

Defense Implemented: Adversarial Training

Metric	Clean	No Defense Attack	w/ Defense Attack
Loss	0.00491	0.217	0.000675
Accuracy	100%	89%	100%

Example Misclassifications:

No misclassified images for stage: Clean  
Attack: cw\_pgd\_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: 7  
Attack: cw\_pgd\_attack  
Dataset: MNIST  
Training Epochs: 10  
Adversarial Training Images: 15108  
Test Images: 140  
Accuracy: 0.89  
Precision: 0.95  
Recall: 0.93  
F1-score: 0.94  
ROC AUC Score: 1.00

Misclassifications:

0 -> 6: 2  
0 -> 5: 1  
0 -> 3: 1  
0 -> 2: 1  
0 -> 9: 2

No misclassified images for stage: w/ Defense Attack  
Attack: cw\_pgd\_attack  
Dataset: MNIST  
Training Epochs: 10

Retrained Clean and Adversarial Images: 69318

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.

-----