#### NeurlPS 2024

# Constrained Adaptive Attack: Effective Adversarial Attack Against Deep Neural Networks for Tabular Data

Thibault Simonetto<sup>1</sup>, Salah Ghamizi<sup>2, 3</sup>, Maxime Cordy<sup>1</sup>

<sup>1</sup>University of Luxembourg, Luxembourg

<sup>2</sup>Luxembourg Institute of Science and Technology, Luxembourg

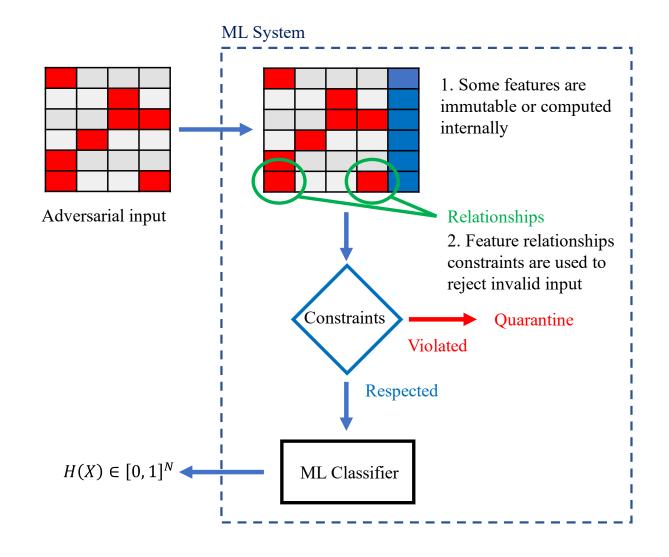
<sup>3</sup>RIKEN Center for Advanced Intelligence Project, Tokyo, Japan







## Adversarial examples in tabular data



## Relation constraints on feature space

Finance:

 $avg\ transaction\ amount\ \leq max\ transaction\ amount$ 

$$installment = loan\_amount \times \frac{int\_rate \times (1+int\_rate)^{term}}{(1+int\_rate)^{term}-1}$$

## **Problem formulation**

Given a classification model H,

a maximum perturbation  $\epsilon$  under a  $L_p$  distance

a set of constraints  $\Omega$ 

Objective of **constrained** adversarial attacks, for clean sample x find perturbation  $\delta$ :

✓ With 
$$H(x) \neq H(x + \delta)$$

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

$$\checkmark x + \delta \vDash \Omega$$

# Feature relation constraints as a penalty function

#### **Objective**

 $penalty(x,\omega)$  how far is x from satisfying  $\omega$ 

 $penalty(x,\omega) = 0$  iff  $\omega(x) =$ True

#### **Example**

 $\omega_1 \equiv avg\_transaction \leq max\_transaction$ 

 $penalty(x, \omega_1) = max(0, avg\_transaction - max\_transaction)$ 

#### **Mapping to penalty functions**

Constraint	Penalty function
$\psi_1 = \psi_2$	$\mid \psi_1 - \psi_2 \mid$
$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
$\psi_1 < \psi_2$	$max(0,\psi_1-\psi_2+\tau)$
$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
$\omega_1 \vee \omega_2$	$\min(\omega_1,\omega_2)$

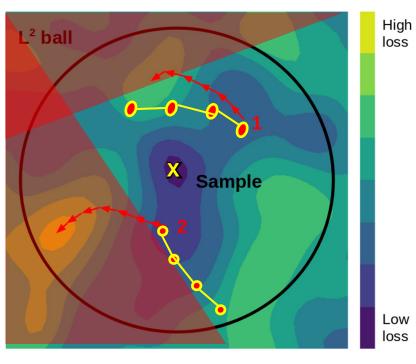
## **Constrained Adaptive PGD**

#### **Projected Gradient Descent (PGD)**

$$x^{(k+1)} = P_{\mathcal{S}}(x^{(k)} + \eta \nabla l(h(x), y))$$

#### **Gradient loss + Constraints regularization**

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



# **Constrained Adaptive PGD**

#### **Gradient loss + Constraints regularization**

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

#### **Constrained Gradient Descent**

$$z^{(k+1)} = P_{\mathcal{S}}(x^{(k)} + \eta^{(k)}(\nabla \mathcal{L}'(x^{(k)}))$$

#### Adaptive Step size $\eta$

$$\eta^{(0)} = 2\epsilon, \ \eta^{(k+1)} = \left\{ \begin{array}{l} \eta^{(k)}/2, & \text{if } \mathcal{L}' \text{ does not decrease} \\ \eta^{(k)}, & \text{otherwise} \end{array} \right\}$$

#### Gradient step momentum $\alpha$

$$x^{(k+1)} = R_{\Omega} (P_{\mathcal{S}}(x^{(k)} + \alpha \cdot (z^{(k+1)} - x^{(k)}) + (1 - \alpha) \cdot (x^{(k)} - x^{(k-1)})))$$

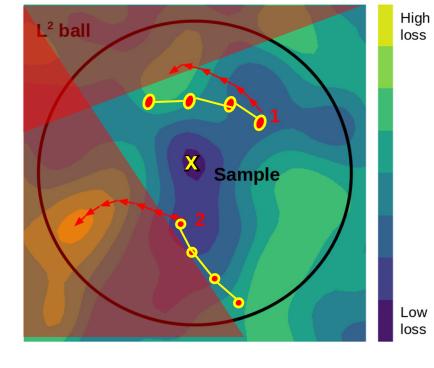
$$+ (1 - \alpha) \cdot (x^{(k)} - x^{(k-1)})))$$

Current gradient direction

Previous perturbation direction

#### Repair operator $R_{\Omega}$

$$\omega_1 \equiv installment = loan\_amount \times \frac{int\_rate \times (1 + int\_rate)^{term}}{(1 + int\_rate)^{term} - 1}$$



## **Experimental settings**

#### **Datasets:**



Credit scoring Lending Club Loan Data



Botnet detection CTU



URL phishing URL

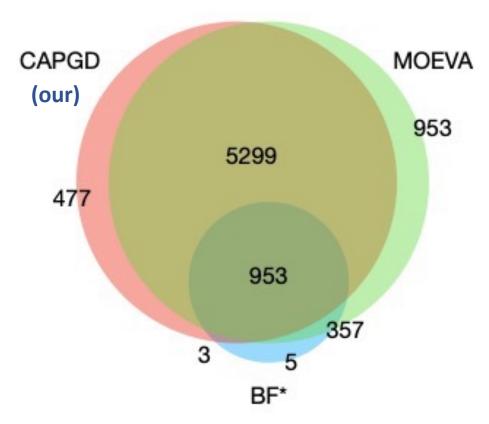


ICU survival WIDS

**Models:** 5 Neural network architectures

- 2 Regularizations
- 2 Transformers
- 1 Semi-supervised

## Are these attacks complementary?



Set and number of successfull attacks across all models and datasets

#### **Insight:**

- 1. MOEVA and CAPGD are complementary
  - 2. Together they subsume BF\*

## **Constrained Adaptive Attack**

We propose CAA

CAPGD

MOEVA

Efficiency

Of examples can be generated by CAA

## Impact of CAA

### Robust accuracy



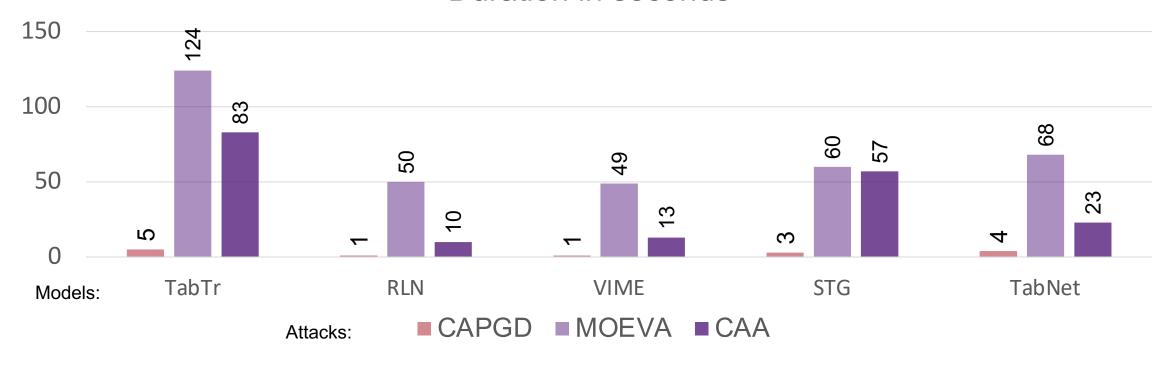
Dataset: Lending Club Loan data

#### **Insight:**

CAA effectively combines the benefits of CAPGD and MOEVA

# **Efficiency of CAA**

#### Duration in seconds

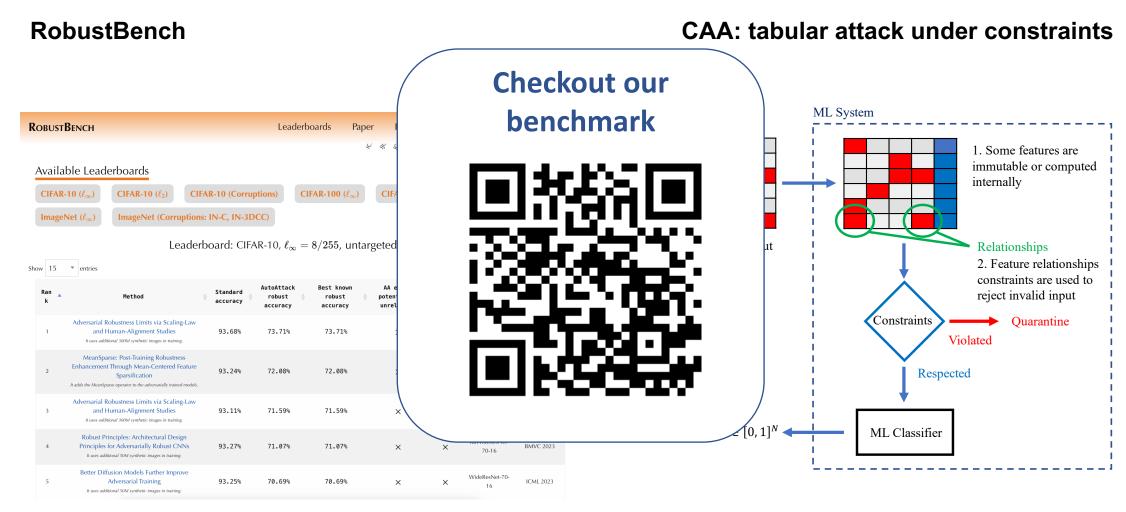


Dataset: Lending Club Loan data

#### **Insight:**

CAA reduces execution costs by up to 5 times compared to MOEVA

## Conclusion



https://robustbench.github.io/

https://github.com/serval-uni-lu/tabularbench