Constrained Adversarial Attacks and Defenses

T. Simonetto, S. Dyrmishi, S. Ghamizi, M. Cordy, Y. Le Traon University of Luxembourg



Adversarial examples are subject to domain constraints













With $H(x) \neq H(x + \delta)$

With $L_p(x, x + \delta) < \epsilon$

 $x + \delta \in \mathcal{X}_{\Omega}$ (domain constraints)

Refused transaction

Noise / Additional transactions

Accepted transaction

Adversarial examples are inputs carefully designed to cause erroneous predictions in machine learning systems.

To fool real-world systems, they must respect domain constraints.

Constraints as penalty functions ...

Example of constraint:

 $installment = loam_amount \times \frac{int_rate \times (1 + int_rate)^{term}}{(1 + int_rate)^{term} - 1}$

Constraints	formul	lae

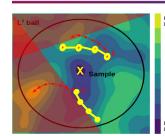
 $\begin{array}{l}
\omega_1 \wedge \omega_2 \\
\omega_1 \vee \omega_2 \\
\psi \in \Psi = \{\psi_1, \dots \psi_k\} \\
\psi_1 \leq \psi_2 \\
\psi_1 < \psi_2 \\
\psi_1 = \psi_2
\end{array}$

Penalty function

 $\begin{array}{l} \overline{\omega_{1} + \omega_{2}} \\ \min(\omega_{1}, \omega_{2}) \\ \min(\{\psi_{i} \in \Psi : \mid \psi - \psi_{i} \mid \}) \\ max(0, \psi_{1} - \psi_{2}) \\ max(0, \psi_{1} - \psi_{2} + \tau) \\ \mid \psi_{1} - \psi_{2} \mid \end{array}$

... used in 2 constrained attacks

Objective: for x find δ



High C-PGD (gradient based)

 $\nabla_{x^t} loss(h(x^t), y) - \nabla_{x^t} penalty(x^t)$

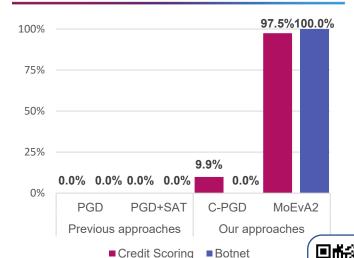
MoEvA2 (genetic based)

 $g_1(x) \equiv h(x)$

 $g_2(x) \equiv \, L_p(x-x_0)$

 $g_3(x) \equiv \sum_{\omega_i \in \Omega} penalty(x, \omega_i)$

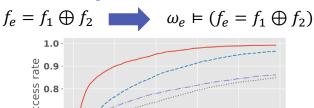
Attacks success rate



Success = Misclassification & Constraints Satisfaction

Defenses

Constraints augmentation:



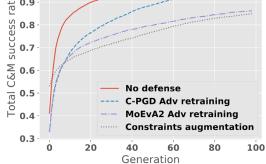


Figure 1: Success rate of MoEvA2 against the original model and the three defended models, over the generations.

