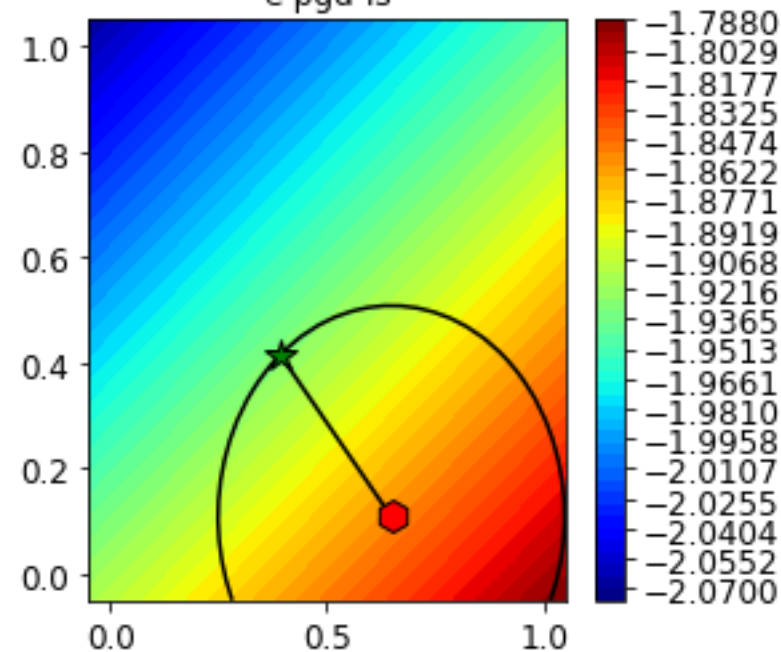
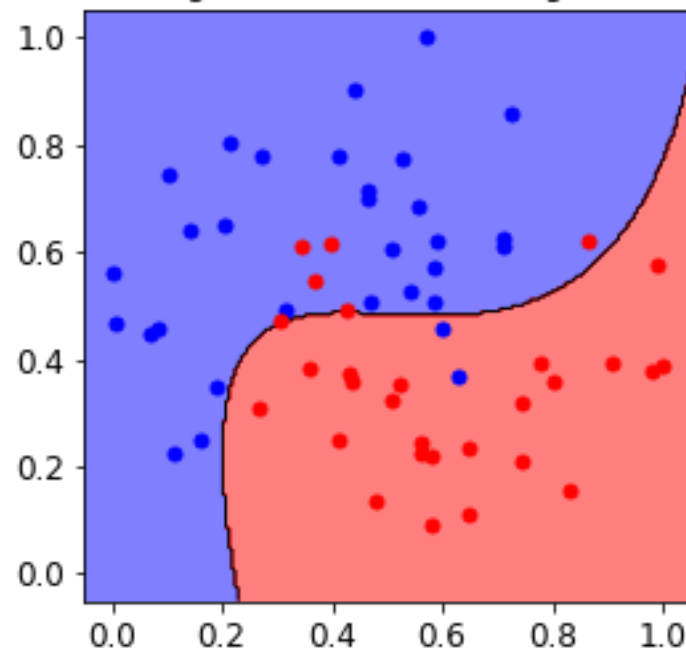


e-pgd-ls

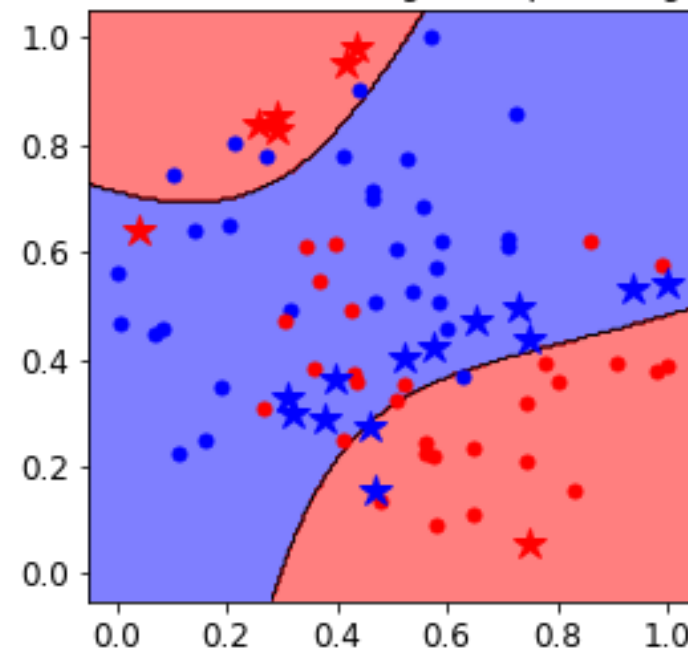


(a)

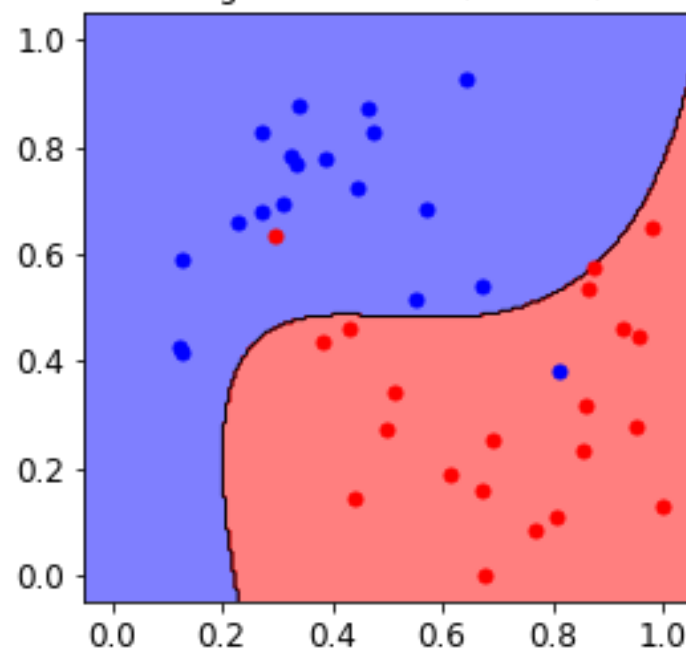
Original classifier (training set)



Poisoned classifier (training set + poisoning points)

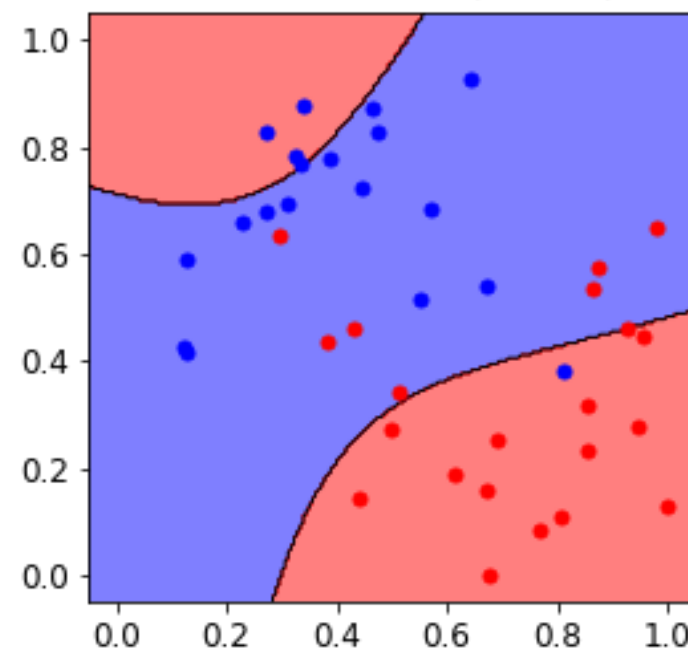


Original classifier (test set)



Accuracy on test set: 95.00%

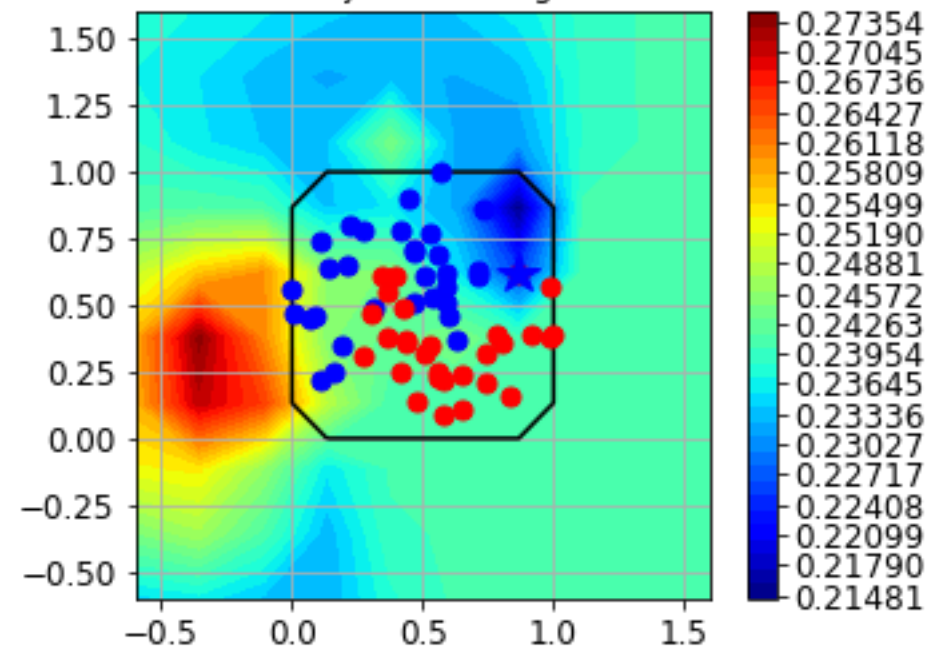
Poisoned classifier (test set)



Accuracy on test set: 67.50%

(c)

Attacker objective and gradients



(b)