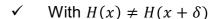
Towards Adaptive Attacks on Constrained Tabular ML

T. Simonetto¹, S. Ghamizi², M. Cordy¹ ¹University of Luxembourg ²RIKEN AI, Luxembourg Institute of Science and Technology

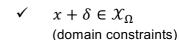


Adversarial examples are subject to domain constraints

Objective: for x find δ



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





Noise / Additional transactions

Accepted transaction

Refused 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 a penalty function

Example of constraint:

$$\omega_1 \equiv installment = loan_amount \times \frac{int_rate \times (1 + int_rate)^{term}}{(1 + int_rate)^{term} - 1}$$

Penalty function:

$$penalty(x, \omega_1) = |installment - \\ (loan_amount \times \frac{int_rate \times (1 + int_rate)^{term}}{(1 + int_rate)^{term} - 1})|$$

Loss function:

TabTr.

RLN

$$\mathcal{L}'(x) = \mathcal{L}(x, y, h, \Omega) = l(h(x), y) - \sum_{\omega_i \in \Omega} penalty(x, \omega_i)$$

Constrained Adaptive PGD

Gradient step momentum:

$$\begin{split} z^{(k+1)} &= P_{\mathcal{S}}(x^{(k)} + \eta^{(k)}(\nabla \mathcal{L}'(x^{(k)})) \\ x^{(k+1)} &= R_{\Omega}(P_{\mathcal{S}}(x^{(k)} + \alpha \cdot (z^{(k+1)} - x^{(k)}) \\ &\qquad \qquad + (1 - \alpha) \cdot (x^{(k)} - x^{(k+1)}))) \end{split}$$

Step size adaptation

$$\sum_{i=w_{j-1}}^{w_j-1} \mathbf{1}_{\mathcal{L}'(x^{(i+1)}) > \mathcal{L}'(x^{(i)})} < \rho \cdot (w_j - w_{j-1})$$

Repair operator: project values of feature in equality constraints.

Robust Accuracy (%)

LCLD 80 69.5 68.3 67 66.4 67.4 40 0 **URL** 92.5 93.3 94.4 93.4 100 60 20

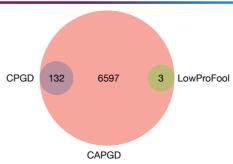
VIME

■Clean NLPF = CPGD ■ CAPGD

STG

TabNet

CAPGD subsumes other gradient attacks



Number of successful adversarial examples

Next steps

CAA: An adaptive attack combining strong complementary attacks.

TabularBench: a comprehensive benchmark of robustness of tabular deep learning classification models.

