

## A. Hardness

In this section we show that finding the minimum adversarial distortion with a certified approximation ratio is hard. We first introduce some basic definitions and theorems in Section A.1. We provide some backgrounds about in-approximability reduction in Section A.2. Section A.3 gives a warmup proof for boolean case and then Section A.4 provides the proof of our main hardness result (for network with real inputs).

### A.1. Definitions

We provide some basic definitions and theorems in this section. First, we define the classic 3SAT problem.

**Definition A.1** (3SAT problem). *Given  $n$  variables and  $m$  clauses in a conjunctive normal form CNF formula with the size of each clause at most 3, the goal is to decide whether there exists an assignment to the  $n$  Boolean variables to make the CNF formula to be satisfied.*

For the 3SAT problem in Definition A.1, we introduce the Exponential Time Hypothesis (ETH), which is a common concept in complexity field.

**Hypothesis A.2** (Exponential Time Hypothesis (ETH) (Impagliazzo et al., 1998)). *There is a  $\delta > 0$  such that the 3SAT problem defined in Definition A.1 cannot be solved in  $O(2^{\delta n})$  time.*

ETH had been used in many different problems, e.g. clustering (Ailon et al., 2018; Cohen-Addad et al., 2018), low-rank approximation (Razenshteyn et al., 2016; Song et al., 2017a;b; 2018). For more details, we refer the readers to a survey (Lokshtanov et al., 2013).

Then we define another classical question in complexity theory, the SET-COVER problem, which we will use in our proof. The exact SET-COVER problem is one of Karp’s 21 NP-complete problems known to be NP-complete in 1972:

**Definition A.3** (SET-COVER). *The inputs are  $U, S$ ;  $U = \{1, 2, \dots, n\}$  is a universe,  $P(U)$  is the power set of  $U$ , and  $S = \{S_1, \dots, S_m\} \subseteq P(U)$  is a family of subsets,  $\cup_{j \in [m]} S_j = U$ . The goal is to give a YES/NO answer to the follow decision problem:*

*Does there exist a set-cover of size  $t$ , i.e.,  $\exists C \subseteq [m]$ , such that  $\cup_{j \in C} S_j = U$  with  $|C| = t$ ?*

Alternatively, we can also state the problem as finding the minimum set cover size  $t_0$ , via a binary search on  $t$  using the answers of the decision problem in A.3. The Approximate SET-COVER problem is defined as follows.

**Definition A.4** (Approximate SET-COVER). *The inputs are  $U, S$ ;  $U = \{1, 2, \dots, n\}$  is a universe,  $P(U)$  is the power set of  $U$ , and  $S = \{S_1, \dots, S_m\} \subseteq P(U)$  is a family of subsets,  $\cup_{j \in [m]} S_j = U$ . The goal is to distinguish between the following two cases:*

- (I): There exists a small set-cover, i.e.,  $\exists C \subseteq [m]$ , such that  $\cup_{j \in C} S_j = U$  with  $|C| \leq t$ .
- (II): Every set-cover is large, i.e., every  $C \subseteq [m]$  with  $\cup_{j \in C} S_j = U$  satisfies that  $|C| > \alpha t$ , where  $\alpha > 1$ .

An oracle that solves the Approximate SET-COVER problem outputs an answer  $t_U \geq t_0$  but  $t_U \leq \alpha t_0$  using a binary search, where  $t_U$  is an upper bound of  $t_0$  with a guaranteed approximation ratio  $\alpha$ . For example, we can use a greedy (rather than exact) algorithm to solve the SET-COVER problem, which cannot always find the smallest size of set cover  $t_0$ , but the size  $t_U$  given by the greedy algorithm is at most  $\alpha$  times as large as  $t_0$ .

In our setting, we want to investigate the hardness of finding the lower bound with a guaranteed approximation ration, but an approximate algorithm for SET-COVER gives us an upper bound of  $t_0$  instead of an lower bound of  $t_0$ . However, in the following proposition, we show that finding an lower bound with an approximation ratio of  $\alpha$  is as hard as finding an upper bound with an approximation ratio of  $\alpha$ .

**Proposition A.5.** *Finding a lower bound  $t_L$  for the size of the minimal set-cover (that has size  $t_0$ ) with an approximation ratio  $\alpha$  is as hard as finding an upper bound  $t_U$  with an approximation ratio  $\alpha$ .*

*Proof.* If we find a lower bound  $t_L$  with  $\frac{t_0}{\alpha} \leq t_L \leq t_0$ , by multiplying both sides by  $\alpha$ , we also find an upper bound  $t_U = \alpha t_L$  which satisfies that  $t_0 \leq t_U \leq \alpha t_0$ . So finding an lower bound with an approximation ratio  $\alpha$  is at least as hard as finding an upper bound with an approximation ratio  $\alpha$ . The converse is also true.  $\square$

SET-COVER is a well-studied problem in the literature. Here we introduce a theorem from (Raz & Safra, 1997; Alon et al., 2006; Dinur & Steurer, 2014) which implies the hardness of approximating SET-COVER.

**Theorem A.6** ((Raz & Safra, 1997; Alon et al., 2006; Dinur & Steurer, 2014)). *Unless  $NP = P$ , there is no polynomial time algorithm that gives a  $(1 - o(1)) \ln n$ -approximation to SET-COVER problem with universe size  $n$ .*

We now formally define our neural network robustness verification problems.

**Definition A.7** (ROBUST-NET( $\mathbb{R}$ )). *Given an  $n$  hidden nodes ReLU neural network  $F(x) : \mathbb{R}^d \rightarrow \mathbb{R}$  where all weights are fixed, for a query input vector  $x \in \mathbb{R}^d$  with  $F(x) \leq 0$ . The goal is to give a YES/NO answer to the following decision problem:*

*Does there exist a  $y$  with  $\|x - y\|_1 \leq r$  such that  $F(y) > 0$ ?*

With an oracle of the decision problem available, we can figure out the smallest  $r$  (defined as  $r_0$ ) such that there exists a vector  $y$  with  $\|x - y\|_1 \leq r$  and  $F(y) > 0$  via a binary search.

We also define a binary variant of the ROBUST-NET problem, denoted as ROBUST-NET( $\mathbb{B}$ ). The proof for this variant is more straightforward than the real case, and will help the reader understand the proof for the real case.

**Definition A.8** (ROBUST-NET( $\mathbb{B}$ )). *Given an  $n$  hidden nodes ReLU neural network  $F(x) : \{0, 1\}^d \rightarrow \{0, 1\}$  where weights are all fixed, for a query input vector  $x \in \{0, 1\}^d$  with  $F(x) = 0$ . The goal is to give a YES/NO answer to the following decision problem:*

*Does there exist a  $y$  with  $\|x - y\|_1 \leq r$  such that  $F(y) = 1$ ?*

Then, we define the approximate version of our neural network robustness verification problems.

**Definition A.9** (Approximate ROBUST-NET( $\mathbb{B}$ )). *Given an  $n$  hidden nodes ReLU neural network  $F(x) : \{0, 1\}^d \rightarrow \{0, 1\}$  where weights are all fixed, for a query input vector  $x \in \{0, 1\}^d$  with  $F(x) = 0$ . The goal is to distinguish the following two cases :*

(I): There exists a point  $y$  such that  $\|x - y\|_1 \leq r$  and  $F(y) = 1$ .

(II): For all  $y$  satisfies  $\|x - y\|_1 \leq \alpha r$ , the  $F(y) = 0$ , where  $\alpha > 1$ .

**Definition A.10** (Approximate ROBUST-NET( $\mathbb{R}$ )). *Given an  $n$  hidden nodes ReLU neural network  $F(x) : \mathbb{R}^d \rightarrow \mathbb{R}$  where weights are all fixed, for a query input vector  $x \in \mathbb{R}^d$  with  $F(x) \leq 0$ . The goal is to distinguish the following two cases :*

(I): There exists a point  $y$  such that  $\|x - y\|_1 \leq r$  and  $F(y) > 0$ .

(II): For all  $y$  satisfies  $\|x - y\|_1 \leq \alpha r$ , the  $F(y) \leq 0$ , where  $\alpha > 1$ .

As an analogy to SET-COVER, an oracle that solves the Approximate ROBUST-NET( $\mathbb{R}$ ) problem can output an answer  $r \geq r_0$  but  $r \leq \alpha r_0$ , which is an upper bound of  $r_0$  with a guaranteed approximation ratio  $\alpha$ . With a similar statement as in Proposition A.5, if we divide the answer  $r$  by  $\alpha$ , then we get a lower bound  $r' = \frac{r}{\alpha}$  where  $r' \geq \frac{r_0}{\alpha}$ , which is a lower bound with a guaranteed approximation ratio. If we can solve Approximate ROBUST-NET( $\mathbb{R}$ ), we can get a lower bound with a guaranteed approximation ratio, which is the desired goal of our paper.

## A.2. Background of the PCP theorem

The famous Probabilistically Checkable Proofs (PCP) theorem is the cornerstone of the theory of computational hardness of approximation, which investigates the inherent difficulty in designing efficient approximation algorithms for various optimization problems.<sup>2</sup> The formal definition can be stated as follows,

**Theorem A.11** ((Arora & Safra, 1998; Arora et al., 1998)). *Given a SAT formula  $\phi$  of size  $n$  we can in time polynomial in  $n$  construct a set of  $M$  tests satisfying the following:*

(I) : Each test queries a constant number  $d$  of bits from a proof, and based on the outcome of the queries it either accepts or rejects  $\phi$ .

(II) : (Yes Case / Completeness) If  $\phi$  is satisfiable, then there exists a proof so that all tests accept  $\phi$ .

(III) : (No Case / Soundness) If  $\phi$  is not satisfiable, then no proof will cause more than  $M/2$  tests to accept  $\phi$ .

Note that PCP kind of reduction is slightly different from NP reduction, for more examples (e.g. maximum edge biclique, sparsest cut) about how to use PCP theorem to prove inapproximability results, we refer the readers to (Ambühl et al., 2011).

<sup>2</sup>[https://en.wikipedia.org/wiki/PCP\\_theorem](https://en.wikipedia.org/wiki/PCP_theorem)

### A.3. Warm-up

We state our hardness result for ROBUST-NET( $\mathbb{B}$ ) (boolean inputs case) in this section. The reduction procedure for network with boolean inputs is more straightforward and easier to understand than the real inputs case.

**Theorem A.12.** *Unless  $\text{NP} = \text{P}$ , there is no polynomial time algorithm to give a  $(1 - o(1)) \ln n$ -approximation to ROBUST-NET( $\mathbb{B}$ ) problem (Definition A.9) with  $n$  hidden nodes.*

*Proof.* Consider a set-cover instance, let  $S$  denote a set of sets  $\{S_1, S_2, \dots, S_d\}$  where  $s_j \subseteq [n], \forall j \in [d]$ .

For each set  $S_j$  we create an input node  $u_j$ . For each element  $i \in [n]$ , we create a hidden node  $v_i$ . For each  $i \in [n]$  and  $j \in [d]$ , if  $i \in S_j$ , then we connect  $u_j$  and  $v_i$ . We also create an output node  $w$ , for each  $i \in [n]$ , we connect node  $v_i$  and node  $w$ .

Let  $\mathbf{1}_{i \in S_j}$  denote the indicator function that it is 1 if  $i \in S_j$  and 0 otherwise. Let  $T_i$  denote the set that  $T_i = \{j \mid i \in S_j, \forall j \in [d]\}$ . For each  $i \in [n]$ , we define an activation function  $\phi_i$  satisfies that

$$\phi_i = \begin{cases} 1, & \text{if } \sum_{j \in T_i} u_j \geq 1, \\ 0, & \text{otherwise.} \end{cases}$$

Since  $u_j \in \{0, 1\}$ ,  $\phi_i$  can be implemented in this way using ReLU activations:

$$\phi_i = 1 - \max \left( 0, 1 - \sum_{j \in T_i} u_j \right).$$

Note that  $\sum_{j=1}^d \mathbf{1}_{i \in S_j} = \sum_{j=1}^d u_j$ , because  $u_j = 1$  indicates choosing set  $S_j$  and  $u_j = 0$  otherwise.

For final output node  $w$ , we define an activation function  $\psi$  satisfies that

$$\psi = \begin{cases} 1, & \text{if } \sum_{i=1}^n v_i \geq n, \\ 0, & \text{otherwise.} \end{cases}$$

Since  $v_i \in [n]$ ,  $\psi$  can be implemented as

$$\psi = \max \left( 0, \sum_{i=1}^n v_i - n + 1 \right).$$

We use vector  $x$  to denote  $\{0\}^d$  vector and it is to easy to see that  $F(x) = 0$ . Let  $\alpha > 1$  denote a fixed parameter. Also, we have  $F(y) > 0$  if and only if  $C = \{j \mid y_j = 1\}$  is a set-cover. According to our construction, we can have the following two claims,

**Claim A.13** (Completeness). *If there exists a set-cover  $C \subseteq [d]$  with  $\cup_{j \in C} S_j = [n]$  and  $|C| \leq r$ , then there exists a point  $y \in \{0, 1\}^d$  such that  $\|x - y\|_1 \leq r$  and  $F(y) > 0$ .*

**Claim A.14** (Soundness). *If for every  $C \subseteq [d]$  with  $\cup_{j \in C} S_j = U$  satisfies that  $|C| > \alpha \cdot t$ , then for all  $y \in \{0, 1\}^d$  satisfies that  $\|x - y\|_1 \leq \alpha r$ ,  $F(y) \leq 0$  holds.*

Therefore, using Theorem A.11, Theorem A.6, Claim A.13 and Claim A.14 completes the proof.  $\square$

### A.4. Main result

With the proof for ROBUST-NET( $\mathbb{B}$ ) as a warm-up, we now prove our main hardness result for ROBUST-NET( $\mathbb{R}$ ) in this section.

**Theorem A.15.** *Unless  $\text{NP} = \text{P}$ , there is no polynomial time algorithm to give an  $(1 - o(1)) \ln n$ -approximation to ROBUST-NET( $\mathbb{R}$ ) problem (Definition A.10) with  $n$  hidden nodes.*

*Proof.* Consider a set-cover instance, let  $S$  denote a set of sets  $\{S_1, S_2, \dots, S_d\}$  where  $S_j \subseteq [n], \forall j \in [d]$ . For each set  $S_j$  we create an input node  $u_j$ . For each  $j \in [d]$ , we create a hidden node  $t_j$  and connect  $u_j$  and  $t_j$ .

For each element  $i \in [n]$ , we create a hidden node  $v_i$ . For each  $i \in [n]$  and  $j \in [d]$ , if  $i \in S_j$ , then we connect  $u_j$  and  $v_i$ . Finally, we create an output node  $w$  and for each  $i \in [n]$ , we connect node  $v_i$  and node  $w$ .

Let  $\delta = 1/d$ . For each  $j \in [d]$ , we apply an activation function  $\phi_{1,j}$  on  $t_j$  such that

$$\phi_{1,j} = -\max(0, \delta - u_j) + \max(0, u_j - 1 + \delta)$$

It is easy to see that

$$t_j = \phi_{1,j} = \begin{cases} u_j - \delta & \text{if } u_j \in [0, \delta] \\ u_j - (1 - \delta) & \text{if } u_j \in [1 - \delta, 1] \\ 0 & \text{otherwise.} \end{cases}$$

Let  $T_i$  denote the set that  $T_i = \{j \mid i \in S_j, \forall j \in [d]\}$ . For each  $i \in [n]$ , we need an activation function  $\phi_{2,i}$  on node  $v_i$  which satisfies that

$$\phi_{2,i} \in \begin{cases} [-\delta, 0], & \text{if } \forall j \in T_i, t_j \in [-\delta, 0], \\ [0, \delta], & \text{if } \exists j \in T_i, t_j \in [0, \delta]. \end{cases}$$

This can be implemented in the following way,

$$\phi_{2,i} = \max_{j \in T_i} t_j.$$

For the final output node  $w$ , we define it as

$$w = \min_{i \in [n]} v_i.$$

We use vector  $x$  to denote  $\{0\}^d$  vector and it is to easy to see that  $F(x) = -\delta < 0$ . Let  $\alpha > 1$  denote a fixed parameter.

According to our construction, we can have the following two claims.

**Claim A.16** (Completeness). *If there exists a set-cover  $C \subseteq [d]$  with  $\cup_{j \in C} S_j = [n]$  and  $|C| \leq r$ , then there exists a point  $y \in [0, 1]^d$  such that  $\|x - y\|_1 \leq r$  and  $F(y) > 0$ .*

*Proof.* Without loss of generality, we let the set cover to be  $\{S_1, S_2, \dots, S_r\}$ . Let  $y_1 = y_2 = \dots = y_r = 1$  and  $y_{r+1} = y_{r+2} = \dots = y_d = 0$ . By the definition of  $t_j$ , we have  $t_1 = t_2 = \dots = t_r = \delta$ . Since  $\{S_1, S_2, \dots, S_r\}$  is a set-cover, we know that  $v_i = \delta$  for all  $i \in [n]$ . Then  $F(y) = w = \min_{i \in [n]} v_i = \delta > 0$ . Since we also have  $\|y\|_1 = r$ , the adversarial point is found.  $\square$

**Claim A.17** (Soundness). *If for every  $C \subseteq [d]$  with  $\cup_{j \in C} S_j = U$  satisfies that  $|C| > \alpha \cdot r$ , then for all  $y \in [0, 1]^d$  satisfies that  $\|x - y\|_1 \leq \alpha r(1 - 1/d)$ ,  $F(y) \leq 0$  holds.*

*Proof.* Proof by contradiction. We assume that there exists  $y$  such that  $F(y) > 0$  and  $\|y\|_1 \leq \alpha r(1 - 1/d)$ . Since  $F(y) > 0$ , we have for all  $i$ ,  $v_i > 0$ . Thus there exists  $j \in T_i$  such that  $t_j > 0$ . Let  $\pi : [n] \rightarrow Q$  denote a mapping ( $Q \subseteq [d]$  will be decided later). This means that for each  $i \in [n]$ , there exists  $j \in T_i$ , such that  $1 - \delta < y_j \leq 1$ , and we let  $\pi(i)$  denote that  $j$ .

We define set  $Q \subseteq [d]$  as follows

$$Q = \{j \mid \exists i \in [n], \text{ s.t. } \pi(i) = j \in T_i \text{ and } t_j > 0\}.$$

Since  $\sum_{j \in [d]} |y_j| = \|y\|_1 \leq \alpha r(1 - 1/d)$ , we have

$$\sum_{j \in Q} |y_j| \leq \sum_{j \in [d]} |y_j| \leq \alpha r(1 - 1/d),$$

where the first step follows by  $|Q| \leq d$ .

Because for all  $j \in Q$ ,  $|y_j| > 1 - \delta = 1 - 1/d$ , we have

$$|Q| \leq \frac{\alpha r(1 - 1/d)}{(1 - 1/d)} = \alpha \cdot r.$$

So  $\{S_j\}_{j \in Q}$  is a set-cover with size less than or equal to  $\alpha \cdot r$ , which is a contradiction. □

Therefore, using Theorem A.11, Theorem A.6, Claim A.16 and Claim A.17 completes the proof. □

By making a stronger assumption of ETH, we can have the following stronger result which excludes all  $2^{o(n^c)}$  time algorithms, where  $c > 0$  is some fixed constant:

**Corollary A.18.** *Assuming Exponential Time Hypothesis (ETH, see Hypothesis A.2), there is no  $2^{o(n^c)}$  time algorithm that gives a  $(1 - o(1)) \ln n$ -approximation to ROBUST-NET problem with  $n$  hidden nodes, where  $c > 0$  is some fixed constant.*

*Proof.* It follows by the construction in Theorem A.15 and (Moshkovitz, 2012a;b). □

Note that in (Moshkovitz, 2012a), an additional conjecture, Projection Games Conjecture (PGC) is required for the proof, but the result was improved in (Moshkovitz, 2012b) and PGC is not a requirement any more.

## B. Proof of Theorem 3.5

For a  $m$ -layer ReLU network, assume we know all the pre-ReLU activation bounds  $\mathbf{l}^{(k)}$  and  $\mathbf{u}^{(k)}$ ,  $\forall k \in [m-1]$  for a  $m$ -layer ReLU network and we want to compute the bounds of the  $j$  th output at  $m$  th layer.

The  $j$  th output can be written as

$$f_j(\mathbf{x}) = \sum_{k=1}^{n_{m-1}} \mathbf{W}_{j,k}^{(m)} [\phi_{m-1}(\mathbf{x})]_k + \mathbf{b}_j^{(m)}, \quad (15)$$

$$= \sum_{k=1}^{n_{m-1}} \mathbf{W}_{j,k}^{(m)} \sigma(\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) + \mathbf{b}_j^{(m)}, \quad (16)$$

$$= \sum_{k \in \mathcal{I}_{m-1}^+, \mathcal{I}_{m-1}^-, \mathcal{I}_{m-1}} \mathbf{W}_{j,k}^{(m)} \sigma(\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) + \mathbf{b}_j^{(m)}. \quad (17)$$

For neurons belonging to category (i), i.e.,  $k \in \mathcal{I}_{m-1}^+$ ,

$$\sigma(\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) = \mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}.$$

For neurons belonging to category (ii), i.e.,  $k \in \mathcal{I}_{m-1}^-$ ,

$$\sigma(\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) = 0.$$

Finally, for neurons belonging to Category (iii), i.e.,  $k \in \mathcal{I}_{m-1}$ , we bound their outputs. If we adopt the linear upper and lower bounds in (1) and let  $\mathbf{d}_k^{(m-1)} := \frac{\mathbf{u}_k^{(m-1)}}{\mathbf{u}_k^{(m-1)} - \mathbf{l}_k^{(m-1)}}$ , we have

$$\mathbf{d}_k^{(m-1)} (\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) \leq \sigma(\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) \leq \mathbf{d}_k^{(m-1)} (\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)} - \mathbf{l}_k^{(m-1)}). \quad (18)$$

### B.1. Upper bound

The goal of this section is to prove Lemma B.1.

**Lemma B.1** (Upper bound with explicit function). *Given an  $m$ -layer ReLU neural network function  $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m}$ , parameters  $p, \epsilon$ , there exists two explicit functions  $f^L : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m}$  and  $f^U : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m}$  (see Definition 3.4) such that  $\forall j \in [n_m]$ ,*

$$f_j(\mathbf{x}) \leq f_j^U(\mathbf{x}), \forall \mathbf{x} \in B_p(\mathbf{x}_0, \epsilon).$$

Notice that (18) can be used to construct an upper bound and lower bound of  $f_j(\mathbf{x})$  by considering the signs of the weights  $\mathbf{W}_{j,k}^{(m)}$ . Let  $f_j^{U,m-1}(\mathbf{x})$  be an upper bound of  $f_j(\mathbf{x})$ ;  $f_j^{U,m-1}(\mathbf{x})$  can be constructed by taking the right-hand-side (RHS) of

(18) if  $\mathbf{W}_{j,k}^{(m)} > 0$  and taking the left-hand-side (LHS) of (18) if  $\mathbf{W}_{j,k}^{(m)} < 0$ :

$$f_j^{U,m-1}(\mathbf{x}) = \sum_{k \in \mathcal{I}_{m-1}^+} \mathbf{W}_{j,k}^{(m)} (\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) \quad (19)$$

$$\begin{aligned} &+ \sum_{k \in \mathcal{I}_{m-1}, \mathbf{W}_{j,k}^{(m)} > 0} \mathbf{W}_{j,k}^{(m)} \mathbf{d}_k^{(m-1)} (\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)} - \mathbf{l}_k^{(m-1)}) \\ &+ \sum_{k \in \mathcal{I}_{m-1}, \mathbf{W}_{j,k}^{(m)} < 0} \mathbf{W}_{j,k}^{(m)} \mathbf{d}_k^{(m-1)} (\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) + \mathbf{b}_j^{(m)} \\ &= \sum_{k=1}^{n_{m-1}} \mathbf{W}_{j,k}^{(m)} \mathbf{d}_k^{(m-1)} (\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \mathbf{b}_k^{(m-1)}) - \sum_{k \in \mathcal{I}_{m-1}, \mathbf{W}_{j,k}^{(m)} > 0} \mathbf{W}_{j,k}^{(m)} \mathbf{d}_k^{(m-1)} \mathbf{l}_k^{(m-1)} + \mathbf{b}_j^{(m)}, \quad (20) \end{aligned}$$

$$\begin{aligned} &= \sum_{k=1}^{n_{m-1}} \mathbf{W}_{j,k}^{(m)} \mathbf{d}_k^{(m-1)} \mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) \\ &+ \left( \sum_{k=1}^{n_{m-1}} \mathbf{W}_{j,k}^{(m)} \mathbf{d}_k^{(m-1)} \mathbf{b}_k^{(m-1)} - \sum_{k \in \mathcal{I}_{m-1}, \mathbf{W}_{j,k}^{(m)} > 0} \mathbf{W}_{j,k}^{(m)} \mathbf{d}_k^{(m-1)} \mathbf{l}_k^{(m-1)} + \mathbf{b}_j^{(m)} \right), \quad (21) \end{aligned}$$

where we set  $\mathbf{d}_k^{(m-1)} = 1$  for  $k \in \mathcal{I}_{m-1}^+$  and set  $\mathbf{d}_k^{(m-1)} = 0$  for  $k \in \mathcal{I}_{m-1}^-$  from (19) to (20) and collect the constant terms (independent of  $\mathbf{x}$ ) in the parenthesis from (20) to (21).

If we let  $\mathbf{A}^{(m-1)} = \mathbf{W}^{(m)} \mathbf{D}^{(m-1)}$ , where  $\mathbf{D}^{(m-1)}$  is a diagonal matrix with diagonals being  $\mathbf{d}_k^{(m-1)}$ , then we can rewrite  $f_j^{U,m-1}(\mathbf{x})$  into the following:

$$f_j^{U,m-1}(\mathbf{x}) = \sum_{k=1}^{n_{m-1}} \mathbf{A}_{j,k}^{(m-1)} \mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x}) + \left( \mathbf{A}_{j,:}^{(m-1)} \mathbf{b}^{(m-1)} - \mathbf{A}_{j,:}^{(m-1)} \mathbf{T}_{:,j}^{(m-1)} + \mathbf{b}_j^{(m)} \right) \quad (22)$$

$$= \sum_{k=1}^{n_{m-1}} \mathbf{A}_{j,k}^{(m-1)} \left( \sum_{r=1}^{n_{m-2}} \mathbf{W}_{k,r}^{(m-1)} [\phi_{m-2}(\mathbf{x})]_r \right) + \left( \mathbf{A}_{j,:}^{(m-1)} \mathbf{b}^{(m-1)} - \mathbf{A}_{j,:}^{(m-1)} \mathbf{T}_{:,j}^{(m-1)} + \mathbf{b}_j^{(m)} \right) \quad (23)$$

$$= \sum_{r=1}^{n_{m-2}} \sum_{k=1}^{n_{m-1}} \mathbf{A}_{j,k}^{(m-1)} \mathbf{W}_{k,r}^{(m-1)} [\phi_{m-2}(\mathbf{x})]_r + \left( \mathbf{A}_{j,:}^{(m-1)} \mathbf{b}^{(m-1)} - \mathbf{A}_{j,:}^{(m-1)} \mathbf{T}_{:,j}^{(m-1)} + \mathbf{b}_j^{(m)} \right) \quad (24)$$

$$= \sum_{r=1}^{n_{m-2}} \widetilde{\mathbf{W}}_{j,r}^{(m-1)} [\phi_{m-2}(\mathbf{x})]_r + \widetilde{\mathbf{b}}_j^{(m-1)}. \quad (25)$$

From (21) to (22), we rewrite the summation terms in the parenthesis into matrix-vector multiplications and for each  $j \in [n_m]$  let

$$\mathbf{T}_{k,j}^{(m-1)} = \begin{cases} \mathbf{l}_k^{(m-1)} & \text{if } k \in \mathcal{I}_{m-1}, \mathbf{A}_{j,k}^{(m-1)} > 0 \\ 0 & \text{otherwise} \end{cases}$$

since  $0 \leq \mathbf{d}_k^{(m-1)} \leq 1$ ,  $\mathbf{W}_{j,k}^{(m)} > 0$  is equivalent to  $\mathbf{A}_{j,k}^{(m-1)} > 0$ .

From (22) to (23), we simply write out the inner product  $\mathbf{W}_{k,:}^{(m-1)} \phi_{m-2}(\mathbf{x})$  into a summation form, and from (23) to (24), we exchange the summation order of  $k$  and  $r$ . From (24) to (25), we let

$$\widetilde{\mathbf{W}}_{j,r}^{(m-1)} = \sum_{k=1}^{n_{m-1}} \mathbf{A}_{j,k}^{(m-1)} \mathbf{W}_{k,r}^{(m-1)} \quad (26)$$

$$\widetilde{\mathbf{b}}_j^{(m-1)} = \left( \mathbf{A}_{j,:}^{(m-1)} \mathbf{b}^{(m-1)} - \mathbf{A}_{j,:}^{(m-1)} \mathbf{T}_{:,j}^{(m-1)} + \mathbf{b}_j^{(m)} \right) \quad (27)$$

and now we have (25) in the same form as (15).

Indeed, in (15), the running index is  $k$  and we are looking at the  $m$  th layer, with weights  $\mathbf{W}_{j,k}^{(m)}$ , activation functions  $\phi_{m-1}(\mathbf{x})$  and bias term  $\mathbf{b}_j^{(m)}$ ; in (25), the running index is  $r$  and we are looking at the  $m - 1$  th layer with *equivalent* weights  $\widetilde{\mathbf{W}}_{j,r}^{(m-1)}$ , activation functions  $\phi_{m-2}(\mathbf{x})$  and *equivalent* bias  $\widetilde{\mathbf{b}}_j^{(m-1)}$ . Thus, we can use the same technique from (15) to (25) and obtain an upper bound on the  $f_j^{U,m-1}(\mathbf{x})$  and repeat this procedure until obtaining  $f_j^{U,1}(\mathbf{x})$ , where

$$f_j(\mathbf{x}) \leq f_j^{U,m-1}(\mathbf{x}) \leq f_j^{U,m-2}(\mathbf{x}) \leq \dots \leq f_j^{U,1}(\mathbf{x}).$$

Let the final upper bound  $f_j^U(\mathbf{x}) = f_j^{U,1}(\mathbf{x})$ , and now we have

$$f_j(\mathbf{x}) \leq f_j^U(\mathbf{x}),$$

where  $f_j^U(\mathbf{x}) = [f^U(\mathbf{x})]_j$ ,

$$f_j^U(\mathbf{x}) = \mathbf{A}_{j,:}^{(0)} \mathbf{x} + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{T}_{:,j}^{(k)})$$

and for  $k = 1, \dots, m - 1$ ,

$$\mathbf{A}^{(m-1)} = \mathbf{W}^{(m)} \mathbf{D}^{(m-1)}, \quad \mathbf{A}^{(k-1)} = \mathbf{A}^{(k)} \mathbf{W}^{(k)} \mathbf{D}^{(k-1)},$$

$$\begin{aligned} \mathbf{D}^{(0)} &= \mathbf{I}_{n_0} \\ \mathbf{D}_{r,r}^{(k)} &= \begin{cases} \frac{\mathbf{u}_r^{(k)}}{\mathbf{u}_r^{(k)} - \mathbf{l}_r^{(k)}} & \text{if } r \in \mathcal{I}_k \\ 1 & \text{if } r \in \mathcal{I}_k^+ \\ 0 & \text{if } r \in \mathcal{I}_k^- \end{cases} \\ \mathbf{T}_{r,j}^{(k)} &= \begin{cases} \mathbf{l}_r^{(k)} & \text{if } r \in \mathcal{I}_k, \mathbf{A}_{j,r}^{(k)} > 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

## B.2. Lower bound

The goal of this section is to prove Lemma B.2.

**Lemma B.2** (Lower bound with explicit function). *Given an  $m$ -layer ReLU neural network function  $f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m}$ , parameters  $p, \epsilon$ , there exists two explicit functions  $f^L : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m}$  and  $f^U : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m}$  (see Definition 3.4) such that  $\forall j \in [n_m]$ ,*

$$f_j^L(\mathbf{x}) \leq f_j(\mathbf{x}), \forall \mathbf{x} \in B_p(\mathbf{x}_0, \epsilon).$$

Similar to deriving the upper bound of  $f_j(\mathbf{x})$ , we consider the signs of the weights  $\mathbf{W}_{j,k}^{(m)}$  to derive the lower bound. Let  $f_j^{L,m-1}(\mathbf{x})$  be a lower bound of  $f_j(\mathbf{x})$ ;  $f_j^{L,m-1}(\mathbf{x})$  can be constructed by taking the right-hand-side (RHS) of (18) if  $\mathbf{W}_{j,k}^{(m)} < 0$  and taking the left-hand-side (LHS) of (18) if  $\mathbf{W}_{j,k}^{(m)} > 0$ . Following the procedure in (19) to (25) (except that now the additional bias term is from the set  $k \in \mathcal{I}_{m-1}$ ,  $\mathbf{W}_{j,k}^{(m)} < 0$ ), the lower bound is similar to the upper bound we have derived but but replace  $\mathbf{T}^{(m-1)}$  by  $\mathbf{H}^{(m-1)}$ , where for each  $j \in [n_m]$ ,

$$\mathbf{H}_{k,j}^{(m-1)} = \begin{cases} \mathbf{l}_k^{(m-1)} & \text{if } k \in \mathcal{I}_{m-1}, \mathbf{A}_{j,k}^{(m-1)} < 0 \\ 0 & \text{otherwise.} \end{cases}$$

It is because the linear upper and lower bounds in (1) has the same slope  $\frac{u}{u-1}$  on both sides (i.e.  $\sigma(y)$  is bounded by two lines with the same slope but different intercept), which gives the same  $\mathbf{A}$  matrix and  $\mathbf{D}$  matrix in computing the upper



bound and lower bound of  $f_j(\mathbf{x})$ . This is the key to facilitate a faster computation under this linear approximation (1). Thus, the lower bound for  $f_j(\mathbf{x})$  is:

$$f_j^L(\mathbf{x}) \leq f_j(\mathbf{x}),$$

where  $f_j^L(\mathbf{x}) = [f^L(\mathbf{x})]_j$ ,

$$f_j^L(\mathbf{x}) = \mathbf{A}_{j,:}^{(0)} \mathbf{x} + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{H}_{:,j}^{(k)})$$

and for  $k = 1, \dots, m-1$ ,

$$\mathbf{H}_{r,j}^{(k)} = \begin{cases} \mathbf{l}_r^{(k)} & \text{if } r \in \mathcal{I}_k, \mathbf{A}_{j,r}^{(k)} < 0 \\ 0 & \text{otherwise.} \end{cases}$$

### C. Proof of Corollary 3.7

By Theorem 3.5, for  $\mathbf{x} \in B_p(\mathbf{x}_0, \epsilon)$ , we have  $f_j^L(\mathbf{x}) \leq f_j(\mathbf{x}) \leq f_j^U(\mathbf{x})$ . Thus,

$$f_j(\mathbf{x}) \leq f_j^U(\mathbf{x}) \leq \max_{\mathbf{x} \in B_p(\mathbf{x}, \epsilon)} f_j^U(\mathbf{x}), \quad (28)$$

$$f_j(\mathbf{x}) \geq f_j^L(\mathbf{x}) \geq \min_{\mathbf{x} \in B_p(\mathbf{x}, \epsilon)} f_j^L(\mathbf{x}). \quad (29)$$

Since  $f_j^U(\mathbf{x}) = \mathbf{A}_{j,:}^{(0)} \mathbf{x} + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{T}_{:,j}^{(k)})$ ,

$$\begin{aligned} \gamma_j^U &:= \max_{\mathbf{x} \in B_p(\mathbf{x}_0, \epsilon)} f_j^U(\mathbf{x}) = \max_{\mathbf{x} \in B_p(\mathbf{x}_0, \epsilon)} \left( \mathbf{A}_{j,:}^{(0)} \mathbf{x} + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{T}_{:,j}^{(k)}) \right) \\ &= \left( \max_{\mathbf{x} \in B_p(\mathbf{x}_0, \epsilon)} \mathbf{A}_{j,:}^{(0)} \mathbf{x} \right) + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{T}_{:,j}^{(k)}) \end{aligned} \quad (30)$$

$$= \epsilon \left( \max_{\mathbf{y} \in B_p(\mathbf{0}, 1)} \mathbf{A}_{j,:}^{(0)} \mathbf{y} \right) + \mathbf{A}_{j,:}^{(0)} \mathbf{x}_0 + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{T}_{:,j}^{(k)}) \quad (31)$$

$$= \epsilon \|\mathbf{A}_{j,:}^{(0)}\|_q + \mathbf{A}_{j,:}^{(0)} \mathbf{x}_0 + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{T}_{:,j}^{(k)}). \quad (32)$$

From (30) to (31), we do a transformation of variable  $\mathbf{y} := \frac{\mathbf{x} - \mathbf{x}_0}{\epsilon}$  and therefore  $\mathbf{y} \in B_p(\mathbf{0}, 1)$ . By the definition of dual norm  $\|\cdot\|_*$ :

$$\|\mathbf{z}\|_* = \{\sup_{\mathbf{y}} \mathbf{z}^\top \mathbf{y} \mid \|\mathbf{y}\| \leq 1\},$$

and the fact that  $\ell_q$  norm is dual of  $\ell_p$  norm for  $p, q \in [1, \infty]$ , the term  $\left( \max_{\mathbf{y} \in B_p(\mathbf{0}, 1)} \mathbf{A}_{j,:}^{(0)} \mathbf{y} \right)$  in (31) can be expressed as  $\|\mathbf{A}_{j,:}^{(0)}\|_q$  in (32). Similarly,

$$\begin{aligned} \gamma_j^L &:= \min_{\mathbf{x} \in B_p(\mathbf{x}_0, \epsilon)} f_j^L(\mathbf{x}) = \min_{\mathbf{x} \in B_p(\mathbf{x}_0, \epsilon)} \left( \mathbf{A}_{j,:}^{(0)} \mathbf{x} + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{H}_{:,j}^{(k)}) \right) \\ &= \left( \min_{\mathbf{x} \in B_p(\mathbf{x}_0, \epsilon)} \mathbf{A}_{j,:}^{(0)} \mathbf{x} \right) + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{H}_{:,j}^{(k)}) \\ &= \epsilon \left( \min_{\mathbf{y} \in B_p(\mathbf{0}, 1)} \mathbf{A}_{j,:}^{(0)} \mathbf{y} \right) + \mathbf{A}_{j,:}^{(0)} \mathbf{x}_0 + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{H}_{:,j}^{(k)}) \\ &= -\epsilon \left( \max_{\mathbf{y} \in B_p(\mathbf{0}, 1)} -\mathbf{A}_{j,:}^{(0)} \mathbf{y} \right) + \mathbf{A}_{j,:}^{(0)} \mathbf{x}_0 + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{H}_{:,j}^{(k)}) \end{aligned} \quad (33)$$

$$= -\epsilon \|\mathbf{A}_{j,:}^{(0)}\|_q + \mathbf{A}_{j,:}^{(0)} \mathbf{x}_0 + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} (\mathbf{b}^{(k)} - \mathbf{H}_{:,j}^{(k)}). \quad (34)$$

Again, from (33) to (34), we simply replace  $\left( \max_{\mathbf{y} \in B_p(\mathbf{0}, 1)} -\mathbf{A}_{j,:}^{(0)} \mathbf{y} \right)$  by  $\|\mathbf{A}_{j,:}^{(0)}\|_q = \|\mathbf{A}_{j,:}^{(0)}\|_q$ . Thus, if we use  $\nu_j$  to denote the common term  $\mathbf{A}_{j,:}^{(0)} \mathbf{x}_0 + \mathbf{b}_j^{(m)} + \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} \mathbf{b}^{(k)}$ , we have

$$\gamma_j^U = \epsilon \|\mathbf{A}_{j,:}^{(0)}\|_q - \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} \mathbf{T}_{:,j}^{(k)} + \nu_j, \quad (\text{upper bound})$$

$$\gamma_j^L = -\epsilon \|\mathbf{A}_{j,:}^{(0)}\|_q - \sum_{k=1}^{m-1} \mathbf{A}_{j,:}^{(k)} \mathbf{H}_{:,j}^{(k)} + \nu_j. \quad (\text{lower bound})$$

## D. Algorithms

We present our full algorithms, **Fast-Lin** in Algorithm 1 and **Fast-Lip** in Algorithm 2.

---

### Algorithm 1 Fast Bounding via Linear Upper/Lower Bounds for ReLU (**Fast-Lin**)

---

**Require:** weights and biases of  $m$  layers:  $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(m)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(m)}$ , original class  $c$ , target class  $j$

- 1: **procedure** FAST-LIN( $\mathbf{x}_0, p, \epsilon_0$ )
- 2:     Replace the last layer weights  $\mathbf{W}^{(m)}$  with a row vector  $\bar{\mathbf{w}} \leftarrow \mathbf{W}_{c,:}^{(m)} - \mathbf{W}_{j,:}^{(m)}$  (see Section 3.3.3)
- 3:     Initial  $\epsilon \leftarrow \epsilon_0$
- 4:     **while**  $\epsilon$  has not achieved a desired accuracy and iteration limit has not reached **do**
- 5:          $\mathbf{l}^{(0)}, \mathbf{u}^{(0)} \leftarrow$  don't care
- 6:         **for**  $k \leftarrow 1$  **to**  $m$  **do**  $\triangleright$  Compute lower and upper bounds for ReLU unis for all  $m$  layers
- 7:              $\mathbf{l}^{(k)}, \mathbf{u}^{(k)} \leftarrow$  COMPUTETWOSIDEBOUNDS( $\mathbf{x}_0, \epsilon, p, \mathbf{l}^{(1:k-1)}, \mathbf{u}^{(1:k-1)}, k$ )
- 8:             **if**  $\mathbf{l}^{(m)} > 0$  **then**  $\triangleright \mathbf{l}^{(m)}$  is a scalar since the last layer weight is a row vector
- 9:                  $\epsilon$  is a lower bound; increase  $\epsilon$  using a binary search procedure
- 10:             **else**
- 11:                  $\epsilon$  is not a lower bound; decrease  $\epsilon$  using a binary search procedure
- 12:      $\tilde{\epsilon}_j \leftarrow \epsilon$
- 13:     **return**  $\tilde{\epsilon}_j$   $\triangleright \tilde{\epsilon}_j$  is a certified lower bound  $\beta_L$
- 14: **procedure** COMPUTETWOSIDEBOUNDS( $\mathbf{x}_0, \epsilon, p, \mathbf{l}^{(1:m'-1)}, \mathbf{u}^{(1:m'-1)}, m'$ )
- 15:      $\mathbf{x}_0 \in \mathbb{R}^{n_0}$  : input data vector,  $p$  :  $\ell_p$  norm,  $\epsilon$  : maximum  $\ell_p$ -norm perturbation
- 16:      $\mathbf{l}^{(k)}, \mathbf{u}^{(k)}, k \in [m']$  : layer-wise bounds  $\triangleright$  Step 1: Form  $\mathbf{A}$  matrices
- 17:     **if**  $m' = 1$  **then**  $\triangleright$  First layer bounds do not depend on  $\mathbf{l}^{(0)}, \mathbf{u}^{(0)}$
- 18:          $\mathbf{A}^{(0)} \leftarrow \mathbf{W}^{(1)}$
- 19:     **else**
- 20:         **for**  $k \leftarrow m' - 1$  **to** 1 **do**
- 21:             **if**  $k = m' - 1$  **then**  $\triangleright$  Construct  $\mathbf{D}^{(m'-1)}, \mathbf{A}^{(m'-1)}, \mathbf{H}^{(m'-1)}, \mathbf{T}^{(m'-1)}$
- 22:                 Construct diagonal matrix  $\mathbf{D}^{(k)} \in \mathbb{R}^{n_k \times n_k}$  using  $\mathbf{l}^{(k)}, \mathbf{u}^{(k)}$  according to Eq. (5).
- 23:                  $\mathbf{A}^{(m'-1)} \leftarrow \mathbf{W}^{(m')} \mathbf{D}^{(m'-1)}$
- 24:             **else**  $\triangleright$  Multiply all saved  $\mathbf{A}^{(k)}$  by  $\mathbf{A}^{(m'-1)}$
- 25:                  $\mathbf{A}^{(k)} \leftarrow \mathbf{A}^{(m'-1)} \mathbf{A}^{(k)}$   $\triangleright$  We save  $\mathbf{A}^{(k)}$  for next function call
- 26:                  $\mathbf{T}^{(k)} \leftarrow \mathbf{0}, \mathbf{H}^{(k)} \leftarrow \mathbf{0}$   $\triangleright$  Initialize  $\mathbf{T}^{(k)}$  and  $\mathbf{H}^{(k)}$
- 27:                 **for all**  $r \in \mathcal{I}_k$  **do**
- 28:                     **for**  $j \leftarrow 1$  **to**  $n_k$  **do**
- 29:                         **if**  $\mathbf{A}_{j,r}^{(k)} > 0$  **then**
- 30:                              $\mathbf{T}_{r,j}^{(k)} \leftarrow \mathbf{l}_r^{(k)}$
- 31:                         **else**
- 32:                              $\mathbf{H}_{r,j}^{(k)} \leftarrow \mathbf{l}_r^{(k)}$
- 33:             **for**  $j = 1$  **to**  $n_{m'}$  **do**  $\triangleright$  Step 2: Compute  $\gamma^U$  and  $\gamma^L$
- 34:                  $\nu_j \leftarrow \mathbf{A}_{j,:}^{(0)} \mathbf{x}_0 + \mathbf{b}_j^{(m')}, \mu_j^+ \leftarrow 0, \mu_j^- \leftarrow 0$   $\triangleright$  Initialize  $\nu_j, \mu_j^+, \mu_j^-$
- 35:                 **for**  $k = 1$  **to**  $m' - 1$  **do**  $\triangleright$  This loop is skipped when  $m' = 1$
- 36:                      $\mu_j^+ \leftarrow \mu_j^+ - \mathbf{A}_{j,:}^{(k)} \mathbf{T}_{:,j}^{(k)}, \mu_j^- \leftarrow \mu_j^- - \mathbf{A}_{j,:}^{(k)} \mathbf{H}_{:,j}^{(k)}$   $\triangleright$  According to Eq. (6)
- 37:                      $\nu_j \leftarrow \nu_j + \mathbf{A}_{j,:}^{(k)} \mathbf{b}^{(k)}$   $\triangleright$  According to Eq. (7)
- 38:                  $\nu_j, \mu_j^+, \mu_j^-$  satisfy Definition 3.6
- 39:                  $\gamma_j^U \leftarrow \mu_j^+ + \nu_j + \epsilon \|\mathbf{A}_{j,:}^{(0)}\|_q$
- 40:                  $\gamma_j^L \leftarrow \mu_j^- + \nu_j - \epsilon \|\mathbf{A}_{j,:}^{(0)}\|_q$   $\triangleright$  Definition 3.6
- 41:     **return**  $\gamma^L, \gamma^U$

---

**Algorithm 2** Fast Bounding via Upper Bounding Local Lipschitz Constant (**Fast-Lip**)

**Require:** Weights of  $m$  layers:  $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(m)}$ , original class  $c$ , target class  $j$

- 1: **procedure** FAST-LIP( $\mathbf{x}_0, p, \epsilon$ )
- 2:     Replace the last layer weights  $\mathbf{W}^{(m)}$  with a row vector  $\bar{\mathbf{w}} \leftarrow \mathbf{W}_{c,:}^{(m)} - \mathbf{W}_{j,:}^{(m)}$  (see Section 3.3.3)
- 3:     Run FAST-LIN to find layer-wise bounds  $\mathbf{l}^{(i)}, \mathbf{u}^{(i)}$ , and form  $\mathcal{I}_i^+, \mathcal{I}_i^-, \mathcal{I}_i$  fo all  $i \in [m]$
- 4:      $\mathbf{C}^{(0)} \leftarrow \mathbf{W}^{(1)}, \mathbf{L}^{(0)} \leftarrow \mathbf{0}, \mathbf{U}^{(0)} \leftarrow \mathbf{0}$
- 5:     **for**  $l \leftarrow 1$  to  $m - 1$  **do**
- 6:          $\mathbf{C}^{(l)}, \mathbf{L}^{(l)}, \mathbf{U}^{(l)} = \text{BOUND LAYER GRAD}(\mathbf{C}^{(l-1)}, \mathbf{L}^{(l-1)}, \mathbf{U}^{(l-1)}, \mathbf{W}^{(l+1)}, n_{l+1}, \mathcal{I}_l^+, \mathcal{I}_l^-, \mathcal{I}_l)$
- 7:              $\triangleright \mathbf{v} \in \mathbb{R}^{n_0}$  because the last layer is replaced with a row vector  $\bar{\mathbf{w}}$
- 8:      $\mathbf{v} \leftarrow \max(|\mathbf{C}^{(m-1)} + \mathbf{L}^{(m-1)}|, |\mathbf{C}^{(m-1)} + \mathbf{U}^{(m-1)}|)$       $\triangleright$  All operations are element-wise;
- 9:      $\tilde{\epsilon}_j \leftarrow \min(\frac{g(\mathbf{x}_0)}{\|\mathbf{v}\|_q}, \epsilon)$       $\triangleright q$  is the dual norm of  $p, \frac{1}{p} + \frac{1}{q} = 1$
- 10:     **return**  $\tilde{\epsilon}_j$       $\triangleright \tilde{\epsilon}_j$  is a certified lower bound  $\beta_L$ . We can also bisect  $\tilde{\epsilon}_j$  (omitted).
- 11: **procedure** BOUND LAYER GRAD( $\mathbf{C}, \mathbf{L}, \mathbf{U}, \mathbf{W}, n', \mathcal{I}^+, \mathcal{I}^-, \mathcal{I}$ )
- 12:     **for**  $k \in [n_0]$  **do**      $\triangleright n_0$  is the dimension of  $\mathbf{x}_0$
- 13:         **for**  $j \in [n']$  **do**
- 14:              $\mathbf{C}'_{j,k} \leftarrow \sum_{i \in \mathcal{I}^+} \mathbf{W}_{j,i} \mathbf{C}_{i,k}$
- 15:              $\mathbf{U}'_{j,k} \leftarrow \sum_{i \in \mathcal{I}^+, \mathbf{W}_{j,i} > 0} \mathbf{W}_{j,i} \mathbf{U}_{i,k} + \sum_{i \in \mathcal{I}^+, \mathbf{W}_{j,i} < 0} \mathbf{W}_{j,i} \mathbf{L}_{i,k} +$
- 16:                  $\sum_{i \in \mathcal{I}, \mathbf{W}_{j,i} < 0, \mathbf{C}_{i,k} + \mathbf{L}_{i,k} < 0} \mathbf{W}_{j,i} (\mathbf{C}_{i,k} + \mathbf{L}_{i,k}) + \sum_{i \in \mathcal{I}, \mathbf{W}_{j,i} > 0, \mathbf{C}_{i,k} + \mathbf{U}_{i,k} > 0} \mathbf{W}_{j,i} (\mathbf{C}_{i,k} + \mathbf{U}_{i,k})$
- 17:              $\mathbf{L}'_{j,k} \leftarrow \sum_{i \in \mathcal{I}^+, \mathbf{W}_{j,i} > 0} \mathbf{W}_{j,i} \mathbf{L}_{i,k} + \sum_{i \in \mathcal{I}^+, \mathbf{W}_{j,i} < 0} \mathbf{W}_{j,i} \mathbf{U}_{i,k} +$
- 18:                  $\sum_{i \in \mathcal{I}, \mathbf{W}_{j,i} > 0, \mathbf{C}_{i,k} + \mathbf{L}_{i,k} < 0} \mathbf{W}_{j,i} (\mathbf{C}_{i,k} + \mathbf{L}_{i,k}) + \sum_{i \in \mathcal{I}, \mathbf{W}_{j,i} < 0, \mathbf{C}_{i,k} + \mathbf{U}_{i,k} > 0} \mathbf{W}_{j,i} (\mathbf{C}_{i,k} + \mathbf{U}_{i,k})$
- 19:     **return**  $\mathbf{C}', \mathbf{L}', \mathbf{U}'$

**E. An alternative bound on the Lipschitz constant**

Using the property of norm, we can derive an upper bound of the gradient norm of a 2-layer ReLU network in the following:

$$\begin{aligned} & \|\nabla f_j(\mathbf{x})\|_q \\ &= \|\mathbf{W}_{j,:}^{(2)} \mathbf{\Lambda}^{(1)} \mathbf{W}^{(1)}\|_q \\ &= \|\mathbf{W}_{j,:}^{(2)} (\mathbf{\Lambda}_a^{(1)} + \mathbf{\Lambda}_u^{(1)}) \mathbf{W}^{(1)}\|_q \end{aligned} \tag{35}$$

$$\leq \|\mathbf{W}_{j,:}^{(2)} \mathbf{\Lambda}_a^{(1)} \mathbf{W}^{(1)}\|_q + \|\mathbf{W}_{j,:}^{(2)} \mathbf{\Lambda}_u^{(1)} \mathbf{W}^{(1)}\|_q \tag{36}$$

$$\leq \|\mathbf{W}_{j,:}^{(2)} \mathbf{\Lambda}_a^{(1)} \mathbf{W}^{(1)}\|_q + \sum_{r \in \mathcal{I}_1} \|\mathbf{W}_{j,r}^{(2)} \mathbf{W}_{r,:}^{(1)}\|_q \tag{37}$$

where with a slight abuse of notation, we use  $\mathbf{\Lambda}_a^{(1)}$  to denote the diagonal activation matrix for neurons who are always activated, i.e. its  $(r, r)$  entry  $\mathbf{\Lambda}_{a(r,r)}^{(1)}$  is 1 if  $r \in \mathcal{I}_1^+$  and 0 otherwise, and we use  $\mathbf{\Lambda}_u^{(1)}$  to denote the diagonal activation matrix for neurons whose status are uncertain because they could possibly be active or inactive, i.e. its  $(r, r)$  entry  $\mathbf{\Lambda}_{u(r,r)}^{(1)}$  is 1 if  $r \in \mathcal{I}_1$  and 0 otherwise. Therefore, we can write  $\mathbf{\Lambda}^{(1)}$  as a sum of  $\mathbf{\Lambda}_a^{(1)}$  and  $\mathbf{\Lambda}_u^{(1)}$ .

Note that (35) to (36) is from the sub-additive property of a norm, and (36) to (37) uses the sub-additive property of a norm again and set the uncertain neurons encoding all to 1 because

$$\|\mathbf{W}_{j,:}^{(2)} \mathbf{\Lambda}_u^{(1)} \mathbf{W}^{(1)}\| = \left\| \sum_{r \in \mathcal{I}_1} \mathbf{W}_{j,r}^{(2)} \mathbf{\Lambda}_{u(r,r)}^{(1)} \mathbf{W}_{r,:}^{(1)} \right\| \leq \sum_{r \in \mathcal{I}_1} \|\mathbf{W}_{j,r}^{(2)} \mathbf{\Lambda}_{u(r,r)}^{(1)} \mathbf{W}_{r,:}^{(1)}\| \leq \sum_{r \in \mathcal{I}_1} \|\mathbf{W}_{j,r}^{(2)} \mathbf{W}_{r,:}^{(1)}\|.$$

Notice that (37) can be used as an upper bound of Lipschitz constant and is applicable to compute a certified lower bound for minimum adversarial distortion of a general  $\ell_p$  norm attack. However, this bound is expected to be less tight because we simply include all the uncertain neurons to get an upper bound on the norm in (37).

## F. Details of Experiments in Section 4

### F.1. Methods

Below, we give detailed descriptions on the methods that we compare in Table 1, Table F.1 and Table F.2:

- **Fast-Lin**: Our proposed method of directly bounding network output via **linear** upper/lower bounds for ReLU, as discussed in Section 3.3 and Algorithm 1;
- **Fast-Lip**: Our proposed method based on bounding local **Lipschitz** constant, in Section 3.4 and Algorithm 2;
- **Reluplex**: **Reluplex** (Katz et al., 2017) is a satisfiability modulo theory (SMT) based solver which delivers a true minimum distortion, but is very computationally expensive;
- **LP-Full**: A linear programming baseline method with formulation borrowed from (Wong & Kolter, 2018). Note that we solve the primal LP formulation *exactly* to get a best possible bound. This variant solves **full** relaxed **LP** problems at *every* layer to give a final “adversarial polytope”. Similar to our proposed methods, it only gives a lower bound. We extend this formulation to  $p = 2$  case, where the input constraint becomes quadratic and requires a quadratic constrained programming (QCP) solver, which is usually slower than LP solvers.
- **LP**: Similar to **LP-Full**, but this variant solves only **one** LP problem for the full network at the output neurons and the layer-wise bounds for the neurons in hidden layers are solved by **Fast-Lin**. We also extend it to  $p = 2$  case with QCP constraints on the inputs. **LP** and **LP-Full** are served as our baselines to compare with **Fast-Lin** and **Fast-Lip**;
- **Attacks**: Any successful adversarial example gives a valid *upper bound* for the minimum adversarial distortion. For larger networks where **Reluplex** is not feasible, we run adversarial attacks and obtain an upper bound of minimal adversarial distortions to compare with. We apply the  $\ell_2$  and  $\ell_\infty$  variants of Carlini and Wagner’s attack (CW) (Carlini & Wagner, 2017c) to find the best  $\ell_2$  and  $\ell_\infty$  distortions. We found that the CW  $\ell_\infty$  attack usually finds adversarial examples with smaller  $\ell_\infty$  distortions than using PGD (projected gradient descent). We use EAD (Chen et al., 2018b), a Elastic-Net regularized attack, to find adversarial examples with small  $\ell_1$  distortions. We run CW  $\ell_2$  and  $\ell_\infty$  attacks for 3,000 iterations and EAD attacks for 2,000 iterations;
- **CLEVER**: **CLEVER** (Weng et al., 2018) is an attack-agnostic robustness score based on local Lipschitz constant estimation and provides an estimated lower-bound. It is capable of performing robustness evaluation for large-scale networks but is not a certified lower bound;
- **Op-norm**: Operator norms of weight matrices were first used in (Szegedy et al., 2013) to give a robustness lower bound. We compute the  $\ell_p$  induced norm of weight matrices of each layer and use their product as the global Lipschitz constant  $L_q^j$ . A valid lower bound is given by  $g(\mathbf{x}_0)/L_q^j$  (see Section 3.4). We only need to pre-compute the operator norms once for all the examples.

### F.2. Setup

We use MNIST and CIFAR datasets and evaluate the performance of each method in MLP networks with up to 7 layers or over 10,000 neurons, which is the largest network size for non-trivial and guaranteed robustness verification to date. We use the same number of hidden neurons for each layer and denote a  $m$ -layer network with  $n$  hidden neurons in each layer as  $m \times [n]$ . Each network is trained with a grid search of learning rates from  $\{0.1, 0.05, 0.02, 0.01, 0.005\}$  and weight decays from  $\{10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}\}$  and we select the network with the best validation accuracy. We consider both targeted and untargeted robustness under  $\ell_p$  distortions ( $p = 1, 2, \infty$ ); for targeted robustness, we consider three target classes: a random class, a least likely class and a runner-up class (the class with second largest probability). The reported average scores are an average of 100 images from the test set, with images classified wrongly skipped. Reported time is per image. We use binary search to find the certified lower bounds in **Fast-Lin**, **Fast-Lip**, **LP** and **LP-Full**, and the maximum number of search iterations is set to 15.

We implement our algorithm using Python (with Numpy and Numba)<sup>3</sup>, while for the LP based method we use the highly efficient Gurobi commercial LP solver with Python Interface. All experiments are conducted in single thread mode (we

<sup>3</sup><https://github.com/huanzhang12/CertifiedReLURobustness>

disable the concurrent solver in Gurobi) on a Intel Xeon E5-2683v3 (2.0 GHz) CPU. Despite the inefficiency of Python, we still achieve two orders of magnitudes speedup compared with **LP**, while achieving a very similar lower bound. Our methods are automatically parallelized by Numba and can gain further speedups on a multi-core CPU, but we disabled this parallelization for a fair comparison to other methods.

### F.3. Discussions

In Table 1a (full Table: Table F.1), we compare the lower bound  $\beta_L$  computed by each algorithm to the true minimum distortion  $r_0$  found by **Reluplex**. We are only able to verify 2 and 3 layer MNIST with 20 neurons per hidden layer within reasonable time using **Reluplex**. It is worth noting that the input dimension (784) is very large compared to the network evaluated in (Katz et al., 2017) with only 5 inputs. Lower bounds found by **Fast-Lin** is very close to **LP**, and the gaps are within 2-3X from the true minimum distortion  $r_0$  found by **Reluplex**. The upper bound given by CW  $\ell_\infty$  are also very close to  $r_0$ .

In Table 1b (full Table: Table F.2), we compare **Fast-Lin**, **Fast-Lip** with **LP** and **Op-norm** on larger networks with up to over ten thousands hidden neurons. **Fast-Lin** and **Fast-Lip** are significantly faster than **LP** and are able to verify much larger networks (**LP** becomes very slow to solve exactly on 4-layer MNIST with 4096 hidden neurons, and is infeasible for even larger CIFAR models). **Fast-Lin** achieves a very similar bound comparing with results of **LP** over all smaller models, but being *over two orders of magnitude faster*. We found that **Fast-Lip** can achieve better bounds when  $p = 1$  in two-layers networks, and is comparable to **Fast-Lin** in shallow networks. Meanwhile, we also found that **Fast-Lin** scales better than **Fast-Lip** for deeper networks, where **Fast-Lin** usually provides a good bound even when the number of layers is large. For deeper networks, neurons in the last few layers are likely to have uncertain activations, making **Fast-Lip** being too pessimistic. However, **Fast-Lip** outperforms the global Lipschitz constant based bound (**Op-norm**) which quickly goes down to 0 when the network goes deeper, as **Fast-Lip** is bounding the *local* Lipschitz constant to compute robustness lower bound. In Table F.2, we also apply our method to MNIST and CIFAR models to compare the minimum distortion for *untargeted* attacks. The computational benefit of **Fast-Lin** and **Fast-Lip** is more significant than **LP** because **LP** needs to solve  $n_m$  objectives (where  $n_m$  is the total number of classes), whereas the cost of our methods stay mostly unchanged as we get the bounds for all network outputs simultaneously.

In Table 2, we compute our two proposed lower bounds on neural networks with defending techniques to evaluate the effects of defending techniques (e.g. how much robustness is increased). We train the network with two defending methods, defensive distillation (DD) (Papernot et al., 2016) and adversarial training (Madry et al., 2018) based on robust optimization. For DD we use a temperature of 100, and for adversarial training, we train the network for 100 epochs with adversarial examples crafted by 10 iterations of PGD with  $\epsilon = 0.3$ . The test accuracy for the adversarially trained models dropped from 98.5% to 97.3%, and from 98.6% to 98.1%, for 3 and 4 layer MLP models, respectively. We observe that both defending techniques can increase the computed robustness lower bounds, however adversarial training is significantly more effective than defensive distillation. The lower bounds computed by **Fast-Lin** are close to the desired robustness guarantee  $\epsilon = 0.3$ .

Table F.1. Comparison of our proposed certified lower bounds **Fast-Lin** and **Fast-Lip**, **LP** and **LP-Full**, the estimated lower bounds by **CLEVER**, the exact minimum distortion by **Reluplex**, and the upper bounds by **Attack** algorithms (CW  $\ell_\infty$  for  $p = \infty$ , CW  $\ell_2$  for  $p = 2$ , and EAD for  $p = 1$ ) on 2, 3 layers toy MNIST networks with *only 20 neurons per layer*. Differences of lower bounds and speedup are measured on the two corresponding **bold** numbers in each row, representing the best answer from our proposed algorithms and LP based approaches. **Reluplex** is designed to verify  $\ell_\infty$  robustness so we omit results for  $\ell_2$  and  $\ell_1$ . Note that **LP-Full** and **Reluplex** are very slow and cannot scale to any practical networks, and the purpose of this table is to show how close our fast bounds are compared to the true minimum distortion provided by **Reluplex** and the bounds that are slightly tighter but very expensive (e.g. **LP-Full**).

Toy Networks			Average Magnitude of Distortions on 100 Images							
Network	$p$	Target	Certified Bounds				difference ours vs. LP(-Full)	Exact	Uncertified	
			Our bounds		Our baselines			Reluplex	CLEVER	Attacks
			Fast-Lin	Fast-Lip	LP	LP-Full	(Katz et al., 2017)	(Weng et al., 2018)	CW/EAD	
MNIST $2 \times [20]$	$\infty$	runner-up	<b>0.0191</b>	0.0167	<b>0.0197</b>	0.0197	-3.0%	0.04145	0.0235	0.04384
		rand	<b>0.0309</b>	0.0270	<b>0.0319</b>	0.0319	-3.2%	0.07765	0.0428	0.08060
		least	<b>0.0448</b>	0.0398	<b>0.0462</b>	0.0462	-3.1%	0.11711	0.0662	0.1224
	2	runner-up	<b>0.3879</b>	0.3677	0.4811	<b>0.5637</b>	-31.2%	-	0.4615	0.64669
		rand	<b>0.6278</b>	0.6057	0.7560	<b>0.9182</b>	-31.6%	-	0.8426	1.19630
		least	<b>0.9105</b>	0.8946	1.0997	<b>1.3421</b>	-32.2%	-	1.315	1.88830
	1	runner-up	2.3798	<b>2.8086</b>	2.5932	<b>2.8171</b>	-0.3%	-	3.168	5.38380
		rand	3.9297	<b>4.8561</b>	4.2681	<b>4.6822</b>	+3.7%	-	5.858	11.4760
		least	5.7298	<b>7.3879</b>	6.2062	<b>6.8358</b>	+8.1%	-	9.250	19.5960
MNIST $3 \times [20]$	$\infty$	runner-up	<b>0.0158</b>	0.0094	0.0168	<b>0.0171</b>	-7.2%	0.04234	0.0223	0.04786
		rand	<b>0.0229</b>	0.0142	0.0241	<b>0.0246</b>	-6.9%	0.06824	0.0385	0.08114
		least	<b>0.0304</b>	0.0196	0.0319	<b>0.0326</b>	-6.9%	0.10449	0.0566	0.11213
	2	runner-up	<b>0.3228</b>	0.2142	0.3809	<b>0.4901</b>	-34.1%	-	0.4231	0.74117
		rand	<b>0.4652</b>	0.3273	0.5345	<b>0.7096</b>	-34.4%	-	0.7331	1.22570
		least	<b>0.6179</b>	0.4454	0.7083	<b>0.9424</b>	-34.4%	-	1.100	1.71090
	1	runner-up	<b>2.0189</b>	1.8819	2.2127	<b>2.5010</b>	-19.3%	-	2.950	6.13750
		rand	<b>2.8550</b>	2.8144	3.1000	<b>3.5740</b>	-20.1%	-	4.990	10.7220
		least	3.7504	<b>3.8043</b>	4.0434	<b>4.6967</b>	-19.0%	-	7.131	15.6850

(a) Comparison of bounds

Toy Networks			Average Running Time per Image					
Network	$p$	Target	Certified Bounds				Exact	Speedup ours vs. LP(-full)
			Our bounds		Our baselines		Reluplex	
			Fast-Lin	Fast-Lip	LP	LP-Full	(Katz et al., 2017)	
MNIST $2 \times [20]$	$\infty$	runner-up	<b>3.09 ms</b>	3.49 ms	<b>217 ms</b>	1.74 s	134 s	70X
		rand	<b>3.25 ms</b>	5.53 ms	<b>234 ms</b>	1.93 s	38 s	72X
		least	<b>3.37 ms</b>	8.90 ms	<b>250 ms</b>	1.97 s	360 s	74X
	2	runner-up	<b>3.00 ms</b>	3.76 ms	1.10 s	<b>20.6 s</b>	-	6864X
		rand	<b>3.37 ms</b>	6.16 ms	1.20 s	<b>23.1 s</b>	-	6838X
		least	<b>3.29 ms</b>	9.89 ms	1.27 s	<b>26.4 s</b>	-	8021X
	1	runner-up	2.85 ms	<b>39.2 ms</b>	1.27 s	<b>16.1 s</b>	-	412X
		rand	3.32 ms	<b>54.8 ms</b>	1.59 s	<b>17.3 s</b>	-	316X
		least	3.46 ms	<b>68.1 ms</b>	1.74 s	<b>17.7 s</b>	-	260X
MNIST $3 \times [20]$	$\infty$	runner-up	<b>5.58 ms</b>	3.64 ms	253 ms	<b>6.12 s</b>	4.7 hrs	1096X
		rand	<b>6.12 ms</b>	5.23 ms	291 ms	<b>7.16 s</b>	11.6 hrs	1171X
		least	<b>6.62 ms</b>	7.06 ms	307 ms	<b>7.30 s</b>	12.6 hrs	1102X
	2	runner-up	<b>5.35 ms</b>	3.95 ms	1.22 s	<b>57.5 s</b>	-	10742X
		rand	<b>5.86 ms</b>	5.81 ms	1.27 s	<b>66.3 s</b>	-	11325X
		least	<b>5.94 ms</b>	7.55 ms	1.34 s	<b>77.3 s</b>	-	13016X
	1	runner-up	<b>5.45 ms</b>	39.6 ms	1.27 s	<b>75.0 s</b>	-	13763X
		rand	<b>5.56 ms</b>	52.9 ms	1.47 s	<b>82.0 s</b>	-	14742X
		least	6.07 ms	<b>65.9 ms</b>	1.68 s	<b>85.9 s</b>	-	1304X

(b) Comparison of time

Towards Fast Computation of Certified Robustness for ReLU Networks

Table F.2. Comparison of our proposed certified lower bounds **Fast-Lin** and **Fast-Lip** with other lower bounds (**LP**, **Op-norm**, **CLEVER**) and upper bounds (**Attack** algorithms: CW for  $p = 2, \infty$ , EAD for  $p = 1$ ) on networks with 2-7 layers, where each layer has 1024 or 2048 nodes. Differences of lower bounds and speedup are measured on the two corresponding **bold** numbers in each row. Note that **LP-Full** and **Reluplex** are *computationally infeasible* for all the networks reported here, and “-” indicates the method is computationally infeasible for that network. For **Op-norm**, computation time for each image is negligible as the operator norms can be pre-computed.

Large Networks			Average Magnitude of Distortion on 100 Images					Average Running Time per Image					
Network	$p$	Target	Certified Bounds			diff ours vs. LP	Uncertified		Certified Bounds			Speedup ours vs. LP	
			Our bounds		LP (Baseline)		Op-norm (Szegegy et al., 2013)	CLEVER (Weng et al., 2018)	Attacks CW/EAD	Our bounds			LP (Baseline)
			Fast-Lin	Fast-Lip					Fast-Lin	Fast-Lip			
MNIST $2 \times [1024]$	$\infty$	runner-up	<b>0.02256</b>	0.01802	<b>0.02493</b>	0.00159	-9.5%	0.0447	0.0856	<b>127 ms</b>	167 ms	<b>19.3 s</b>	151X
		rand	<b>0.03083</b>	0.02512	<b>0.03386</b>	0.00263	-8.9%	0.0708	0.1291	<b>156 ms</b>	219 ms	<b>20.8 s</b>	133X
		least	<b>0.03854</b>	0.03128	<b>0.04281</b>	0.00369	-10.0%	0.0925	0.1731	<b>129 ms</b>	377 ms	<b>22.2 s</b>	172X
	2	runner-up	<b>0.46034</b>	0.42027	<b>0.55591</b>	0.24327	-17.2%	0.8104	1.1874	<b>127 ms</b>	196 ms	<b>419 s</b>	3305X
		rand	<b>0.63299</b>	0.59033	<b>0.75164</b>	0.40201	-15.8%	1.2841	1.8779	<b>128 ms</b>	234 ms	<b>195 s</b>	1523X
		least	<b>0.79263</b>	0.73133	<b>0.94774</b>	0.56509	-16.4%	1.6716	2.4556	<b>163 ms</b>	305 ms	<b>156 s</b>	956X
	1	runner-up	2.78786	<b>3.46500</b>	<b>3.21866</b>	0.20601	+7.7%	4.5970	9.5295	117 ms	<b>1.17 s</b>	<b>38.9 s</b>	33X
		rand	3.88241	<b>5.10000</b>	<b>4.47158</b>	0.35957	+14.1%	7.4186	17.259	139 ms	<b>1.40 s</b>	<b>48.1 s</b>	34X
		least	4.90809	<b>6.36600</b>	<b>5.74140</b>	0.48774	+10.9%	9.9847	23.933	151 ms	<b>1.62 s</b>	<b>53.1 s</b>	33X
MNIST $3 \times [1024]$	$\infty$	runner-up	<b>0.01830</b>	0.01021	<b>0.02013</b>	0.00004	-9.1%	0.0509	0.1037	<b>1.20 s</b>	1.81 s	<b>50.4 s</b>	42X
		rand	<b>0.02216</b>	0.01236	<b>0.02428</b>	0.00007	-8.7%	0.0717	0.1484	<b>1.12 s</b>	1.11 s	<b>52.7 s</b>	47X
		least	<b>0.02432</b>	0.01384	<b>0.02665</b>	0.00009	-8.7%	0.0825	0.1777	<b>1.02 s</b>	924 ms	<b>54.3 s</b>	53X
	2	runner-up	<b>0.35867</b>	0.22120	<b>0.41040</b>	0.06626	-12.6%	0.8402	1.3513	<b>898 ms</b>	1.59 s	<b>438 s</b>	487X
		rand	<b>0.43892</b>	0.26980	<b>0.49715</b>	0.10233	-11.7%	1.2441	2.0387	<b>906 ms</b>	914 ms	<b>714 s</b>	788X
		least	<b>0.48361</b>	0.30147	<b>0.54689</b>	0.13256	-11.6%	1.4401	2.4916	<b>925 ms</b>	1.01 s	<b>858 s</b>	928X
	1	runner-up	<b>2.08887</b>	1.80150	<b>2.36642</b>	0.00734	-11.7%	4.8370	10.159	<b>836 ms</b>	3.16 s	<b>91.1 s</b>	109X
		rand	<b>2.59898</b>	2.25950	<b>2.91766</b>	0.01133	-10.9%	7.2177	17.796	<b>863 ms</b>	3.84 s	<b>109 s</b>	126X
		least	<b>2.87560</b>	2.50000	<b>3.22548</b>	0.01499	-10.8%	8.3523	22.395	<b>900 ms</b>	4.20 s	<b>122 s</b>	136X
MNIST $4 \times [1024]$	$\infty$	runner-up	<b>0.00715</b>	0.00219	-	0.00001	-	0.0485	0.08635	<b>1.90 s</b>	4.58 s	-	-
		rand	<b>0.00823</b>	0.00264	-	0.00001	-	0.0793	0.1303	<b>2.25 s</b>	3.08 s	-	-
		least	<b>0.00899</b>	0.00304	-	0.00001	-	0.1028	0.1680	<b>2.15 s</b>	3.02 s	-	-
	2	runner-up	<b>0.16338</b>	0.05244	-	0.11015	-	0.8689	1.2422	<b>2.23 s</b>	3.50 s	-	-
		rand	<b>0.18891</b>	0.06487	-	0.17734	-	1.4231	1.8921	<b>2.37 s</b>	2.72 s	-	-
		least	<b>0.20672</b>	0.07440	-	0.23710	-	1.8864	2.4451	<b>2.56 s</b>	2.77 s	-	-
	1	runner-up	<b>1.33794</b>	0.58480	-	0.00114	-	5.2685	10.079	<b>2.42 s</b>	2.71 s	-	-
		rand	<b>1.57649</b>	0.72800	-	0.00183	-	8.9764	17.200	<b>2.42 s</b>	2.91 s	-	-
		least	<b>1.73874</b>	0.82800	-	0.00244	-	11.867	23.910	<b>2.54 s</b>	3.54 s	-	-
CIFAR $5 \times [2048]$	$\infty$	runner-up	<b>0.00137</b>	0.00020	-	0.00000	-	0.0062	0.00950	<b>24.2 s</b>	60.4 s	-	-
		rand	<b>0.00170</b>	0.00030	-	0.00000	-	0.0147	0.02351	<b>26.2 s</b>	78.1 s	-	-
		least	<b>0.00188</b>	0.00036	-	0.00000	-	0.0208	0.03416	<b>27.8 s</b>	79.0 s	-	-
	2	runner-up	<b>0.06122</b>	0.00951	-	0.00156	-	0.2712	0.3778	<b>34.0 s</b>	60.7 s	-	-
		rand	<b>0.07654</b>	0.01417	-	0.00333	-	0.6399	0.9497	<b>36.8 s</b>	49.4 s	-	-
		least	<b>0.08456</b>	0.01778	-	0.00489	-	0.9169	1.4379	<b>37.4 s</b>	49.8 s	-	-
	1	runner-up	<b>0.93835</b>	0.22632	-	0.00000	-	4.0755	7.6529	<b>36.5 s</b>	70.6 s	-	-
		rand	<b>1.18928</b>	0.31984	-	0.00000	-	9.7145	21.643	<b>37.5 s</b>	53.6 s	-	-
		least	<b>1.31904</b>	0.38887	-	0.00001	-	12.793	34.497	<b>38.3 s</b>	48.6 s	-	-
CIFAR $6 \times [2048]$	$\infty$	runner-up	<b>0.00075</b>	0.00005	-	0.00000	-	0.0054	0.00770	<b>37.2 s</b>	106 s	-	-
		rand	<b>0.00090</b>	0.00007	-	0.00000	-	0.0131	0.01866	<b>37.0 s</b>	119 s	-	-
		least	<b>0.00095</b>	0.00008	-	0.00000	-	0.0199	0.02868	<b>37.2 s</b>	126 s	-	-
	2	runner-up	<b>0.03463</b>	0.00228	-	0.00476	-	0.2394	0.2979	<b>56.1 s</b>	99.5 s	-	-
		rand	<b>0.04129</b>	0.00331	-	0.01079	-	0.5860	0.7635	<b>60.2 s</b>	95.6 s	-	-
		least	<b>0.04387</b>	0.00385	-	0.01574	-	0.8756	1.2111	<b>61.8 s</b>	88.6 s	-	-
	1	runner-up	<b>0.59638</b>	0.05647	-	0.00000	-	3.3569	6.0112	<b>57.2 s</b>	108 s	-	-
		rand	<b>0.72178</b>	0.08212	-	0.00000	-	8.2507	17.160	<b>61.4 s</b>	88.2 s	-	-
		least	<b>0.77179</b>	0.09397	-	0.00000	-	12.603	28.958	<b>62.1 s</b>	65.1 s	-	-
CIFAR $7 \times [1024]$	$\infty$	runner-up	<b>0.00119</b>	0.00006	-	0.00000	-	0.0062	0.0102	<b>10.5 s</b>	27.3 s	-	-
		rand	<b>0.00134</b>	0.00008	-	0.00000	-	0.0112	0.0218	<b>10.6 s</b>	29.2 s	-	-
		least	<b>0.00141</b>	0.00010	-	0.00000	-	0.0148	0.0333	<b>11.2 s</b>	30.9 s	-	-
	2	runner-up	<b>0.05279</b>	0.00308	-	0.00020	-	0.2661	0.3943	<b>16.3 s</b>	28.2 s	-	-
		rand	<b>0.05938</b>	0.00407	-	0.00029	-	0.5145	0.9730	<b>16.9 s</b>	27.3 s	-	-
		least	<b>0.06249</b>	0.00474	-	0.00038	-	0.6253	1.3709	<b>17.4 s</b>	27.6 s	-	-
	1	runner-up	<b>0.76647</b>	0.07028	-	0.00000	-	4.815	7.9987	<b>16.9 s</b>	27.8 s	-	-
		rand	<b>0.86467</b>	0.09239	-	0.00000	-	8.630	22.180	<b>17.6 s</b>	26.7 s	-	-
		least	<b>0.91127</b>	0.10639	-	0.00000	-	11.44	31.529	<b>17.5 s</b>	23.5 s	-	-
MNIST $3 \times [1024]$	$\infty$	untargeted	<b>0.01808</b>	0.01016	<b>0.01985</b>	0.00004	-8.9%	0.0458	0.0993	<b>915 ms</b>	2.17 s	<b>227 s</b>	248X
			<b>0.35429</b>	0.21833	-	0.06541	-	0.7413	1.1118	<b>950 ms</b>	2.02 s	-	-
			<b>2.05645</b>	1.78300	<b>2.32921</b>	0.00679	-11.7%	3.9661	9.0044	<b>829 ms</b>	4.41 s	<b>537 s</b>	648X
CIFAR $5 \times [2048]$	$\infty$	untargeted	<b>0.00136</b>	0.00020	-	0.00000	-	0.0056	0.00950	<b>24.1 s</b>	72.9 s	-	-
			<b>0.06097</b>	0.00932	-	0.00053	-	0.2426	0.3702	<b>34.2 s</b>	77.0 s	-	-
			<b>0.93429</b>	0.22535	-	0.00000	-	3.6704	7.3687	<b>35.6 s</b>	90.2 s	-	-