
Test Images: 140
Accuracy: 1.00
Precision: 1.00
Recall: 1.00
F1-score: 1.00
ROC AUC score is not defined for a single class.

Number of misclassified images for No Defense Attack: 50
Attack: fgsm_bim_attack
Dataset: MNIST
Training Epochs: 10
Adversarial Training Images: 27105
Test Images: 140
Accuracy: 0.22
Precision: 0.38
Recall: 0.33
F1-score: 0.35
ROC AUC Score: 1.00

Misclassifications:

- 0 -> 6: 13
- 0 -> 5: 7
- 0 -> 4: 3
- 0 -> 7: 2
- 0 -> 2: 15
- 0 -> 9: 8
- 0 -> 8: 1
- 0 -> 3: 1

Number of misclassified images for w/ Defense Attack: 1
Attack: fgsm_bim_attack
Dataset: MNIST
Training Epochs: 10
Retrained Clean and Adversarial Images: 81315
Test Images: 140
Accuracy: 0.98
Precision: 0.99
Recall: 0.99
F1-score: 0.99
ROC AUC Score: 1.00

Misclassifications:
0 -> 6: 1