

Verifying Neural Networks ReluPlex

Verifying cyberphysical systems

Sayan Mitra

mitras@illinois.edu

Reference : Introduction to Neural Network Verification by Aws Albarghouthi

Outline

Modeling neural networks

Requirements of neural networks

Reluplex algorithm

References

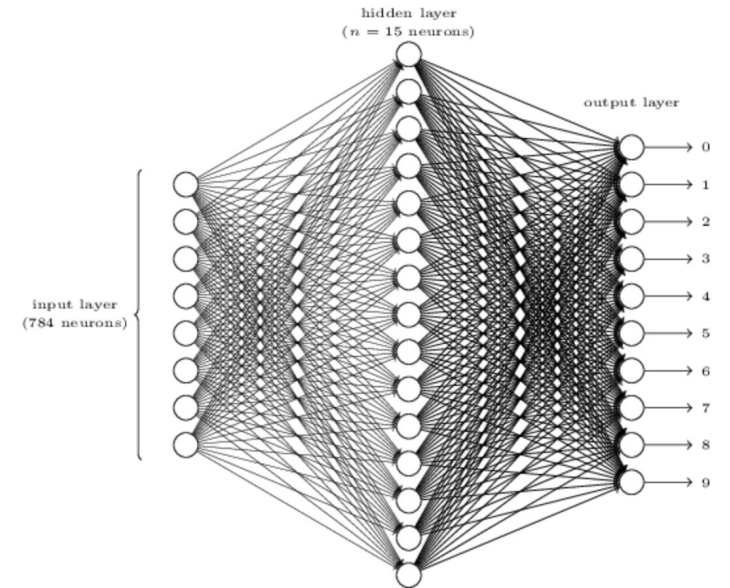
- Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks by Guy Katz, Clark Barrett, David Dill, Kyle Julian, Mykel Kochenderfer, CAV 2017
- Book: Introduction to Neural Network Verification by Aws Albarghouthi, 2021
- Further exploration: Neural Network Verification competition: <https://sites.google.com/view/vnn20/>
- [Neural Networks and Deep Learning](#)

Neural networks as graphs

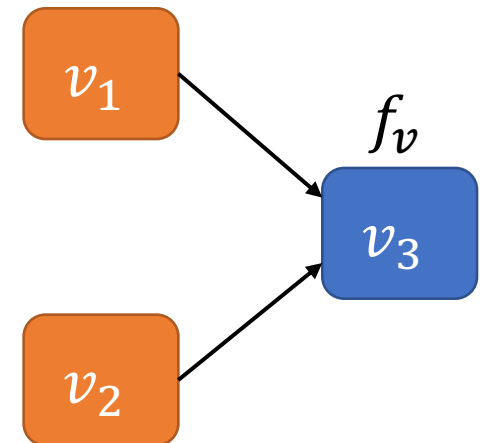
- Consider a *network* $G = (V, E)$
- Let $V^{in}, V^o \subseteq V$ be nonempty sets of input, output vertices
- Each $v \in V \setminus V^{in}$ is associated with a function $f_v: \mathbb{R}^{ind(v)} \rightarrow \mathbb{R}$ where $ind(v)$ is the indegree of v
- Typically these functions are differentiable (almost everywhere)

Assume w.l.o.g

- All nodes are reachable from some V^{in}
- Every node reach some V^o
- There is a total ordering on V and another on E

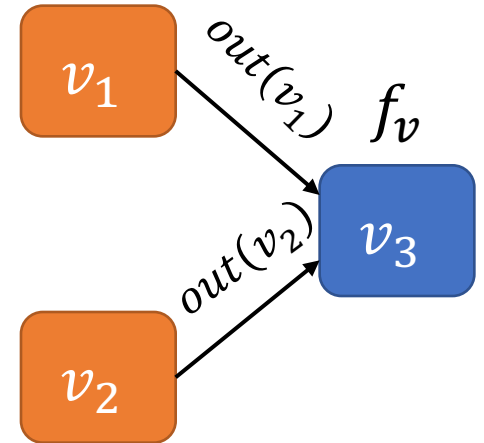


<http://neuralnetworksanddeeplearning.com/chap1.html>



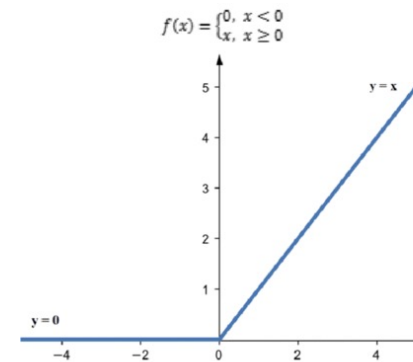
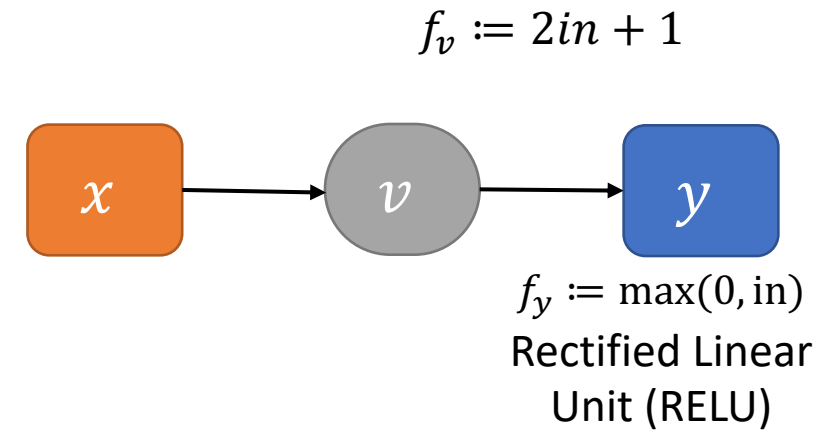
Computation performed by a Network

- $|V^{in}| = n, |V^o| = m$
- The network $G = (V, E)$ defines a function $\mathbb{R}^n \rightarrow \mathbb{R}^m$ defined as follows
- For any non input node $v \in V$ the output from v is recursively defined as follows: let $(v_1, v), (v_2, v), \dots, (v_{ind(v)}, v)$ be the ordered set of all input edges to v
$$out(v) := f_v(out(v_1), \dots, out(v_{ind(v)}))$$
 - Base case: For input vector $\mathbf{x} \in \mathbb{R}^n$, for any input node $v \in V^{in}$, output $out(v) := \mathbf{x}_i$

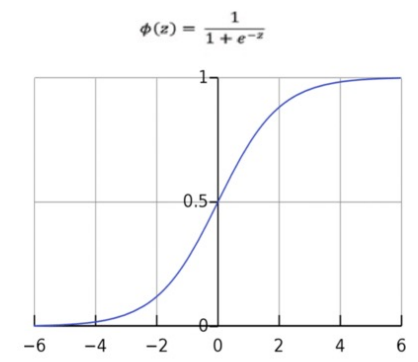


Example functions

$$\begin{aligned} out(y) &= f_y(f_v(out(x))) \\ &= f_y(f_v(x)) \\ &= f_y(2x + 1) \\ &= \max(0, 2x + 1) \end{aligned}$$



RELU



Sigmoid

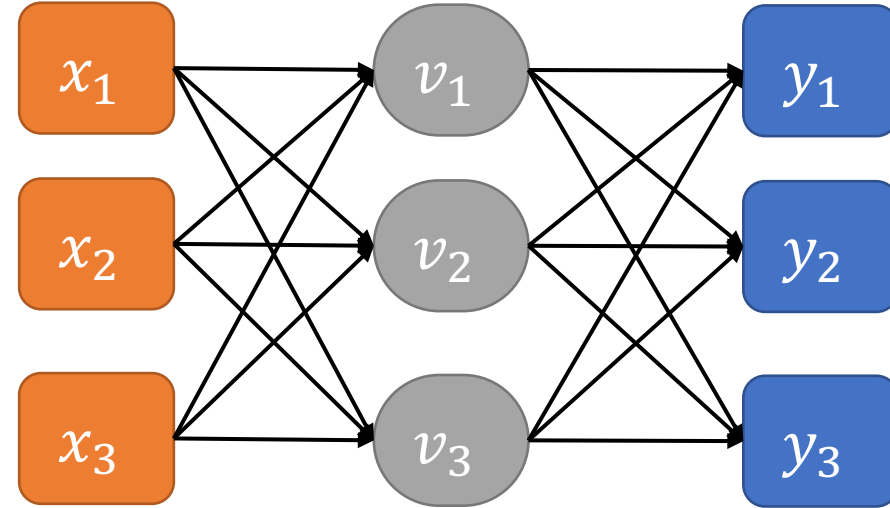
$$f_y := \sigma(in) = \frac{1}{1 + \exp(-in)}$$

Example functions

$$\begin{aligned} \text{out}(y_1) &= f_y \left(\sum_{i=1}^3 \text{out}(v_i) w_i \right) \\ &= f_y \left(\sum_{i=1}^3 \left(\sum_{j=1}^3 x_j u_j \right) w_i \right) \end{aligned}$$

Multi-layered Perceptron

Hidden layer



Softmax

$$f_{y_i} := \frac{\exp(v_i)}{\sum_{i=1}^3 \exp(v_i)}$$

Typical requirements studied for neural networks

For any input x (image, text, program, controller) $f(x)$ meets <some condition>

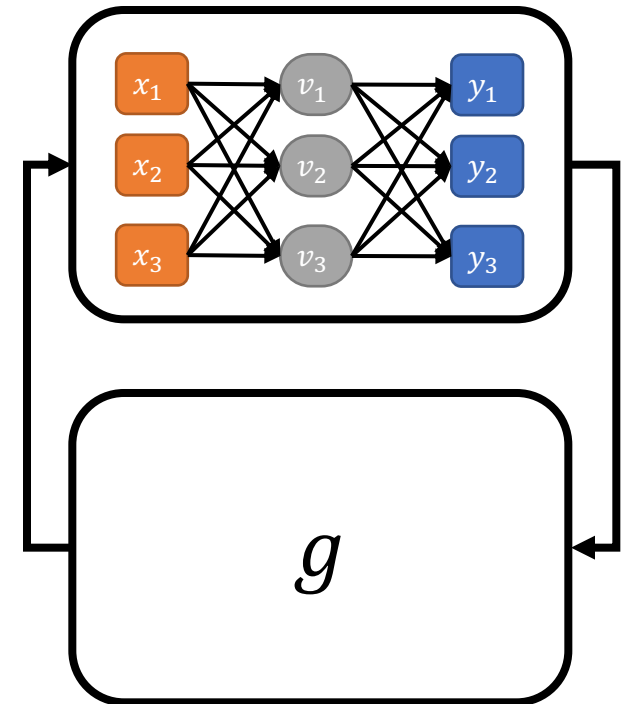
To check the validity of this statement, we can check the satisfiability of the negation, i.e., does there exist some input x such that $f(x)$ does not meet <some condition>

For any inputs x, y (image, text, program, controller) $f(x)$ and $f(y)$ meets <some condition>

Pros. Same form as pre-postconditions of automata transitions

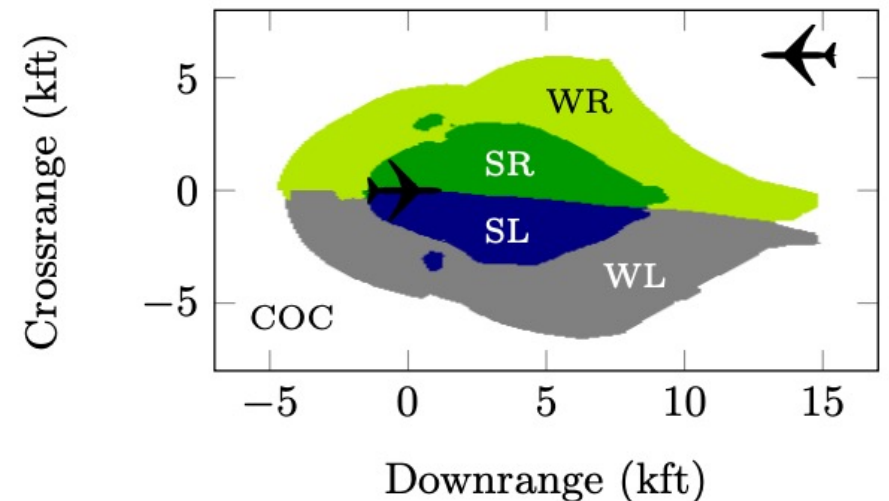
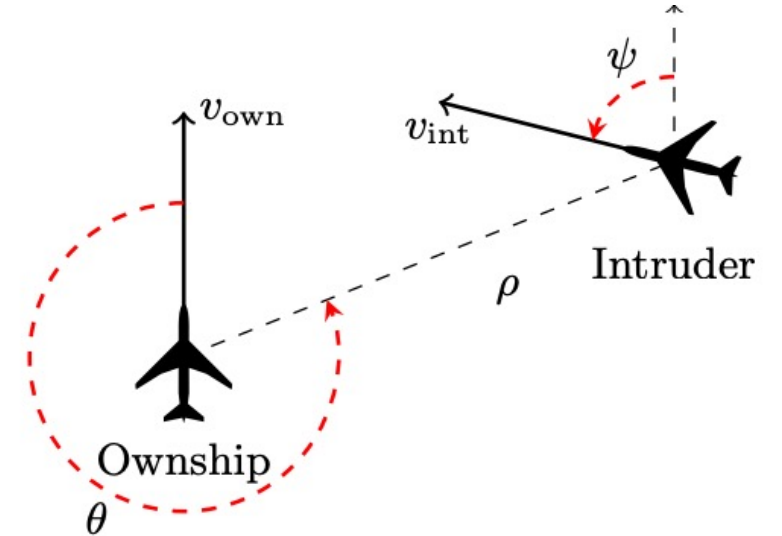
Cons. with this framing

- <some condition> usually requires domain knowledge
- can only be specified *locally*
- Does not cover requirements that use the NN in closed-loop system



Application: ACAS Xu system

- 3MB DNN represents a large (2GB) lookup table for collision avoidance of unmanned aircraft
- Input: $\mathbf{x} \in \mathbb{R}^7$ ρ : Distance; θ : relative angle; ψ : relative heading; v_{own} , v_{in} : Speeds, τ : Time until loss of vertical separation; a_{prev} : Previous advisory.
- Output: Clear of Conflict (COC) or advisor weak/strong left/right.
- Network: 45 networks produced by discretizing τ and a_{prev} , each with 5 inputs and 5 outputs
- Requirement: E.g. If the intruder is far then the score for COC should be above some threshold
 $\forall \mathbf{x} \in \mathbb{R}^7 \mathbf{x}.d > 55947, \mathbf{x}.v_{own} > 1145, \mathbf{x}.v_{in} \leq 60 \Rightarrow f(\mathbf{x}) \geq 1500$



General formulation of correctness

(precondition \wedge neural network) \Rightarrow postcondition

$$\left(\bigwedge_{i=1}^n |x_i - c_i| \leq 0.1 \right) \wedge F_G \wedge \left(\bigwedge_{i=1}^n x_i = v_i^o \right) \wedge \left(\bigwedge_{i=k}^{k+m} r_i = v_i^o \right) \wedge \left(\bigwedge_{i=1}^m r_1 > r_{1+i} \right)$$

Precondition P

$$r_1 == f_{G1}(x_1)$$

$$r_2 == f_{G2}(x_2)$$

...

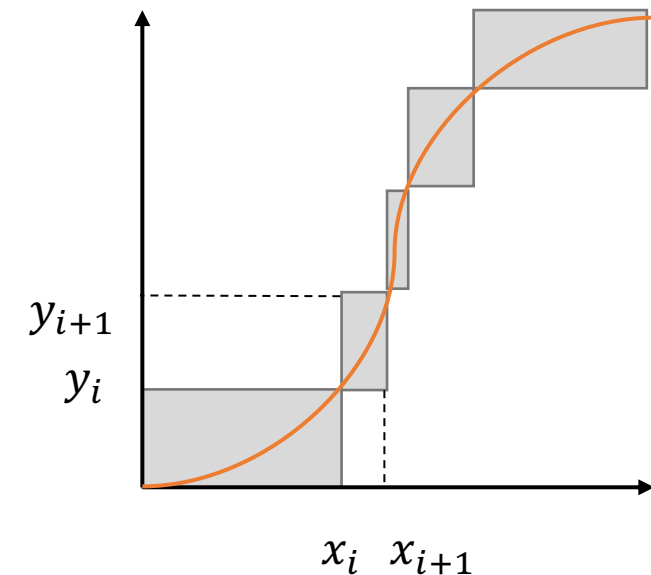
$$r_l == f_{Gl}(x_l)$$

Postcondition Q

Precondition and Post-condition can be encoded in LRA. Solve using DPLL^{LRA}

Sigmoids and ReLUs and exponential blow-up

- Sigmoids can be approximated by piece-wise relations
- Others have considered MILP encodings
- Without disjunctions the problem is an LP
 - Polynomial Time solvable, Simplex
- If for all the inputs, the ReLUs are either active (x) or inactive (0)
 - We can replace each ReLU with either $x_i = f(x)$ or $x_i = 0$ based on lightweight analysis of which is which
- Otherwise, DPLL has to handle the disjunctions and in the worst case will have to consider every possible case (active or inactive) for every possible ReLU leading to exponential number of calls to Simplex
- Reluplex [Katz et al. 2017] addresses this problem



Reluplex Decision Procedure for ReLU NNs

Input: F in Reluplex form

Output: $\exists \mathbf{x} \in \mathbb{R}^m$ such that $\mathbf{x} \models F$?

- Delay case splitting on ReLUs
- In the worst case it is still exponential, but has been shown to be empirically better than DPLL^{LRA}

Reluplex Decision Procedure for ReLU NNs

Input: F in Reluplex form

Output: $\exists \mathbf{x} \in \mathbb{R}^m$ such that $\mathbf{x} \models F$?

Reluplex form

- Equations (same as Simplex)
- Bounds (same as Simplex)
- Relu: $x_i = \text{relu}(x_j)$

Given conjunction of inequalities and ReLU constraints we can convert them to Reluplex form and add $x_i \geq 0$.

Reluplex

Input: A formula F in Reluplex form

Output: $x \models F$ or UNSAT

$x := \langle x_i \mapsto 0 \rangle$; $F' :=$ non ReLU part of F

while true do

$r := \text{Simplex}(F', x)$

if r is unsat **then return UNSAT**

else if $r \models F$ **then return** r

// Handle violated ReLU constraints

Let $x_i = \text{relu}(x_j) \in F$ such that $x[x_i \neq \text{relu}(x[x_j])$

// Case split

if $u_j > 0, l_j < 0$ and $x_i = \text{relu}(x_j)$ considered $> \tau$ times

Reluplex

Input: A formula F in Reluplex form

Output: $x \models F$ or UNSAT

$x := \langle x_i \mapsto 0 \rangle$; $F' :=$ non ReLU part of F

while true do

$r := \text{Simplex}(F', x)$

if r is unsat **then return UNSAT**

else if $r \models F$ **then return** r

// Handle violated ReLU constraints

Let $x_i = \text{relu}(x_j) \in F$ such that $x[x_i \neq \text{relu}(x[x_j])$

if x_i is basic **then** pivot x_i with non-basic variable x_k where $k \neq j$ and $c_{ik} \neq 0$

if x_j is basic **then** pivot x_j with non-basic variable x_k where $k \neq i$ and $c_{jk} \neq 0$

$x[x_i := \text{relu}(x[x_j])$ **OR** $x[x_j := x[x_i$

// Case split

if $u_j > 0, l_j < 0$ and $x_i = \text{relu}(x_j)$ considered $> \tau$ times

$r_1 = \text{Reluplex}(F \wedge x_j \geq 0 \wedge x_j = x_j)$

$r_2 = \text{Reluplex}(F \wedge x_j \leq 0 \wedge x_i = 0)$

if $r_1 = r_2 = \text{unsat}$ **then return UNSAT**

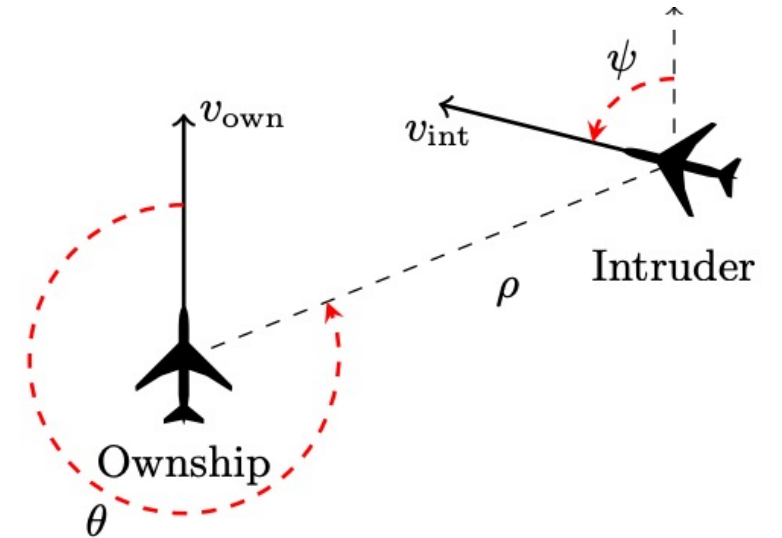
else if $r_1 = \text{unsat}$ **then return** r_1 **else return** r_2

Termination

- The last part of the algorithm ensures termination
- Otherwise, the algorithm could loop forever between fixing relu constraints and Simplex
- $F \equiv x_i = \text{relu}(x_j)$ is split into two cases
 - $F_1 \equiv x_j \geq 0 \wedge x_i = x_j$
 - $F_2 \equiv x_j \leq 0 \wedge x_i = 0$
 - $F' \equiv (F \wedge F_1) \vee (F \wedge F_2)$

Performance of Reluplex on ACAS

	φ_1	φ_2	φ_3	φ_4	φ_5	φ_6	φ_7	φ_8
CVC4	-	-	-	-	-	-	-	-
Z3	-	-	-	-	-	-	-	-
Yices	1	37	-	-	-	-	-	-
MathSat	2040	9780	-	-	-	-	-	-
Gurobi	1	1	1	-	-	-	-	-
Reluplex	8	2	7	7	93	4	7	9



Summary

- Input-output and robustness requirements of neural networks can be expressed as satisfiability queries
- Without ReLU and disjunctions the requirements can be checked efficiently, e.g., Simplex
- Reluplex implements smarter case splitting