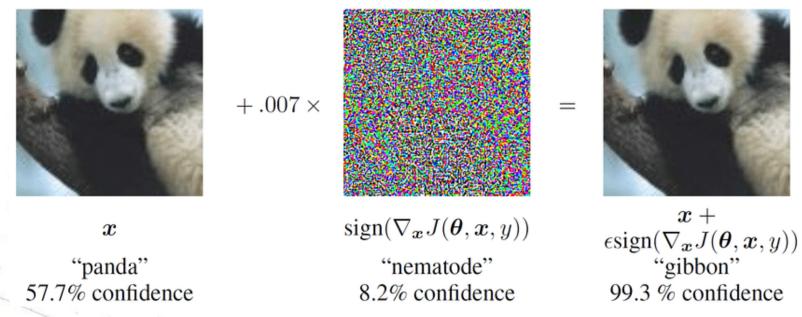
IJCAI-ECAI 2022

# A Unified Framework for Adversarial Attack and Defense in Constrained Feature Space

Thibault Simonetto, Salijona Dyrmishi, Salah Ghamizi, Maxime Cordy, Yves Le Traon University of Luxembourg



### **Evaluating the robustness with adversarial example**



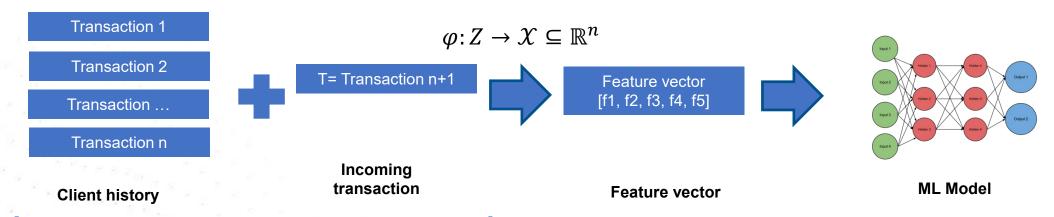
J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," 2014





### **Motivation**

ML model is integrated in a larger software system that takes as input domain objects.



domain object

Domain object Space Z respects some **natural condition** 

Feature space  $\mathcal{X}_{\Omega}$  respects a **set of constraints**  $\Omega$ 



### **Constraints on feature space**

### Related features linked by a constraint:

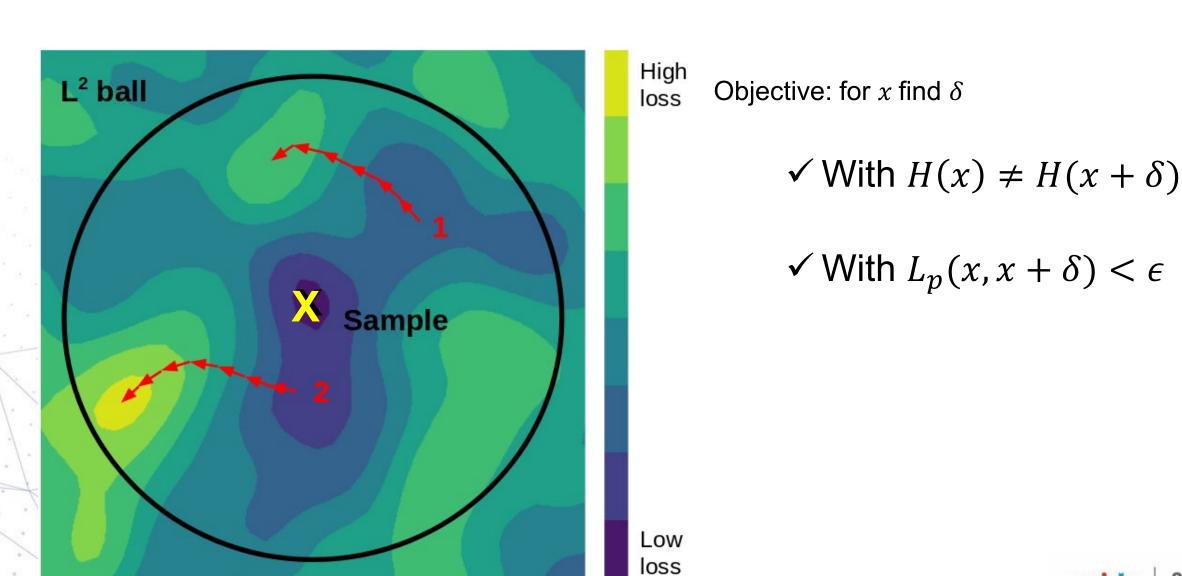
$$open\_acc \leq total\_acc$$

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





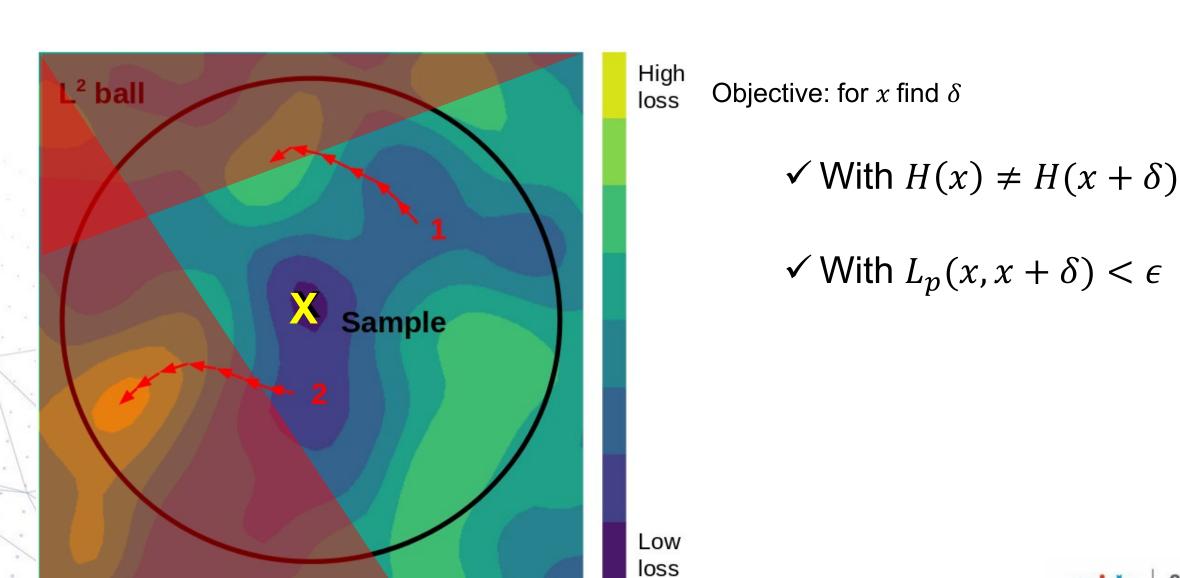
# Existing attacks fails to generate constrained examples







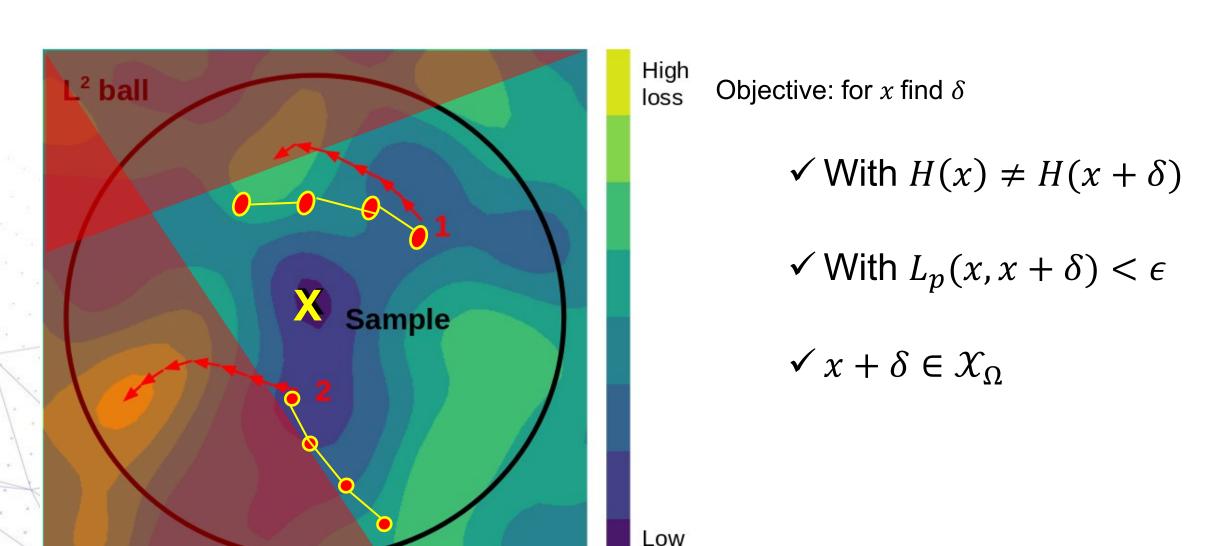
# Existing attacks fails to generate constrained examples







### Existing attacks fails to generate constrained examples



loss



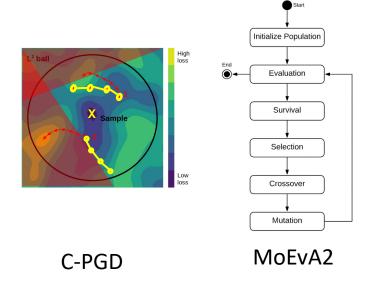
### **Contributions**

### First generic framework for adversarial attacks under domain constraints

$$\omega \coloneqq \omega_1 \wedge \omega_2 \mid \omega_1 \vee \omega_2 \mid \psi_1 \succeq \psi_2 \mid f \in \{\psi_1 \dots \psi_k\}$$
  
$$\psi \coloneqq c \mid f \mid \psi_1 \oplus \psi_2 \mid x_i$$

3	Constraints formulae	Penalty function
* 1	$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
	$\omega_1 ee \omega_2$	$\min(\omega_1,\omega_2)$
i	$\psi \in \Psi = \{\psi_1, \dots \psi_k\}$	$\min(\{\psi_i \in \Psi : \mid \psi - \psi_i \mid \})$
1	$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
. `	$\psi_1 < \psi_2$	$max(0,\psi_1-\psi_2+\tau)$
	$\psi_1=\psi_2$	$\mid \psi_1 - \psi_2 \mid$

Generic constraints language



2 Attacks





### Constraint grammar

$$\omega \coloneqq \omega_1 \wedge \omega_2 \mid \omega_1 \vee \omega_2 \mid \psi_1 \succeq \psi_2 \mid f \in \{\psi_1 \dots \psi_k\}$$
  
$$\psi \coloneqq c \mid f \mid \psi_1 \oplus \psi_2 \mid x_i$$

 $f \in F$  is the value of feature f for a given input x', c is a constant real value,  $\omega, \omega_1, \omega_2$  are constraint formulae,  $g \in \{<, \leq, =, \neq, \geq, >\}$ ,  $\psi, \psi_1, \dots, \psi_k$  are numeric expressions,  $\{<, =, +, -, *, \setminus\}$ , and  $\{x_i\}$  is the value of the i<sup>th</sup> feature of the clean input x

### Constraint grammar

$$\omega \coloneqq \omega_1 \wedge \omega_2 \mid \omega_1 \vee \omega_2 \mid \psi_1 \succeq \psi_2 \mid f \in \{\psi_1 \dots \psi_k\}$$
  
$$\psi \coloneqq c \mid f \mid \psi_1 \oplus \psi_2 \mid x_i$$

 $f \in F$  is the value of feature f for a given input x', c is a constant real value,  $\omega, \omega_1, \omega_2$  are constraint formulae,  $g \in \{<, \leq, =, \neq, \geq, >\}$ ,  $\psi, \psi_1, ..., \psi_k$  are numeric expressions,  $\{<, +, -, *, \setminus\}$ , and  $\{x_i\}$  is the value of the i<sup>th</sup> feature of the clean input x



### Constraint grammar

$$\omega \coloneqq \omega_1 \wedge \omega_2 \mid \omega_1 \vee \omega_2 \mid \psi_1 \succeq \psi_2 \mid f \in \{\psi_1 \dots \psi_k\}$$
  
$$\psi \coloneqq c \mid f \mid \psi_1 \oplus \psi_2 \mid x_i$$

 $f \in F$  is the value of feature f for a given input x', c is a constant real value,  $\omega, \omega_1, \omega_2$  are constraint formulae,  $g \in \{<, \leq, =, \neq, \geq, >\}$ ,  $\psi, \psi_1, \dots, \psi_k$  are numeric expressions,  $\{ \in \{+, -, *, \setminus \} \}$ , and  $\{x_i\}$  is the value of the i<sup>th</sup> feature of the clean input  $\{x\}$ 

### Mapping to penalty functions

Constraints formulae	Penalty function
$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
$\omega_1 \vee \omega_2$	$\min(\omega_1,\omega_2)$
$\psi \in \Psi = \{\psi_1, \dots \psi_k\}$	$\min(\{\psi_i \in \Psi : \mid \psi - \psi_i \mid \})$
$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
$\psi_1 < \psi_2$	$max(0,\psi_1-\psi_2+\tau)$
$\psi_1=\psi_2$	$\mid \psi_1 - \psi_2 \mid$

Table 1: From constraint formulae to penalty functions.  $\tau$  is an infinitesimal value.

Constraint is satisfied if and only if  $g(\omega, x) = 0$ 



### Constraint grammar

$$\omega \coloneqq \omega_1 \wedge \omega_2 \mid \omega_1 \vee \omega_2 \mid \psi_1 \succeq \psi_2 \mid f \in \{\psi_1 \dots \psi_k\}$$
  
$$\psi \coloneqq c \mid f \mid \psi_1 \oplus \psi_2 \mid x_i$$

 $f \in F$  is the value of feature f for a given input x', c is a constant real value,  $\omega, \omega_1, \omega_2$  are constraint formulae,  $g \in \{<, \leq, =, \neq, \geq, >\}$ ,  $\psi, \psi_1, ..., \psi_k$  are numeric expressions,  $\{<, +, -, *, \setminus\}$ , and  $\{x_i\}$  is the value of the i<sup>th</sup> feature of the clean input x

### Mapping to penalty functions

Constraints formulae	Penalty function
$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
$\omega_1 \vee \omega_2$	$\min(\omega_1,\omega_2)$
$\psi \in \Psi = \{\psi_1, \dots \psi_k\}$	$\min(\{\psi_i \in \Psi :  \psi - \psi_i \})$
$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
$\psi_1 < \psi_2$	$max(0,\psi_1-\psi_2+\tau)$
$\psi_1 = \psi_2$	$\mid \psi_1 - \psi_2 \mid$

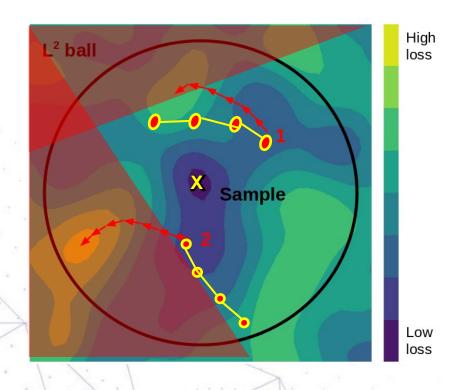
Table 1: From constraint formulae to penalty functions.  $\tau$  is an infinitesimal value.

Constraint is satisfied if and only if  $g(\omega, x) = 0$ 

# Sufficient expressiveness to instantiate constraints in different domains



# **Approach 1: C-PGD: Gradient evaluation of the constraints**



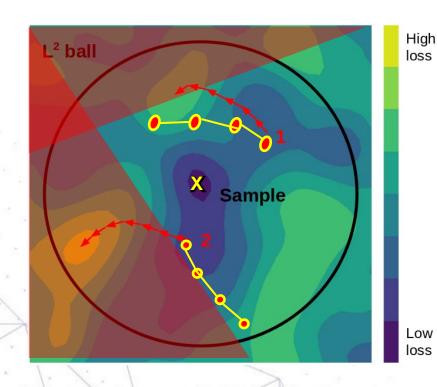
### **Projected Gradient Descent (PGD)**

$$x^{t+1} = x^t + \alpha \operatorname{sign}(\nabla_{x^t} loss(h(x^t), y))$$





### **Approach 1: C-PGD: Gradient evaluation of the constraints**



**Projected Gradient Descent (PGD)** 

$$x^{t+1} = x^t + \alpha \operatorname{sign}(\nabla_{x^t} loss(h(x^t), y))$$

**Constraints regularization** 

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

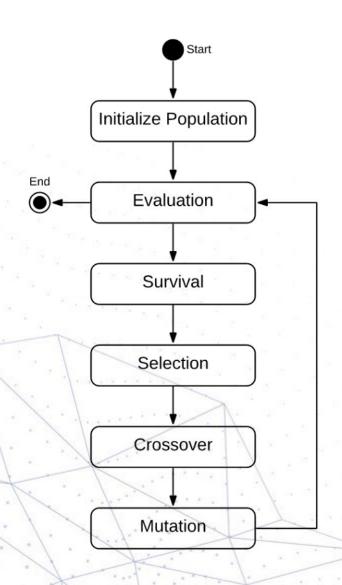
**Constrained Projected Gradient Descent (C-PGD)** 

$$x^{t+1} = x^t + \alpha \, sign\left(\nabla_{x^t} loss(h(x^t), y) - \nabla_{x^t} penalty(x^t)\right)$$





### **Approach 2: MoEvA2: Evolutionary approach**



Multi-objective genetic algorithm (NSGA-III)

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

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

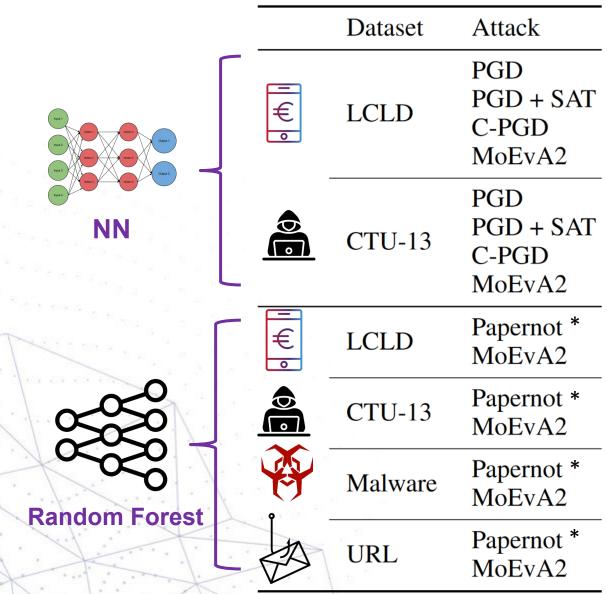
minimise 
$$g_3(x) \equiv \sum_{\omega_i \in \Omega} penalty(x, \omega_i)$$



# How effective are our approaches at generating adversarial examples?



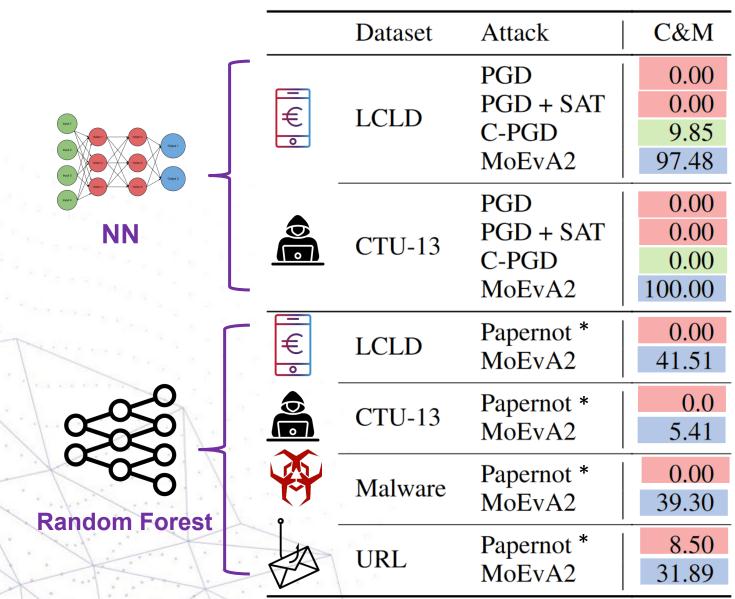
### How effective are our approaches at generating adversarial examples?







### How effective are our approaches at generating adversarial examples?



Attacks unaware of domain constraints often fail.

C-PGD is effective on a single dataset.

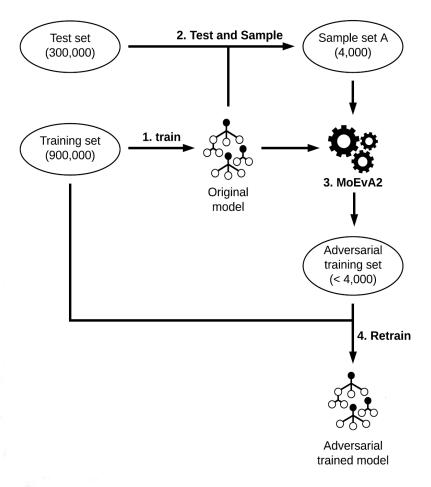
MoEvA2 is always successful.





<sup>\*</sup> Extended to random forest





**Adversarial retraining** 





We hypothesize that augmenting  $\Omega$  with a set of engineered constraints can robustify a model.



We hypothesize that augmenting  $\Omega$  with a set of engineered constraints can robustify a model.

We engineer a new feature

$$f_e = f_1 \oplus f_2$$



We hypothesize that augmenting  $\Omega$  with a set of engineered constraints can robustify a model.

We engineer a new feature

$$f_e = f_1 \oplus f_2$$

We have the **new constraint** 

$$\omega_e \vDash (f_e = f_1 \oplus f_2)$$





### How effective are defense techniques?





Defense	Attack	LCLD	CTU-13
None	C-PGD	9.85	0.00
None	MoEvA2	97.48	
C-PGD Adv. retraining * C-PGD Adv. retraining *	C-PGD	8.78	NA
	MoEvA2	94.90	NA
MoEvA2 Adv. retraining * MoEvA2 Adv. retraining *	C-PGD	2.70	NA
	MoEvA2	85.20	0.8
Constraints augment. Constraints augment.	C-PGD MoEvA2	0.00 80.43	NA 0.00
MoEvA2 Adv. retrain. † Combined defenses †	MoEvA2	82.00	NA
	MoEvA2	77.43	NA

Table 3: Success rate of C-PGD and MoEvA2 after adversarial retraining and constraint augmentation (on neural networks). For a fair comparison, the model denoted by the same symbols (\* or †) are trained with the same number of adversarial examples, generated from the same original samples.

Constraint augmentation is an effective alternative defense to adversarial retraining.

Constraint augmentation and adversarial retraining have complementary effects.



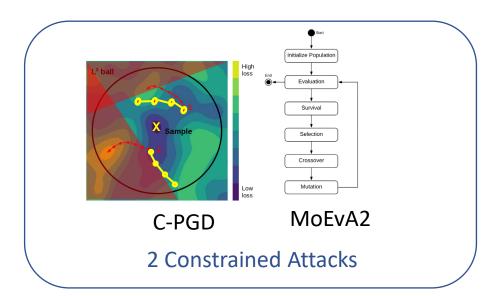


Constraints formulae	Penalty function
$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
$\omega_1 \vee \omega_2$	$\min(\omega_1,\omega_2)$
$\psi \in \Psi = \{\psi_1, \dots \psi_k\}$	$\min(\{\psi_i \in \Psi :  \psi - \psi_i \})$
$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
$\psi_1 < \psi_2$	$max(0,\psi_1-\psi_2+\tau)$
$\psi_1=\psi_2$	$\mid \psi_1 - \psi_2 \mid$

Generic constraints language

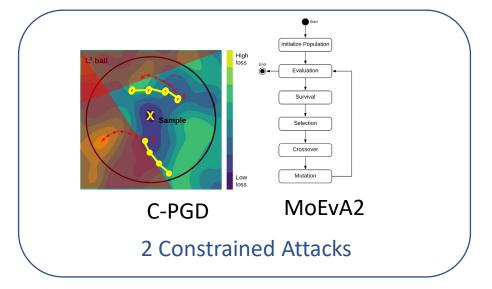
Constraints formulae	Penalty function
$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
$\omega_1 \vee \omega_2$	$\min(\omega_1,\omega_2)$
$\psi \in \Psi = \{\psi_1, \dots \psi_k\}$	$\min(\{\psi_i \in \Psi :  \psi - \psi_i \})$
$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
$\psi_1 < \psi_2$	$max(0,\psi_1-\psi_2+\tau)$
$\psi_1 = \psi_2$	$\mid \psi_1 - \psi_2 \mid$

Generic constraints language



Constraints formulae	Penalty function
$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
$\omega_1 \vee \omega_2$	$\min(\omega_1,\omega_2)$
$\psi \in \Psi = \{\psi_1, \dots \psi_k\}$	$\min(\{\psi_i \in \Psi :   \psi - \psi_i   \})$
$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
$\psi_1 < \psi_2$	$max(0,\psi_1-\psi_2+\tau)$
$\psi_1 = \psi_2$	$\mid \psi_1 - \psi_2 \mid$

Generic constraints language











Attacks unaware of domain constraints often fail.

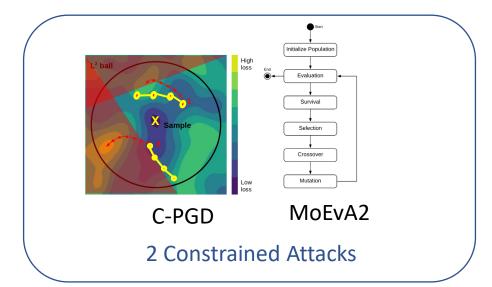
MoEvA2 is always successful.



Constraints formulae	Penalty function
$\omega_1 \wedge \omega_2$	$\omega_1 + \omega_2$
$\omega_1 \vee \omega_2$	$\min(\omega_1,\omega_2)$
$\psi \in \Psi = \{\psi_1, \dots \psi_k\}$	$\min(\{\psi_i \in \Psi :  \psi - \psi_i \})$
$\psi_1 \le \psi_2$	$max(0,\psi_1-\psi_2)$
$y/y_1 < y/y_2$	$max(0, 1/2 - 1/2 + \tau)$

 $|\psi_1 - \psi_2|$ 

Generic constraints language





 $\psi_1 = \psi_2$ 







Attacks unaware of domain constraints often fail.

MoEvA2 is always successful.

$$\omega_e \vDash (f_e = f_1 \oplus f_2)$$

**Constrained Augmentation** 

New defense method as effective as adversarial retraining





Constraints formulae Penalty function

 $\begin{array}{lll} \omega_{1} \wedge \omega_{2} & \omega_{1} + \omega_{2} \\ \omega_{1} \vee \omega_{2} & \min(\omega_{1}, \omega_{2}) \\ \psi \in \Psi = \{\psi_{1}, \dots \psi_{k}\} & \min(\{\psi_{i} \in \Psi \\ \psi_{1} \leq \psi_{2} & \max(0, \psi_{1} - \psi_{1}) \\ \psi_{1} = \psi_{2} & |\psi_{1} - \psi_{2}| \end{array}$ 

Generic constraints lang







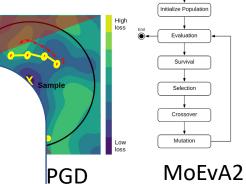
Attacks unaware of do constraints often fa

MoEvA2 is always successful.

**Checkout our framework** 



Stand 149 Row 5



Constrained Attacks

$$\vDash (f_e = f_1 \oplus f_2)$$

trained Augmentation

New defense method as effective as adversarial retraining



