Failure Rate Against Adversarial Attacks CWDeep FGM-step size:0.001-norm:1 FGM-step size:0.001-norm:2 - 1.00e+01 FGM-step size:0.001-norm:inf FGM-step size:0.01-norm:1 FGM-step size:0.01-norm:2 FGM-step size:0.01-norm:inf Failure Rate (mislabelled samples per second) FGM-step size:0.1-norm:1 - 1.00e-01 FGM-step size:0.1-norm:2 FGM-step size:0.1-norm:inf HSJ PGD-step size:0.001-norm:1 PGD-step size:0.001-norm:2 PGD-step size:0.001-norm:inf - 1.00e-03 PGD-step size:0.01-norm:1 PGD-step size:0.01-norm:2 PGD-step size:0.01-norm:inf PGD-step size:0.1-norm:1 PGD-step size:0.1-norm:2 PGD-step size:0.1-norm:inf - 1.00e-05 Patch-scale max:0.03 Patch-scale max:0.1 Patch-scale max:0.25 Patch-scale max:0.5 Patch-scale max:1.0 - 1.00e-07 Pixel-th:1 Pixel-th:16 Pixel-th:2 Pixel-th:4 Pixel-th:8 Thresh 1.00e-09 Conf-cutoff:0.9 Conf-cutoff:0.99 Conf-cutoff:0.5 Sigmoid- γ :100- β :0.0 TVM-prob:0. Gauss-In-sigma:0.9 Round-decimals: Round-decimals: abel-sigma:0.9 Round-decimals: Round-decimals: Gauss-Out-scale:0. Gauss-Out-scale:0. Label-sigma:0. Gauss-Out-scale:0 Label-sigma:0 Sigmoid- γ :1- β : $\bar{1}$ - γ :100

Defenses

Attacks