Defense Implemented: Input Transformation

Input Transformation Approach: Differential Privacy

| Metric | Clean | No Defense Attack | w/ Defense Attack |
|----------|---------|--------------------|--------------------|
| Loss | 0.00659 | 0.118 | 0.00358 |
| Accuracy | 100% | 92% | 100% |

Example Misclassifications:

---------------------------------------------------------
No misclassified images for stage: Clean
Attack: pgd_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: 5
Attack: pgd_bim_attack
Dataset: MNIST
Training Epochs: 10
Adversarial Training Images: 11463
Test Images: 140
Accuracy: 0.92
Precision: 0.96
Recall: 0.96
F1-score: 0.96
ROC AUC Score: 1.00
---------------------------------------------------------

Misclassifications:
0 -> 6: 2
0 -> 7: 1
0 -> 9: 1
0 -> 4: 1


---------------------------------------------------------
No misclassified images for stage: w/ Defense Attack
Attack: pgd_bim_attack
Dataset: MNIST

Training Epochs: 10
Retrained Clean and Adversarial Images: 65673
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.
---------------------------------------------------------