Defense Implemented: Input Transformation

Input Transformation Approach: Differential Privacy

| Metric   | Clean   | No Defense Attack | w/ Defense Attack |
|==========|=========|===================|===================|
| Loss     | 0.00302 | 0.16              | 0.0195            |
| Accuracy | 100%    | 89%               | 98%               |

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: 15154
Test Images: 140
Accuracy: 0.89
Precision: 0.95
Recall: 0.93
F1-score: 0.94
ROC AUC Score: 1.00
--------------------------------------------------------

Misclassifications:
0 -> 9: 1
0 -> 6: 4
0 -> 2: 2


--------------------------------------------------------
Number of misclassified images for w/ Defense Attack: 1
Attack: cw_pgd_attack
Dataset: MNIST
Training Epochs: 10

Retrained Clean and Adversarial Images: 69364
Test Images: 140
Accuracy: 0.98
Precision: 0.99
Recall: 0.99
F1-score: 0.99
ROC AUC Score: 1.00
----------------------------------------------------------

Misclassifications:
0 -> 9: 1