

Defense Implemented: Adversarial Training

Metric	Clean	No Defense Attack	w/ Defense Attack
Loss	0.0156	4.04	0.0298
Accuracy	98%	20%	98%

Example Misclassifications:

Number of misclassified images for Clean: 1  
Attack: fgsm\_pgd\_attack  
Dataset: MNIST  
Training Epochs: 10  
Trained Clean Images: 54210  
Test Images: 140  
Accuracy: 0.98  
Precision: 0.99  
Recall: 0.99  
F1-score: 0.99  
ROC AUC Score: 1.00

Misclassifications:  
0 -> 7: 1

Number of misclassified images for No Defense Attack: 51  
Attack: fgsm\_pgd\_attack  
Dataset: MNIST  
Training Epochs: 10  
Adversarial Training Images: 27105  
Test Images: 140  
Accuracy: 0.20  
Precision: 0.35  
Recall: 0.31  
F1-score: 0.33  
ROC AUC Score: 1.00

Misclassifications:  
0 -> 6: 9  
0 -> 8: 6  
0 -> 2: 26  
0 -> 5: 6  
0 -> 7: 2  
0 -> 9: 2

Number of misclassified images for w/ Defense Attack: 1

Attack: fgsm\_pgd\_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 -> 8: 1