

Difference-of-Convex Algorithm applied to adversarial robustness verification

Ismaila Seck ^{1,2,3} Gaelle Loosli ^{3,4} Stephane Canu ^{2,3} Yi-Shuai Niu ⁵

¹Normandie Univ, INSA Rouen ²UNIROUEN, UNIHAVRE, LITIS ³UCA, LIMOS

⁴PobRun ⁵School of Mathematical Sciences, Shanghai Jiao Tong University

October 2, 2019

Adversarial examples



Figure 1: Illustration of the use of adversarial examples.

- x, y : original image and its class
- x' : adversarial image we are looking for
- $f_k(.)$: the *k*-th output of the network

$$\begin{cases} \min & \|\mathbf{x} - \mathbf{x}'\| \\ \text{s.t.} & \operatorname*{argmax}_{k=1,\dots,c} f_k(\mathbf{x}') \neq y, \\ & \mathbf{x}' \in [0,1]^d. \end{cases}$$
(1)

- x, y : original image and its class
- \mathbf{x}' : adversarial image we are looking for
- $f_k(.)$: the *k*-th output of the network

$$\begin{cases} \min & \|\mathbf{x} - \mathbf{x}'\| \\ \text{s.t.} & \underset{\substack{k=1,\dots,c \\ \mathbf{x}' \in [0,1]^d}. \end{cases} \end{cases}$$

(1)

(1)
$$\iff \begin{cases} \min \|\mathbf{x} - \mathbf{x}'\| \\ \text{s.t.} & m \le f_k(\mathbf{x}') + (1 - a_k)M_m, \quad k = 1, \dots, c \\ & m \ge f_k(\mathbf{x}'), & k = 1, \dots, c \\ & \sum_{k=1}^{c} a_k = 1, \\ & a_y = 0, \\ & m \in \mathbb{R}, \\ & \mathbf{a} \in \{0, 1\}^c, \\ & \mathbf{x}' \in [0, 1]^d. \end{cases}$$
(2)

Using DC to get rid of the binary variables

$$\begin{array}{ll} \min & \|\mathbf{x} - \mathbf{x}'\| + \sum_{k=1}^{c} a_{k}(1 - a_{k}) \\ \text{s.t.} & m \leq f_{k}(\mathbf{x}') + (1 - a_{k})M_{m}, \quad k = 1, \dots, c \\ & m \geq f_{k}(\mathbf{x}'), \qquad \qquad k = 1, \dots, c \\ & \sum_{k=1}^{c} a_{k} = 1, \\ & a_{y} = 0, \\ & m \in \mathbb{R}, \\ & a \in [0, 1]^{c}, \\ & \mathbf{x}' \in [0, 1]^{d}. \end{array}$$

(3)

Thanks for your attention !