

Verification of Binarized Deep Neural Networks

Mădălina Eraşcu

Institute e-Austria Timișoara, Romania

West University of Timișoara, Romania

`madalina.erascu@e-uvvt.ro`

Based on the papers: (1) *Verifying Properties of Binarized Deep Neural Networks*
– N. Narodytska et al, AAAI-18, (2) *Formal Analysis of Deep Binarized Neural
Networks* – N. Narodytska, IJCAI-18

November 27th, 2019



Outline

Motivation

Preliminaries

- Properties of neural networks
- Binarized Neural Networks (BNNs)

Encoding the BNNs

- Mixed Integer Linear Program (MILP) Encoding
- Integer Linear Programming (ILP) Encoding
- SAT Encoding
- Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.
- ▶ **Question:** Can we trust the decisions that neural networks make?

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.
- ▶ **Question:** Can we trust the decisions that neural networks make?

Define the properties that the neural networks should have and verify if they hold for the network.

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.
- ▶ **Question:** Can we trust the decisions that neural networks make?

Define the properties that the neural networks should have and verify if they hold for the network.

Contributions of the papers:

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.
- ▶ **Question:** Can we trust the decisions that neural networks make?

Define the properties that the neural networks should have and verify if they hold for the network.

Contributions of the papers:

1. Define interesting properties of the neural networks to analyze: *robustness*, *equivalence*, *invertibility*.

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.
- ▶ **Question:** Can we trust the decisions that neural networks make?

Define the properties that the neural networks should have and verify if they hold for the network.

Contributions of the papers:

1. Define interesting properties of the neural networks to analyze: *robustness*, *equivalence*, *invertibility*.
2. Analyze Binarized Neural Networks (BNNs) using MILP, ILP and SAT methods.

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.
- ▶ **Question:** Can we trust the decisions that neural networks make?

Define the properties that the neural networks should have and verify if they hold for the network.

Contributions of the papers:

1. Define interesting properties of the neural networks to analyze: *robustness, equivalence, invertibility*.
2. Analyze Binarized Neural Networks (BNNs) using MILP, ILP and SAT methods.
3. Exploit the structure of the neural network for speeding-up the timings of analysis using the *counterexample-guided search procedure*.

Motivation

- ▶ Deep learning is used on safety-critical systems, however, the main criticism and concern is on the lack of understanding the decision process behind the networks.
- ▶ **Question:** Can we trust the decisions that neural networks make?

Define the properties that the neural networks should have and verify if they hold for the network.

Contributions of the papers:

1. Define interesting properties of the neural networks to analyze: *robustness*, *equivalence*, *invertibility*.
2. Analyze Binarized Neural Networks (BNNs) using MILP, ILP and SAT methods.
3. Exploit the structure of the neural network for speeding-up the timings of analysis using the *counterexample-guided search procedure*.
4. The methods were tested for image classification tasks for medium-size BNNs.

Why Verification/Analysis of Deep Neural Networks (DNNs) is important?

- ▶ Is there a way to analyze deep neural networks?

Why Verification/Analysis of Deep Neural Networks (DNNs) is important?

- ▶ Is there a way to analyze deep neural networks?
- ▶ Can we explain their decisions?

Why Verification/Analysis of Deep Neural Networks (DNNs) is important?

- ▶ Is there a way to analyze deep neural networks?
- ▶ Can we explain their decisions?
- ▶ How robust are these networks to perturbations of inputs?

Why Verification/Analysis of Deep Neural Networks (DNNs) is important?

- ▶ Is there a way to analyze deep neural networks?
- ▶ Can we explain their decisions?
- ▶ How robust are these networks to perturbations of inputs?
- ▶ How critical is the choice of one architecture over an other?

Why Verification/Analysis of Deep Neural Networks (DNNs) is important?

- ▶ Is there a way to analyze deep neural networks?
- ▶ Can we explain their decisions?
- ▶ How robust are these networks to perturbations of inputs?
- ▶ How critical is the choice of one architecture over an other?

Verification problem: encode the (exact representation of the) network and the property we aim to verify as a formal statement, using, e.g. ILP, SMT or SAT.

Contents

Motivation

Preliminaries

Properties of neural networks
Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding
Integer Linear Programming (ILP) Encoding
SAT Encoding
Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Notations

- ▶ $[m] = \{1, 2, \dots, m\}$
- ▶ $\mathbf{v} = (v_1, \dots, v_m)$, $\mathbf{v} \in \mathbb{R}$
- ▶ $\|\mathbf{v}\|_p$, $p \geq 1$ is the L_p -norm \mathbf{v}

Notations

- ▶ $[m] = \{1, 2, \dots, m\}$
- ▶ $\mathbf{v} = (v_1, \dots, v_m)$, $\mathbf{v} \in \mathbb{R}$
- ▶ $\|\mathbf{v}\|_p$, $p \geq 1$ is the L_p -norm \mathbf{v}
- ▶ Formula A is satisfiable/unsatisfiable ...

Notations

- ▶ $[m] = \{1, 2, \dots, m\}$
- ▶ $\mathbf{v} = (v_1, \dots, v_m)$, $\mathbf{v} \in \mathbb{R}$
- ▶ $\|\mathbf{v}\|_p$, $p \geq 1$ is the L_p -norm \mathbf{v}
- ▶ Formula A is satisfiable/unsatisfiable ...
- ▶ *Image classification problem*: we are given (1) a set of training images drawn from an unknown distribution ν over $X = \mathbb{Z}^n$, where n is the size of individual images, (2) a label for each image $L : \mathbb{Z}^n \rightarrow [s]$.

Notations

- ▶ $[m] = \{1, 2, \dots, m\}$
- ▶ $\mathbf{v} = (v_1, \dots, v_m)$, $\mathbf{v} \in \mathbb{R}$
- ▶ $\|\mathbf{v}\|_p$, $p \geq 1$ is the L_p -norm \mathbf{v}
- ▶ Formula A is satisfiable/unsatisfiable ...
- ▶ **Image classification problem:** we are given (1) a set of training images drawn from an unknown distribution ν over $X = \mathbb{Z}^n$, where n is the size of individual images, (2) a label for each image $L : \mathbb{Z}^n \rightarrow [s]$.
 - ▶ *Training:* given a labeled training set, learn a neural network classifier that can be used as inference engine.
 - ▶ During the inference the network is *fixed*.

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Properties of neural networks

Let F be a general feedforward neural, $F(\mathbf{x})$ be the output of F on input \mathbf{x} and $l_{\mathbf{x}} = L(\mathbf{x})$ be the ground truth label of \mathbf{x} .

Properties of neural networks

Let F be a general feedforward neural, $F(\mathbf{x})$ be the output of F on input \mathbf{x} and $l_{\mathbf{x}} = L(\mathbf{x})$ be the ground truth label of \mathbf{x} .

Properties:

1. **Robustness:** small perturbations on inputs do not affect the output.
 - ▶ **global robustness:** for any valid input, there is no small perturbation that can change the decision of the network on this input.

Definition (Global Robustness)

A feedforward neural network F is globally-robust if for any input \mathbf{x} , $\mathbf{x} \in X$ and τ , $\|\tau\|_{\infty} \leq \epsilon$ we have that $F(\mathbf{x} + \tau) = l_{\mathbf{x}}$.

- ▶ **local robustness:** is defined for a single input \mathbf{x} .

Definition (Local Robustness)

A feedforward neural network F is locally-robust for an input \mathbf{x} , $\mathbf{x} \in X$, if there does not exist τ , $\|\tau\|_{\infty} \leq \epsilon$ such that $F(\mathbf{x} + \tau) \neq l_{\mathbf{x}}$.

- 2.
- 3.

Properties of neural networks

Let F be a general feedforward neural, $F(\mathbf{x})$ be the output of F on input \mathbf{x} and $l_{\mathbf{x}} = L(\mathbf{x})$ be the ground truth label of \mathbf{x} .

Properties:

1. **Robustness:** small perturbations on inputs do not affect the output.
2. **Invertibility:** explore a set of inputs that map to a given output (example: what the inputs of the network are, if exist, that map to a given output).

Definition (Local Invertibility)

A feedforward neural network F is locally invertible for an output \mathbf{s} if there exists \mathbf{x} , $\mathbf{x} \in C(X)$, such that $F(\mathbf{x}) = \mathbf{s}$, where $C(X)$ denotes the constrained domain of inputs.

Related problem: how to enumerate multiple, preferably diverse by some measure, inputs of the network that map to a given output.

3.

Properties of neural networks

Let F be a general feedforward neural, $F(\mathbf{x})$ be the output of F on input \mathbf{x} and $l_{\mathbf{x}} = L(\mathbf{x})$ be the ground truth label of \mathbf{x} .

Properties:

- 1.
- 2.
3. **Network equivalence:** two networks F_1 and F_2 are equivalent if they generate same outputs on all inputs drawn from the domain X .

Definition (Network Equivalence)

Two feedforward neural networks F_1 and F_2 are equivalent if for all $\mathbf{x} \in X$, $F_1(\mathbf{x}) = F_2(\mathbf{x})$.

Application: network alteration.

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

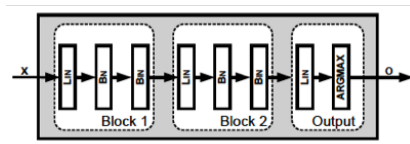
Experimental Results

Binarized Neural Networks (BNNs)

Definition (Binarized Neural Network)

A binarized neural network $BNN : \{-1, 1\}^n \rightarrow [s]$ is a feedforward network that is composed of d blocks, $BLK_1, \dots, BLK_{d-1}, O$. Formally, given an input \mathbf{x} , $BNN(\mathbf{x}) = O(BLK_{d_1}, \dots, (BLK_1(\mathbf{x})))$

Schematic view of a binarized neural network



Structure of internal and outputs blocks, which stacked together form a BNN.

A_k and b_k – parameters of the LIN layer; $\alpha_{k_i}, \gamma_{k_i}, \mu_{k_i}, \sigma_{k_i}$ – parameters of the BN layer. μ and σ correspond to mean and standard deviation computed in the training phase. The BIN layer is parameter free.

Structure of k th Internal block, $BLK_k : \{-1, 1\}^{n_k} \rightarrow \{-1, 1\}^{n_{k+1}}$ on input $x_k \in \{-1, 1\}^{n_k}$	
LIN	$y = A_k x_k + b_k$, where $A_k \in \{-1, 1\}^{n_{k+1} \times n_k}$ and $b_k \in \mathbb{R}^{n_{k+1}}$
BN	$z_i = \alpha_{k_i} \left(\frac{y_i - \mu_{k_i}}{\sigma_{k_i}} \right) + \gamma_{k_i}$, where $y = (y_1, \dots, y_{n_{k+1}})$, and $\alpha_{k_i}, \gamma_{k_i}, \mu_{k_i}, \sigma_{k_i} \in \mathbb{R}$. Assume $\sigma_{k_i} > 0$.
BIN	$x_{k+1} = \text{sign}(z)$ where $z = (z_1, \dots, z_{n_{k+1}}) \in \mathbb{R}^{n_{k+1}}$ and $x_{k+1} \in \{-1, 1\}^{n_{k+1}}$
Structure of Output Block, $O : \{-1, 1\}^{n_d} \rightarrow [s]$ on input $x_d \in \{-1, 1\}^{n_d}$	
LIN	$w = A_d x_d + b_d$, where $A_d \in \{-1, 1\}^{s \times n_d}$ and $b_d \in \mathbb{R}^s$
ARGMAX	$o = \text{argmax}(w)$, where $o \in [s]$

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Encodings of BNNs

... into Boolean formulae.

- ▶ The encoding of the BNN is a conjunction of encodings of its blocks.
- ▶ $BINBLK_k(\mathbf{x}_k, \mathbf{x}_{k+1})$ a Boolean function that encodes the k th block (BLK_k) with an input \mathbf{x}_k and an output \mathbf{x}_{k+1} .
- ▶ $BINO(\mathbf{x}_d, \mathbf{o})$ be a Boolean function that encodes O that takes an input \mathbf{x}_d and outputs \mathbf{o} .
- ▶ The entire *BNN* on input \mathbf{x} can be encoded as a Boolean formula, with \mathbf{x}_1 (first layer) = \mathbf{x} (input):

$$\left(\bigwedge_{k=1}^{d-1} BINBLK_k(x_k, x_{k+1}) \right) \wedge BINO(x_d, \mathbf{o})$$

- ▶ Encodings: MILP \rightsquigarrow ILP \rightsquigarrow SAT.

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Mixed Integer Linear Program (MILP) Encoding

- ▶ *Encoding of BLK_k* : encode each layer in BLK_k to MILP separately. Let a_i be the i -th row of the matrix A_k , $x_k \in \{-1, 1\}^{n_k}$ denote the input to BLK_k .
- ▶ *Linear Transformation*. Transformation for fully connected layer (suitable for convolutions also as they are linear operations). We have:

$$y_i = \langle \mathbf{a}_i, \mathbf{x}_k \rangle + b_i, \quad i = \overline{1, n_{k+1}} \quad (1)$$

where $\mathbf{y} = (y_1, \dots, y_{n_{k+1}}) \in \mathbb{R}^{n_{k+1}}$.

- ▶ *Batch Normalization*: takes the output of the linear layer as an input. By definition, we have:

$$\begin{aligned} z_i &= \alpha_{k_i} \left(\frac{y_i - \mu_{k_i}}{\sigma_{k_i}} \right) + \gamma_{k_i}, \quad i = \overline{1, n_{k+1}} \\ \sigma_{k_i} z_i &= \alpha_{k_i} y_i - \alpha_{k_i} \mu_{k_i} + \sigma_{k_i} \gamma_{k_i} \end{aligned} \quad (2)$$

- ▶ *Binarization*. For the BIN operation, which implements a sign function, we need to deal with conditional constraints.

$$z_i \geq 0 \Rightarrow v_i = 1 \quad (3)$$

$$z_i < 0 \Rightarrow v_i = -1 \quad i = \overline{1, n_{k+1}} \quad (4)$$

Mixed Integer Linear Program Encoding (cont'd)

- ▶ *Encoding of O* : $w_i = \langle \mathbf{a}_i, \mathbf{x}_d \rangle + b_i$, $i = 1, \dots, s$, where \mathbf{a}_i represents the i -th column in A_d and $\mathbf{w} = (w_1, \dots, w_s)$. To encode ARGMAX, an ordering relation between w_i 's must be imposed:

$$\begin{aligned} w_i \geq w_j &\iff d_{ij} = 1 && d_{ij} \text{ newly introduced vars} \\ \sum_{j=1}^s d_{ij} = s &\implies o = i && i, j = \overline{1, s} \end{aligned} \quad (5)$$

Mixed Integer Linear Program Encoding (cont'd)

- *Encoding of O* : $w_i = \langle \mathbf{a}_i, \mathbf{x}_d \rangle + b_i$, $i = 1, \dots, s$, where \mathbf{a}_i represents the i -th column in A_d and $\mathbf{w} = (w_1, \dots, w_s)$. To encode ARGMAX, an ordering relation between w_i 's must be imposed:

$$\begin{aligned} w_i \geq w_j &\iff d_{ij} = 1 && d_{ij} \text{ newly introduced vars} \\ \sum_{j=1}^s d_{ij} = s &\implies o = i && i, j = \overline{1, s} \end{aligned} \quad (5)$$

Example

Consider an internal block with two inputs and one output. Suppose we have the following parameters: $A_k = [1, -1]$, $b_k = [-0.5]$, $\alpha_k = [0.12]$, $\mu_k = [-0.1]$, $\sigma_k = [2]$, $\delta_k = [0.1]$.

1. apply the linear transformation: $y_1 = x_{k1} - x_{k2} - 0.5$
2. apply batch normalization: $2z_1 = 0.12y_1 + (-0.12) * (-0.1) + 2 * 0.1$
3. apply binarization: $z_1 \geq 0 \implies v_1 = 1 \wedge z_1 < 0 \implies v_1 = -1$. ($x_{k+1} = v_1$).

Mixed Integer Linear Program Encoding (cont'd)

- ▶ *Encoding of O* : $w_i = \langle \mathbf{a}_i, \mathbf{x}_d \rangle + b_i$, $i = 1, \dots, s$, where \mathbf{a}_i represents the i -th column in A_d and $\mathbf{w} = (w_1, \dots, w_s)$. To encode ARGMAX, an ordering relation between w_i 's must be imposed:

$$\begin{aligned} w_i \geq w_j &\iff d_{ij} = 1 && d_{ij} \text{ newly introduced vars} \\ \sum_{j=1}^s d_{ij} = s &\implies o = i && i, j = \overline{1, s} \end{aligned} \quad (5)$$

Example

Consider an output block with two inputs and two outputs. We have the following parameters for this block $A_d = [1, -1; -1, 1]$ and $b = [-0.5, 0.2]$.

1. encoding of the linear transformation

$$w_1 = x_{d_1} - x_{d_2} - 0.5 \quad w_2 = -x_{d_1} + x_{d_2} + 0.2$$

2. Two outputs \rightsquigarrow 4 Boolean variables d_{ij} , $i, j = 1, 2$; $d_{11} = d_{22} = 1 \rightsquigarrow$ consider only non-diagonal variables.

$$w_1 \geq w_2 \iff d_{12} = 1 \quad \wedge \quad w_2 < w_1 \iff d_{21} = 1$$

3. we compute the output o of the neural network as:

$$d_{11} + d_{12} = 2 \implies o = 1 \quad \wedge \quad d_{21} + d_{22} = 2 \implies o = 2$$

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Integer Linear Programming (ILP) Encoding

ILP encoding is smaller than the MILP one.

- ▶ *Encoding of BLK_k* : z and y are functional variables of \mathbf{x}_k . We can substitute them in (3) and (4) based on (1) and (2) respectively:

$$\frac{\alpha_{k_i}}{\sigma_{k_i}} (\langle \mathbf{a}_i, \mathbf{x}_k \rangle + b_i) - \frac{\alpha_{k_i}}{\sigma_{k_i}} \mu_{k_i} + \gamma_{k_i} \implies v_i = 1$$

- ▶ *Linear and batch normalization:*

- ▶ Case $\alpha_{k_i} > 0$, we have:

$$\langle \mathbf{a}_i, \mathbf{x}_k \rangle \geq -\frac{\sigma_{k_i}}{\alpha_{k_i}} \gamma_{k_i} + \mu_{k_i} - b_i \implies x'_i = 1 \text{ (see p. 6618 right column, bottom)}$$

Consider $C_i = \left[-\frac{\sigma_{k_i}}{\alpha_{k_i}} \gamma_{k_i} + \mu_{k_i} - b_i \right]$. We encode (3) and (4):

$$\langle \mathbf{a}_i, \mathbf{x}_k \rangle \geq C_i \implies v_i = 1$$

$$\langle \mathbf{a}_i, \mathbf{x}_k \rangle < C_i \implies v_i = -1 \quad i = \overline{1, n_{k+1}}$$

- ▶ Case $\alpha_{k_i} < 0$. Same as above but $C_i = \left[-\frac{\sigma_{k_i}}{\alpha_{k_i}} \gamma_{k_i} + \mu_{k_i} - b_i \right]$
- ▶ Case $\alpha_{k_i} = 0$, we have: $\gamma_{k_i} \implies v_i = 1$

Integer Linear Programming (ILP) Encoding (cont'd)

- ▶ *Encoding of O* : introduce the Boolean variables d_{ij} avoiding the intermediate variables w_i

$$\langle \mathbf{a}_i, \mathbf{x}_d \rangle + b_i \geq \langle \mathbf{a}_j, \mathbf{x}_d \rangle + b_j \iff d_{ij} = 1, \quad i, j = \overline{1, s}$$

$$\iff$$

$$\langle \mathbf{a}_i - \mathbf{a}_j, \mathbf{x}_d \rangle \geq \lceil b_j - b_i \rceil \iff d_{ij} = 1$$

where \mathbf{a}_i and \mathbf{a}_j denote the i th and j th rows in the matrix A_d .

Further, constraints (5) can be used.

Integer Linear Programming (ILP) Encoding (cont'd)

- ▶ *Encoding of O* : introduce the Boolean variables d_{ij} avoiding the intermediate variables w_i

$$\langle \mathbf{a}_i, \mathbf{x}_d \rangle + b_i \geq \langle \mathbf{a}_j, \mathbf{x}_d \rangle + b_j \iff d_{ij} = 1, \quad i, j = \overline{1, s}$$

$$\iff$$

$$\langle \mathbf{a}_i - \mathbf{a}_j, \mathbf{x}_d \rangle \geq \lceil b_j - b_i \rceil \iff d_{ij} = 1$$

where \mathbf{a}_i and \mathbf{a}_j denote the i th and j th rows in the matrix A_d .

Further, constraints (5) can be used.

Example

Recall the internal block with two inputs and one output: $A_k = [1, -1]$, $b_k = [-0.5]$, $\alpha_k = [0.12]$, $\mu_k = [-0.1]$, $\sigma_k = [2]$, $\delta_k = [0.1]$. We have:

$$x_{k_1} - x_{k_2} \geq \left\lceil \frac{-2}{0.1} * 0.1 - 0.1 - (-0.5) \right\rceil = -1 \Rightarrow v_1 = 1$$

$$x_{k_1} - x_{k_2} < \left\lceil \frac{-2}{0.1} * 0.1 - 0.1 - (-0.5) \right\rceil = -1 \Rightarrow v_1 = -1$$

Integer Linear Programming (ILP) Encoding (cont'd)

- ▶ *Encoding of O* : introduce the Boolean variables d_{ij} avoiding the intermediate variables w_i

$$\langle \mathbf{a}_i, \mathbf{x}_d \rangle + b_i \geq \langle \mathbf{a}_j, \mathbf{x}_d \rangle + b_j \iff d_{ij} = 1, \quad i, j = \overline{1, s}$$

$$\iff$$

$$\langle \mathbf{a}_i - \mathbf{a}_j, \mathbf{x}_d \rangle \geq \lceil b_j - b_i \rceil \iff d_{ij} = 1$$

where \mathbf{a}_i and \mathbf{a}_j denote the i th and j th rows in the matrix A_d .

Further, constraints (5) can be used.

Example

Recall the output block with two inputs and two outputs: $A_d = [1, -1; -1, 1]$ and $b = [-0.5, 0.2]$. We have

$$x_{d_1} - x_{d_2} - 0.5 \geq -x_{d_1} + x_{d_2} + 0.2 \iff d_{12} = 1 \text{ equiv. to } x_{d_1} - x_{d_2} \geq \left\lceil \frac{0.7}{2} \right\rceil \iff d_{12} = 1$$

$$x_{d_2} - x_{d_1} + 0.2 \geq -x_{d_2} + x_{d_1} - 0.5 \iff d_{21} = 1 \text{ equiv. to } x_{d_1} - x_{d_2} \leq \left\lceil \frac{0.7}{2} \right\rceil \iff d_{21} = 1$$

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

SAT Encoding

Naive approach: encode the *BLK* and *O* blocks into CNF is to directly translate their ILP encoding defined above into SAT.

SAT Encoding

Naive approach: encode the *BLK* and *O* blocks into CNF is to directly translate their ILP encoding defined above into SAT.

Drawback: inefficient – the resulting encoding will be very large.

SAT Encoding

Naive approach: encode the *BLK* and *O* blocks into CNF is to directly translate their ILP encoding defined above into SAT.

Drawback: inefficient – the resulting encoding will be very large.

Solution: exploit the properties of BNNs.

SAT Encoding

Naive approach: encode the *BLK* and *O* blocks into CNF is to directly translate their ILP encoding defined above into SAT.

Drawback: inefficient – the resulting encoding will be very large.

Solution: exploit the properties of BNNs.

Sequential counters for encoding cardinality constraints

SAT Encoding

Naive approach: encode the *BLK* and *O* blocks into CNF is to directly translate their ILP encoding defined above into SAT.

Drawback: inefficient – the resulting encoding will be very large.

Solution: exploit the properties of BNNs.

Sequential counters for encoding cardinality constraints

Consider a *cardinality constraint*: $\sum_{i=1}^m l_i \geq C$, where $l_i \in \{0, 1\}$ is a Boolean variable and C is a constant. This can be compiled into CNF using *sequential counters* $SQ(l, C)$, $l = (l_1, \dots, l_m)$.

SAT Encoding

Naive approach: encode the *BLK* and *O* blocks into CNF is to directly translate their ILP encoding defined above into SAT.

Drawback: inefficient – the resulting encoding will be very large.

Solution: exploit the properties of BNNs.

Sequential counters for encoding cardinality constraints

Consider a *cardinality constraint*: $\sum_{i=1}^m l_i \geq C$, where $l_i \in \{0, 1\}$ is a Boolean variable and C is a constant. This can be compiled into CNF using *sequential counters* $SQ(l, C)$, $l = (l_1, \dots, l_m)$.

Then $SQ(l, C)$ is equivalent to:

$$(l_1 \Leftrightarrow r_{(1,1)}) \wedge (\neg r_{(1,j)}, j = \overline{2, C})$$

$$r_{(i,1)} \Leftrightarrow l_i \vee r_{(i-1,1)} \quad \wedge$$

$$r_{(i,j)} \Leftrightarrow l_i \wedge r_{(i-1,j-1)} \vee r_{(i-1,j)}, j = \overline{2, C}$$

where $i = \overline{2, m}$, $r_{(j,p)} = \mathbb{T} \iff \sum_{i=1}^j l_i \geq p$

SAT Encoding (cont'd)

► Encoding of BLK_k :

$$\langle \mathbf{a}_i, \mathbf{x}_k \rangle \geq C_i \iff \sum_{j=1}^{n_k} a_{ij} x_{k_j} \geq C_i \iff$$

Variable replacement: $x_{k_j} \in \{0, 1\} \rightsquigarrow x_{k_j}^{(b)} \in \{0, 1\}$ with $x_{k_j} = 2x_{k_j}^{(b)} - 1$. We have:

$$\sum_{j=1}^{n_k} a_{ij} (2x_{k_j}^{(b)} - 1) \geq C_i \Rightarrow v_i = 1 \iff$$

Denote: $\mathbf{a}_i^+ = \{j | a_{ij} = 1\}$ and $\mathbf{a}_i^- = \{j | a_{ij} = -1\}$. We have:

$$\sum_{j \in \mathbf{a}_i^+} x_{k_j}^{(b)} - \sum_{j \in \mathbf{a}_i^-} x_{k_j}^{(b)} \geq \left\lceil \frac{C_i}{2} + \sum_{j=1}^{n_k} \frac{a_{ij}}{2} \right\rceil \iff \sum_{j \in \mathbf{a}_i^+} x_{k_j}^{(b)} - \sum_{j \in \mathbf{a}_i^-} (1 - \overline{x_{k_j}^{(b)}}) \geq C_i' \iff$$

$$\sum_{j \in \mathbf{a}_i^+} x_{k_j}^{(b)} - \sum_{j \in \mathbf{a}_i^-} \overline{x_{k_j}^{(b)}} \geq C_i' + |\mathbf{a}_i^-|. \text{ Hence } \sum_{j=1}^{n_k} l_{k_j} \geq D_i \Rightarrow v_i^{(b)} = 1, \quad i = \overline{1, n_{k+1}}$$

Further we have: $\bigwedge_{i=1}^{n_{k+1}} SQ(I, D_i) \wedge \bigwedge_{i=1}^{n_{k+1}} (r_{i(n_k, D_i)} \iff v_i^{(b)})$

Similarly for

$$\langle \mathbf{a}_i, \mathbf{x}_k \rangle < C_i \Rightarrow v_i = -1 \quad i = \overline{1, n_{k+1}}$$

SAT Encoding (cont'd)

Example

Recall the internal block with two inputs and one output: $A_k = [1, -1]$, $b_k = [-0.5]$, $\alpha_k = [0.12]$, $\mu_k = [-0.1]$, $\sigma_k = [2]$, $\delta_k = [0.1]$. We have:

$$x_{k_1} - x_{k_2} \geq -1 \Rightarrow v_1 = 1 \iff 2x_{k_1}^{(b)} - 1 - (2x_{k_2}^{(b)} - 1) \geq -1 \Rightarrow v_1^{(b)} = 1 \iff$$

$$x_{k_1}^{(b)} - x_{k_2}^{(b)} \geq [0.5]$$

$$x_{k_1} - x_{k_2} < -1 \Rightarrow v_1 = -1$$

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Speeding-up the SAT Encoding

- ▶ Takes advantage of the modular structure of BNNs.

Speeding-up the SAT Encoding

- ▶ Takes advantage of the modular structure of BNNs.
- ▶ The approach works for all properties, we exemplify for adversarial robustness.

Speeding-up the SAT Encoding

- ▶ Takes advantage of the modular structure of BNNs.
- ▶ The approach works for all properties, we exemplify for adversarial robustness.
- ▶ The network can be encoded as a conjunction of two Boolean formulas: *Gen* (generator) encodes the first block of the network, and *Ver* (verifier) encodes the rest of the network:

$$BNN_{A_d}(\mathbf{x} + \tau, o, l_x) = Gen(\mathbf{x} + \tau, \mathbf{y}) \wedge Ver(\mathbf{y}, \mathbf{z}, o, l_x)$$

where

$$Gen(\mathbf{x} + \tau, \mathbf{y}) = CNF(\|\tau\|_\infty \leq \epsilon) \wedge \bigwedge_{i=1}^n CNF(\mathbf{x}_i + \tau_i) \in [LB, UB] \wedge \\ BINBLK_1(\mathbf{x} + \tau, \mathbf{y})$$

$$Ver(\mathbf{y}, \mathbf{z}, o, l_x) = BINBLK_2(\mathbf{y}, \mathbf{z}) \wedge BINO(\mathbf{z}, o) \wedge CNF(o \neq l_x).$$

Speeding-up the SAT Encoding

- ▶ Takes advantage of the modular structure of BNNs.
- ▶ The approach works for all properties, we exemplify for adversarial robustness.
- ▶ The network can be encoded as a conjunction of two Boolean formulas: *Gen* (generator) encodes the first block of the network, and *Ver* (verifier) encodes the rest of the network:

$$BNN_{A_d}(\mathbf{x} + \tau, o, l_x) = Gen(\mathbf{x} + \tau, \mathbf{y}) \wedge Ver(\mathbf{y}, \mathbf{z}, o, l_x)$$

where

$$Gen(\mathbf{x} + \tau, \mathbf{y}) = CNF(\|\tau\|_\infty \leq \epsilon) \wedge \bigwedge_{i=1}^n CNF(\mathbf{x}_i + \tau_i) \in [LB, UB] \wedge BINBLK_1(\mathbf{x} + \tau, \mathbf{y})$$

$$Ver(\mathbf{y}, \mathbf{z}, o, l_x) = BINBLK_2(\mathbf{y}, \mathbf{z}) \wedge BINO(\mathbf{z}, o) \wedge CNF(o \neq l_x).$$

- ▶ *Gen* and *Ver* share only \mathbf{y} \rightsquigarrow use *Craig interpolants* to build efficient search procedure.

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.
2. Let $\tilde{\mathbf{y}}$ denote this assignment to \mathbf{y} .

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.
2. Let $\tilde{\mathbf{y}}$ denote this assignment to \mathbf{y} .
3. Check if the assignment $\mathbf{y} = \tilde{\mathbf{y}}$ can be extended to a satisfying assignment for the Ver formula.

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.
2. Let $\tilde{\mathbf{y}}$ denote this assignment to \mathbf{y} .
3. Check if the assignment $\mathbf{y} = \tilde{\mathbf{y}}$ can be extended to a satisfying assignment for the Ver formula.
4. If yes (assign. makes Ver SAT) \rightsquigarrow adversarial perturbation τ found

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.
2. Let $\tilde{\mathbf{y}}$ denote this assignment to \mathbf{y} .
3. Check if the assignment $\mathbf{y} = \tilde{\mathbf{y}}$ can be extended to a satisfying assignment for the Ver formula.
4. If yes (assign. makes Ver SAT) \rightsquigarrow adversarial perturbation τ found
5. If no \rightsquigarrow generate an interpolant I of $\text{Gen}(\mathbf{x} + \tau, \mathbf{y}) \wedge \text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, \sigma, l_x)$ by extracting an UNSAT core of $\text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, \sigma, l_x)$

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.
2. Let $\tilde{\mathbf{y}}$ denote this assignment to \mathbf{y} .
3. Check if the assignment $\mathbf{y} = \tilde{\mathbf{y}}$ can be extended to a satisfying assignment for the Ver formula.
4. If yes (assign. makes Ver SAT) \rightsquigarrow adversarial perturbation τ found
5. If no \rightsquigarrow generate an interpolant I of $\text{Gen}(\mathbf{x} + \tau, \mathbf{y}) \wedge \text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, o, l_x)$ by extracting an UNSAT core of $\text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, o, l_x)$
6. Use assumptions, which are assignments of $\tilde{\mathbf{y}}$, in the SAT solver to obtain a core. Since none of the satisfying assignments to I can be extended to a valid satisfying assignment of $\text{BNN}_{A_d}(l_x + \tau, o, l_x)$, we block them all in Gen by redefining $\text{Gen} := \text{Gen} \wedge \neg I$.

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.
2. Let $\tilde{\mathbf{y}}$ denote this assignment to \mathbf{y} .
3. Check if the assignment $\mathbf{y} = \tilde{\mathbf{y}}$ can be extended to a satisfying assignment for the Ver formula.
4. If yes (assign. makes Ver SAT) \rightsquigarrow adversarial perturbation τ found
5. If no \rightsquigarrow generate an interpolant I of $\text{Gen}(\mathbf{x} + \tau, \mathbf{y}) \wedge \text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, o, l_x)$ by extracting an UNSAT core of $\text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, o, l_x)$
6. Use assumptions, which are assignments of $\tilde{\mathbf{y}}$, in the SAT solver to obtain a core. Since none of the satisfying assignments to I can be extended to a valid satisfying assignment of $\text{BNN}_{A_d}(l_x + \tau, o, l_x)$, we block them all in Gen by redefining $\text{Gen} := \text{Gen} \wedge \neg I$.
7. Repeat from step 1. The procedure terminates since the solution space is reduced.

Speeding-up the SAT Encoding (cont'd)

Definition (Craig Interpolants)

Let A and B be Boolean formulas such that the formula $A \wedge B$ is UNSAT. Then there exists a formula I , called interpolant, such that $\text{vars}(I) = \text{vars}(A) \cap \text{vars}(B)$, $B \wedge I$ is UNSAT and $A \Rightarrow I$. In general, there exist multiple interpolants for the given A and B .

Idea of the approach:

1. first generate a satisfying assignment to variables τ and \mathbf{y} for $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$.
2. Let $\tilde{\mathbf{y}}$ denote this assignment to \mathbf{y} .
3. Check if the assignment $\mathbf{y} = \tilde{\mathbf{y}}$ can be extended to a satisfying assignment for the Ver formula.
4. If yes (assign. makes Ver SAT) \rightsquigarrow adversarial perturbation τ found
5. If no \rightsquigarrow generate an interpolant I of $\text{Gen}(\mathbf{x} + \tau, \mathbf{y}) \wedge \text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, o, l_x)$ by extracting an UNSAT core of $\text{Ver}(\mathbf{y} = \tilde{\mathbf{y}}, \mathbf{z}, o, l_x)$
6. Use assumptions, which are assignments of $\tilde{\mathbf{y}}$, in the SAT solver to obtain a core. Since none of the satisfying assignments to I can be extended to a valid satisfying assignment of $\text{BNN}_{A_d}(l_x + \tau, o, l_x)$, we block them all in Gen by redefining $\text{Gen} := \text{Gen} \wedge \neg I$.
7. Repeat from step 1. The procedure terminates since the solution space is reduced.
8. If the formula $\text{Gen}(\mathbf{x} + \tau, \mathbf{y})$ becomes UNSAT, then there is no valid perturbation τ , i.e., the network is ϵ -robust on image \mathbf{x} .

Contents

Motivation

Preliminaries

Properties of neural networks

Binarized Neural Networks (BNNs)

Encoding the BNNs

Mixed Integer Linear Program (MILP) Encoding

Integer Linear Programming (ILP) Encoding

SAT Encoding

Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Encoding the Properties

1. Adversarial Constraint:

$$BNN_{A_d}(\mathbf{x} + \tau, o, l_x) = CNF(\|\tau\|_\infty \leq \epsilon) \vee \bigwedge_{i=1}^n CNF((\mathbf{x}_i + \tau_i) \in [L, U]) \wedge BNN(\mathbf{x} + \tau, o) \wedge CNF(o \neq l_x)$$

2. Verifying Universal Adversarial Robustness

$$\bigwedge_{i=1}^{|S|} BNN_{A_d}(\mathbf{x}_i + \tau, o_i, l_{x_i}) \iff q_j \wedge CNF\left(\sum_{i=1}^{|S|} q_j \geq \rho|S|\right)$$

3. Verifying Network Equivalence: if

$$\bigwedge_{i=1}^n CNF(x_i \in [L, U]) \wedge BNN_1(\mathbf{x}, o_1) \wedge BNN_2(\mathbf{x}, o_2) \wedge o_1 \neq o_2$$

is UNSAT then the networks are equivalent. If SAT then we obtain a witness image \mathbf{x} .

Contents

Motivation

Preliminaries

- Properties of neural networks

- Binarized Neural Networks (BNNs)

Encoding the BNNs

- Mixed Integer Linear Program (MILP) Encoding

- Integer Linear Programming (ILP) Encoding

- SAT Encoding

- Speeding-up the SAT Encoding

Encoding the Properties

Experimental Results

Experimental Results

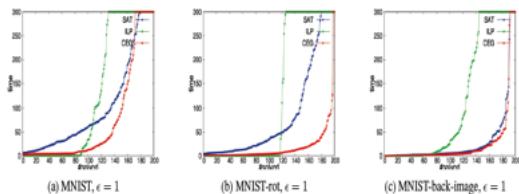
- ▶ Torch framework; Tital Pascal X GPU
- ▶ Datasets: MNIST, MNIST-rot (MINIST where the digits were rotated uniformly in $[0, 2\pi]$ radians), MNIST-back-image (MINIST with a patch from a black-and-white image was used as the background for the digit image)
- ▶ focus on **adversarial robustness**
- ▶ **Architecture**
 - ▶ 4 internal blocks with each block containing a linear layer (LIN) and a final output block.
 - ▶ LIN layer in the first block contains 200 neurons, the LIN layers in other blocks contain 100 neurons
 - ▶ BN and BIN layers in each block were used; additionally a hard *tanh* layer in each internal block was used only during training
 - ▶ For inputs processing, two layers (BN and BIN) were added to the BNN, as the first 2 layers in the network to perform binarization of the grayscale inputs \rightsquigarrow (+) network architecture simplification and search space reduction; (-) lower accuracy of the original BNN by approx.1 %
- ▶ Accuracy of the resulting network on the MNIST, MNISTrot, and MNIST-back-image datasets were 95.7%, 71%, resp. 70%

Experimental Results (cont'd)

Checking **adversarial robustness**:

- ▶ from each dataset randomly picked 20 images correctly classified by the network for each of the 10 classes (corresponding to digits) \rightsquigarrow 200 images
- ▶ for search space reduction, focus on important pixels as defined by saliency map: perturb the top 50% of highly salient pixels in an image; if a valid perturbation that leads to misclassification among this set of pixels can not be found then search again over all pixels of the image.
- ▶ experimented with 3 different maximum perturbation values $\epsilon \in \{1, 3, 5\}$
- ▶ timeout: 300 seconds for each instance
- ▶ compare three methods of searching for adversarial perturbations. (1) ILP method with SCIP solver, (2) pure SAT method for the sequential counters method using Glucose SAT solver, (3) the SAT method from (2) augmented with the counter-example-guided.
- ▶ complexity of the SAT formulae: 1.4 million variables and 5 million clauses: MNIST-rot – approx 7 million clauses; MNIST and MNIST-back – approx 5 and 3 million clauses on average, respectively.
- ▶ the largest instance contains: 3 million variables and 12 million clauses.

Experimental Results (cont'd)



	Solved instances (out of 200)									Certifiably ϵ -robust		
	MNIST			MNIST-rot			MNIST-back-image			#	#	#
	SAT	ILP	CBC	SAT	ILP	CBC	SAT	ILP	CBC			
	#solved (t)	#solved (t)	#solved (t)	#solved (t)	#solved (t)	#solved (t)	#solved (t)	#solved (t)	#solved (t)			
$\epsilon = 1$	180 (77.3)	130 (31.5)	171 (34.1)	179 (57.4)	125 (10.9)	197 (13.5)	191 (18.3)	143 (40.8)	191 (12.8)	138	96	138
$\epsilon = 3$	187 (77.6)	148 (29.0)	181 (35.1)	193 (61.5)	155 (8.3)	198 (13.7)	107 (43.8)	67 (52.7)	119 (44.6)	20	5	21
$\epsilon = 5$	191 (79.5)	165 (29.1)	188 (36.3)	196 (62.7)	170 (11.3)	198 (13.7)	104 (48.8)	70 (53.8)	116 (47.4)	3	-	4

- ▶ **Advantages** of the method: complete search procedure \rightsquigarrow certify ϵ -robustness \rightsquigarrow there exists no adversarial perturbation technique that can fool the network on these images.
- ▶ existing methods: incomplete
- ▶ with increasing ϵ , the number of images on which the network is ϵ -robust decreases as the adversary can leverage the larger value to construct