

ECE598HZ: Advanced Topics in Machine Learning and Formal Methods

Lecture 13: Neural Network Verification

falsification, training verifiable NNs, and practical verifiers

Prof. Huan Zhang

huan@huan-zhang.com

Deadlines

Project proposal due **3/5** 11:59 pm

Discuss your class projects on Canvas, try to find a teammate with common interests.

- To be effective, *don't* just post your name and say you are looking for a teammate
- Be sure to introduce your technical strengths, research interests, and your thoughts about the project to find people with similar interests
- Do it as soon as possible! Discuss with your teammates to finalize project ideas

Deadlines

Project proposal due **3/5** 11:59 pm

Up to 4 pages, template provided

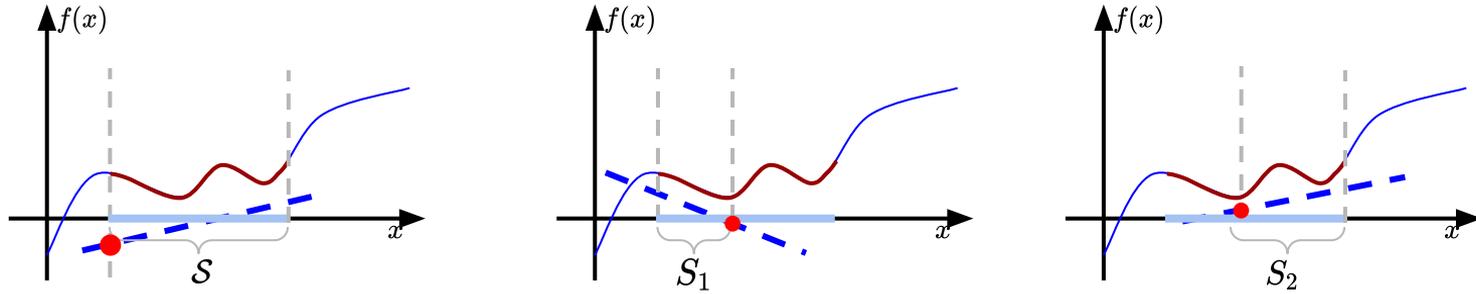
- Introduction of the problem or system under study and why it is important
 - Give clear mathematical description of your problem
 - Related work (What has been done before? do a thorough literature review!)
 - Proposed methodology (what is your planned technique to solve this problem? What are the risks?)
 - Timeline and targets (what goals do you aim to achieve?)
-
- **Some project ideas can be found at the end of Lecture 1 slides**

Deadlines

HW2 due **3/17** 11:59 pm

- Two programming questions
- A lot of bonus questions to help you get better grads
- **Start early!**

Review: bound propagation & branch and bound

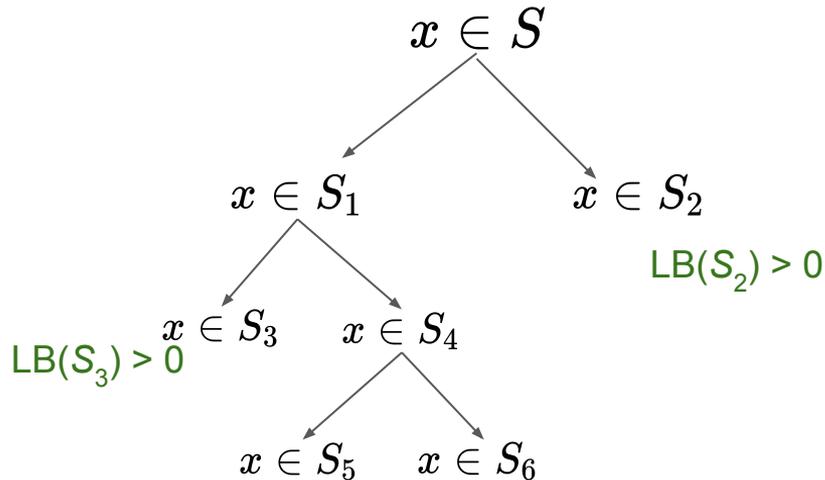


Goal: improve the loose lower bound

Review: Branch and bound

If $LB(S_i) > 0$, it can be removed from our problem since the property is verified on this subdomain S_i ; branch and bound is needed for unverified subdomains only.

$$y^* = \min_{x \in \mathcal{S}} f(x)$$



List of unverified subproblems

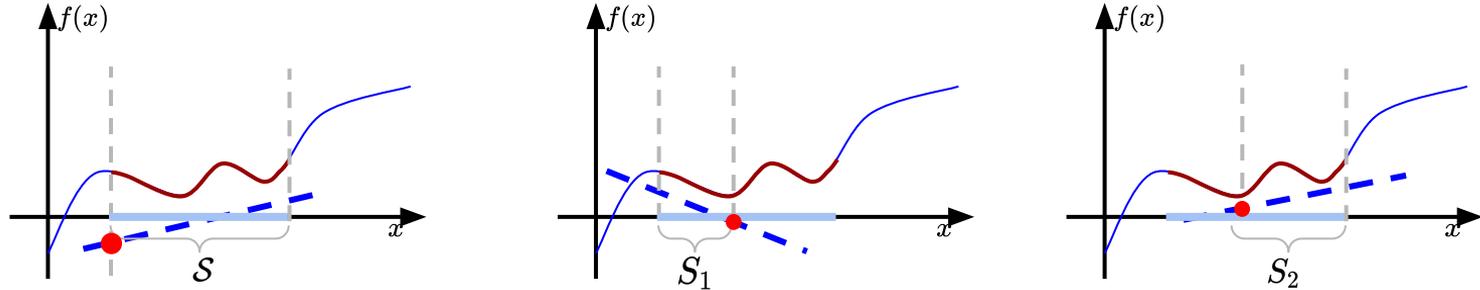
$\{S\}$

~~$\{S_1, S_2\}$~~

~~$\{S_3, S_4\}$~~

$\{S_5, S_6\}$

Review: Branch and bound on input



Split each into domain S , typically by

$$S = \{x_1 \in [-1, 1], x_2 \in [-1, 1]\} \Rightarrow$$

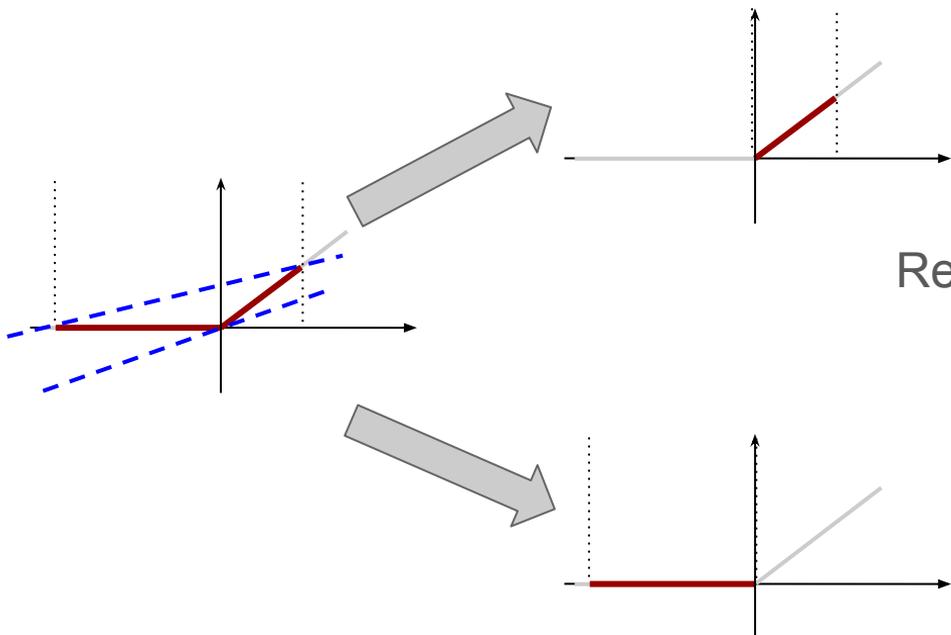
$$S_1 = \{x_1 \in [-1, 0], x_2 \in [-1, 1]\}, S_2 = \{x_2 \in [0, 1], x_2 \in [-1, 1]\}$$

Implementation is easy

Does not work well when input dimension is very high (e.g., image inputs)

Review: Branch and bound on ReLU

Implicitly split input domain S by considering a ReLU neuron in two cases: active and inactive.



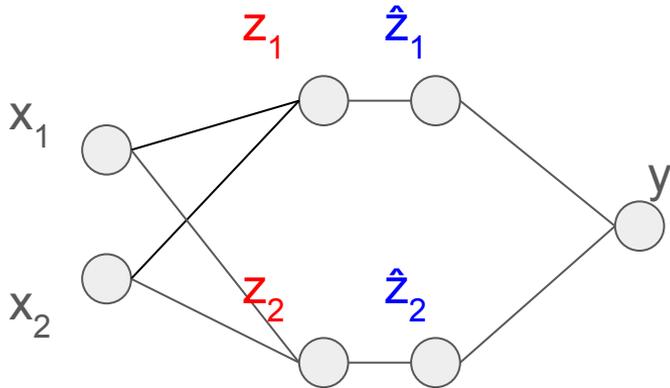
ReLU becomes linear in both subproblems with split constraint (handled using **β -CROWN**)

Review: CROWN with neuron split

Let's look at the **lower bound** only (since only lower bound is needed)

Step 1: bound y using linear functions of y (base case): $y \leq y$, $y \geq y$

Step 2: bound y using linear functions of $\hat{\mathbf{z}}$: simply plugin the definition of the second linear layer: $y = \hat{z}_1 - \hat{z}_2$



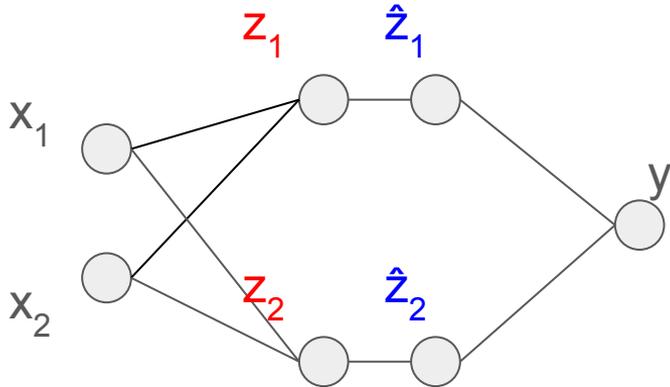
$$y \geq \hat{z}_1 - \hat{z}_2$$

Linear layer: simple substitution

Review: β -CROWN with neuron split

Step 3: bound y using linear functions of z : need linear bounds for ReLU functions, which allows us to replace \hat{z} with z

ReLU layer: use linear bound
Check **sign** of coefficients and take the lower or upper bound



Change in bound propagation: add β for each split constraint

$$y \geq 1 \cdot \hat{z}_1 + (-1) \cdot \hat{z}_2$$

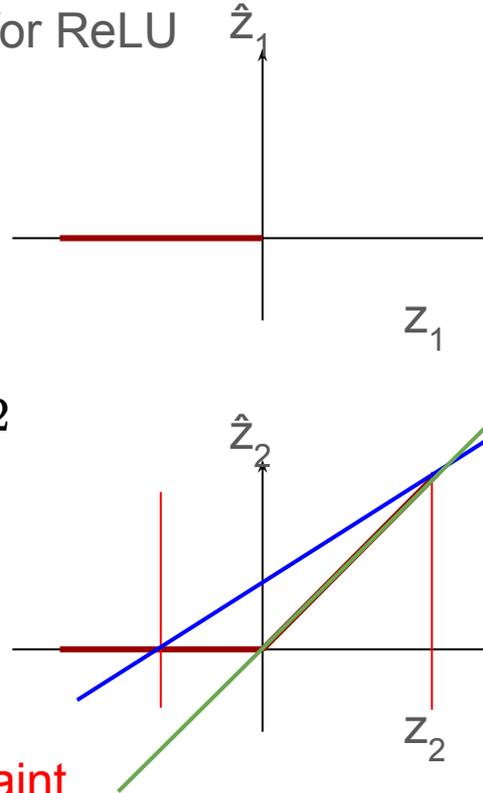
$$0 \leq \text{ReLU}(z_1) \leq 0$$

$$z_2 \leq \text{ReLU}(z_2) \leq \frac{2}{3}z_2 + 2$$



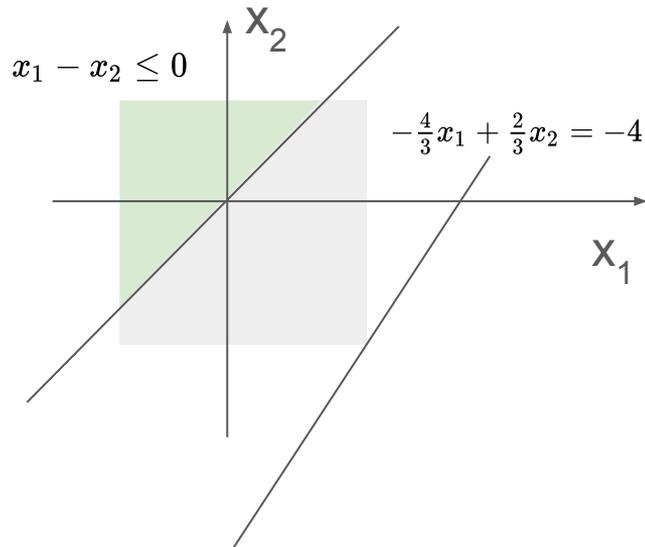
$$y \geq -\frac{2}{3}z_2 - 2 + \beta z_1$$

$\beta \geq 0$

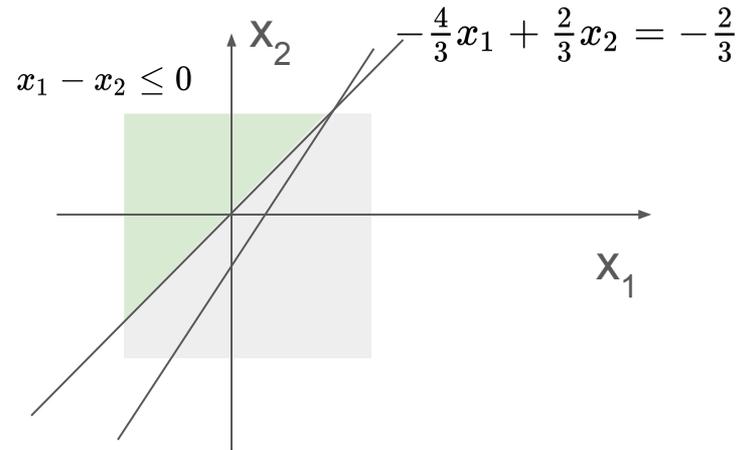


Review: Geometric interpretation

$$\min_{x_1, x_2} -\frac{4}{3}x_1 + \frac{2}{3}x_2 - 2$$



No constraint, obj = -6



With constraint, obj = $-\frac{8}{3}$, improved!

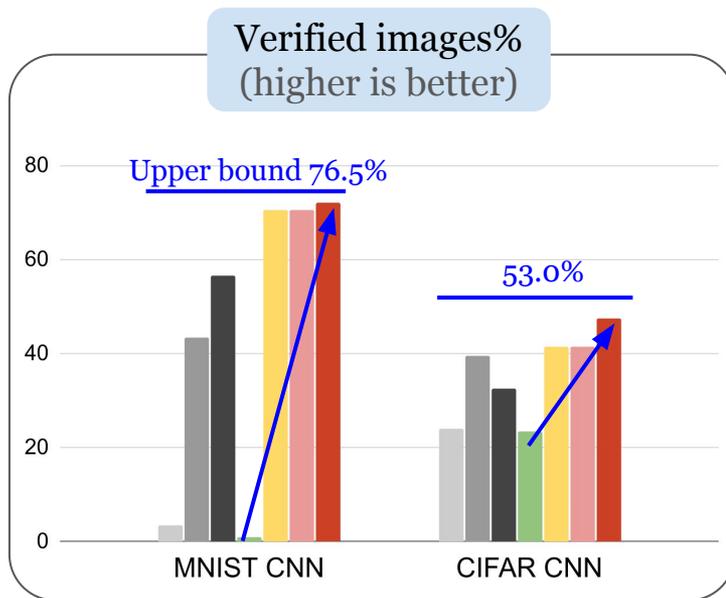
What we have learned so far about NN verification

- Verification as optimization problems
- Mixed Integer programming formulation for verifying ReLU networks
- Linear programming formulation
- Interval bound propagation (IBP)
- Linear bound propagation algorithm (CROWN)
- Bound optimization to improve tightness (α -CROWN)
- Branch and bound to further improve tightness (β -CROWN)

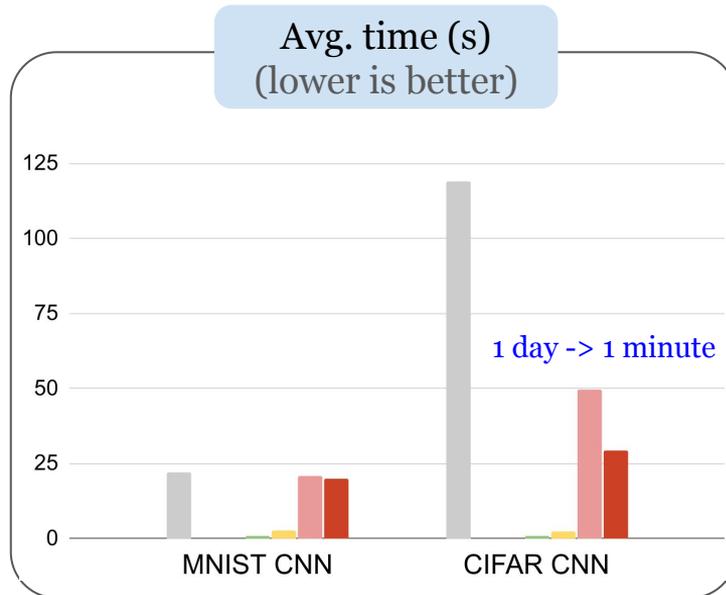
What are missing to solve practical NN verification problems?

Benchmarks: CROWN-family bound propagation algorithms

- Linear Programming (Salman et al. 2019)
- Semidefinite Programming (Dathathri et al. 2020)
- Integer Programming (Tjeng et al. 2017)
- CROWN
- α -CROWN
- β -CROWN
- GCP-CROWN



Model size: ~5k neurons

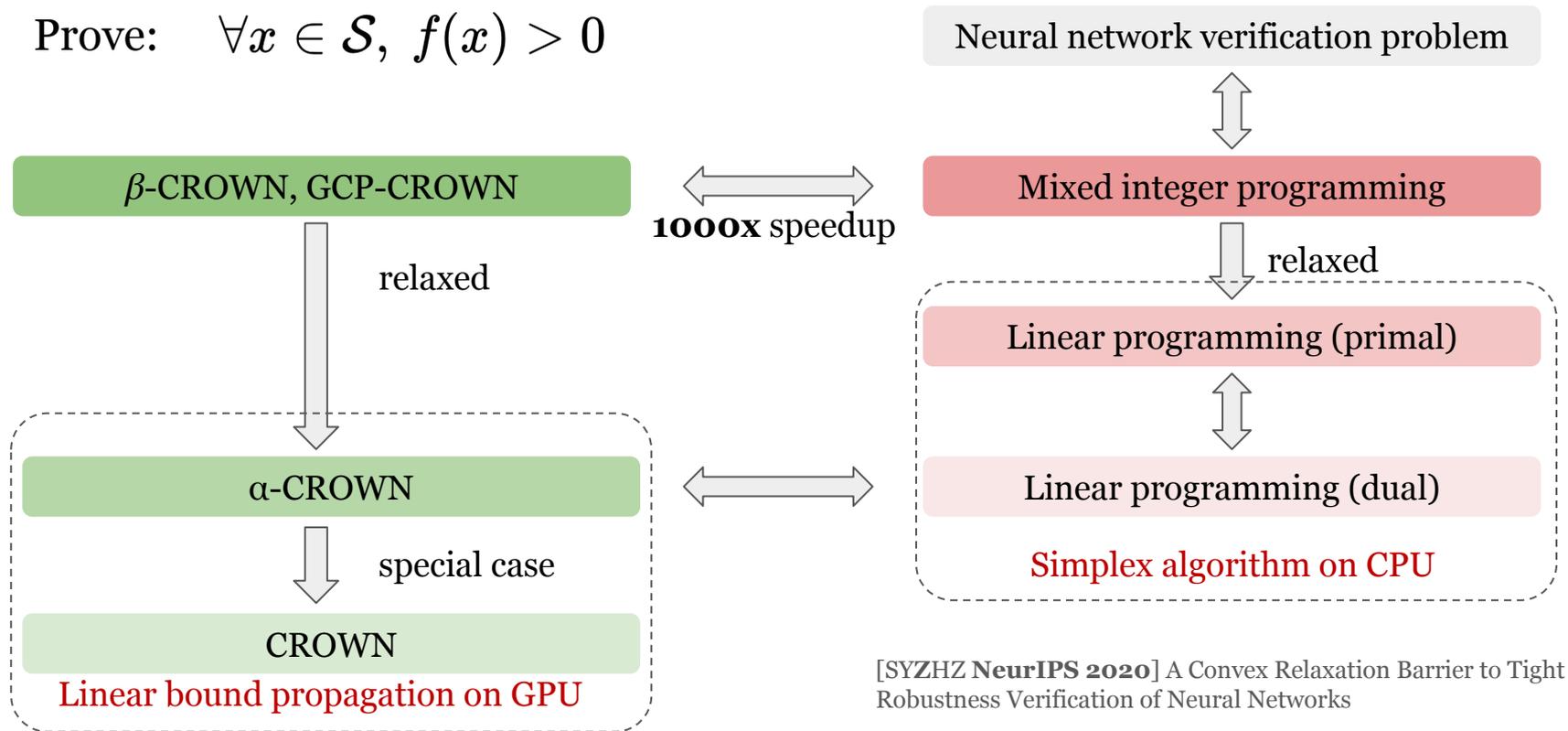


Integer programming and semidefinite programming **not plotted** (~1 day)

Key enablers: specialized bound propagation solver + GPU acceleration + BaB

Theoretical Connections: CROWN vs MIP/LP

Prove: $\forall x \in \mathcal{S}, f(x) > 0$



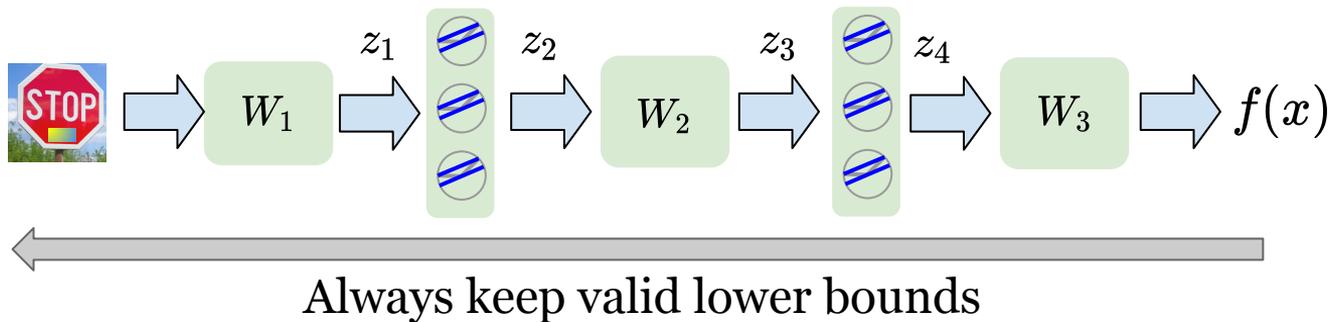
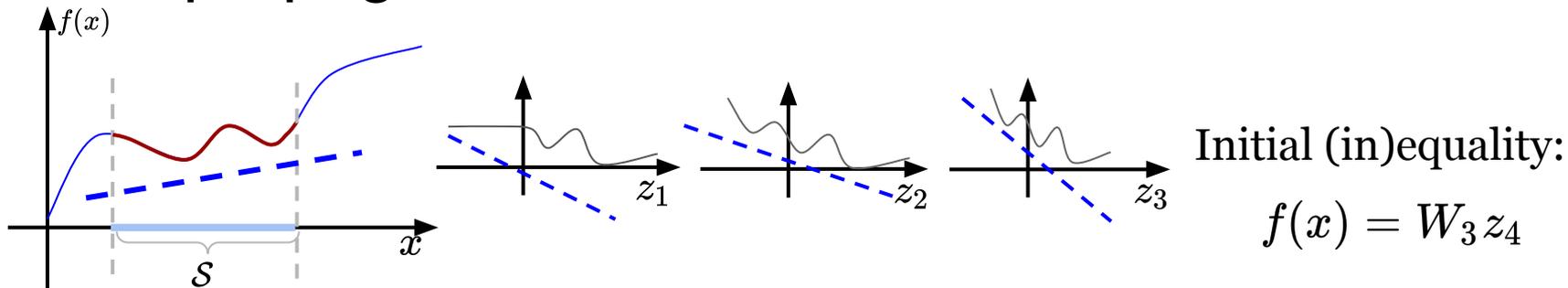
Several additional topics we will discuss today

Bound propagation on general computation graph

Falsification methods

Adversarial training and verification-friendly networks

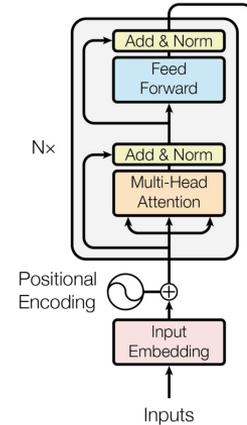
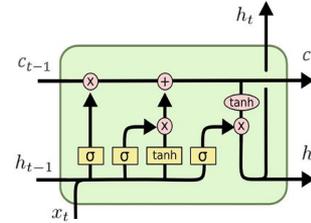
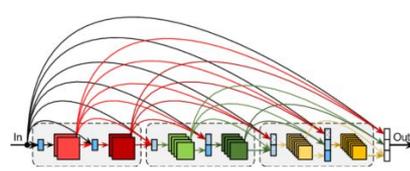
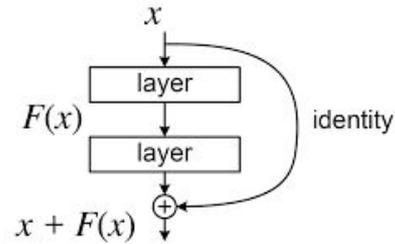
Bound propagation for feedforward neural networks



CROWN main theorem (simplified): $f(x) \geq a_{\text{CROWN}}^\top x + b_{\text{CROWN}} \quad \forall x \in \mathcal{S}$

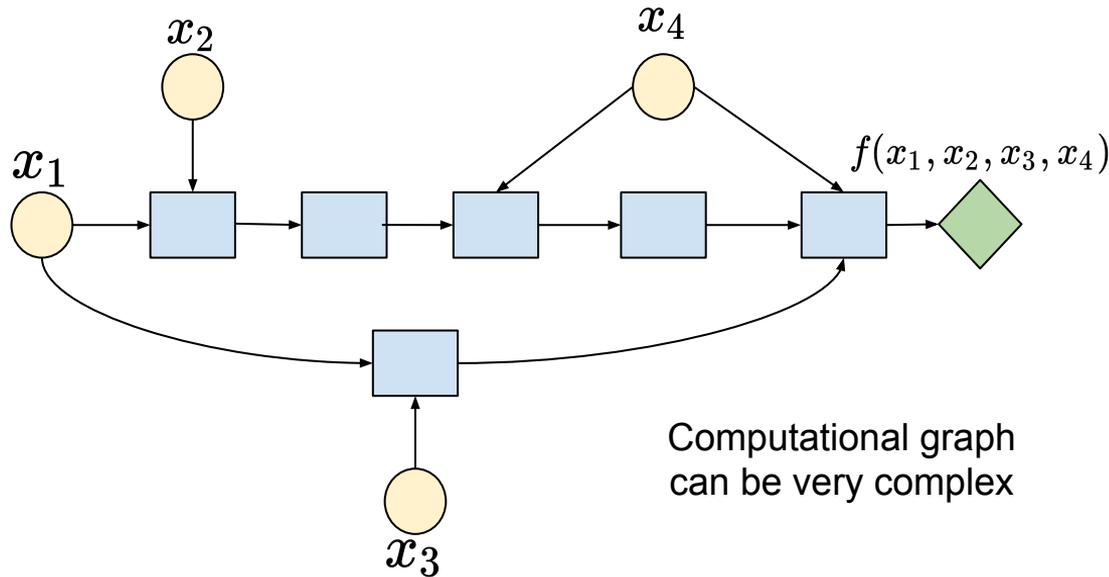
How about more complex networks?

Most modern neural networks have more than the “linear” feedforward structure



Verification on general computation graphs

The idea of bound propagation can be generalized to general computation graphs, as a graph algorithm

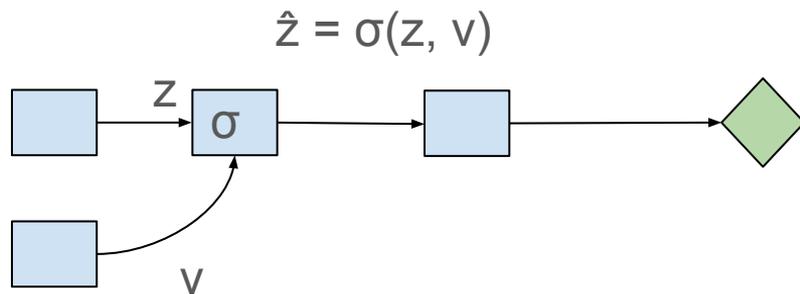


Computational graph
can be very complex

Requirement:
Each computation  can be
bounded using linear
hyperplanes w.r.t. its inputs

Bound propagation on computation graphs

One compute node can have multiple inputs



For example:

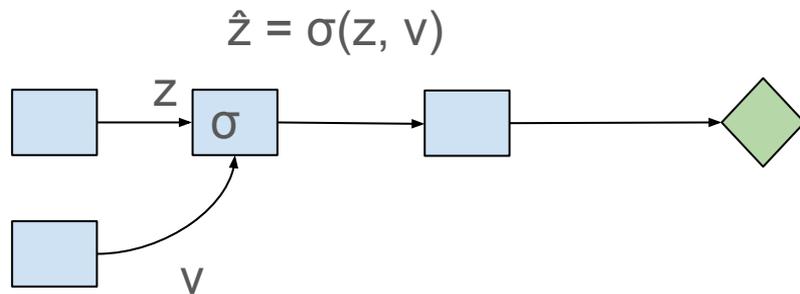
$$\hat{z} = z + v$$

$$\hat{z} = z \times v$$

Bound propagation on computation graphs

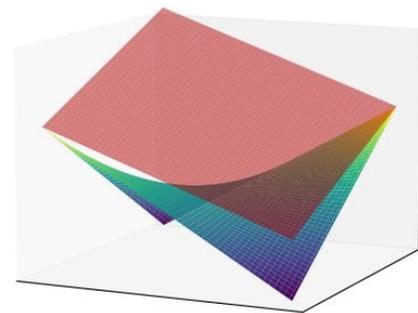
First step: bound σ using linear bounds of z and v

$$\underline{a}^\top z + \underline{b}^\top v + \underline{c} \leq \sigma(z, v) \leq \bar{a}^\top z + \bar{b}^\top v + \bar{c}$$



For $\hat{z} = z + v$, the function is already linear in z and v

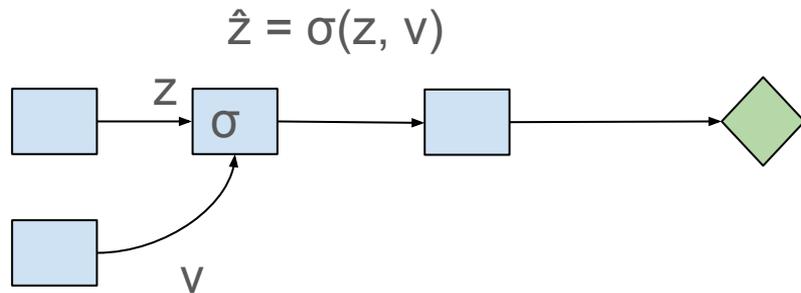
For $\hat{z} = z \times v$, need intermediate layer bounds for z and v



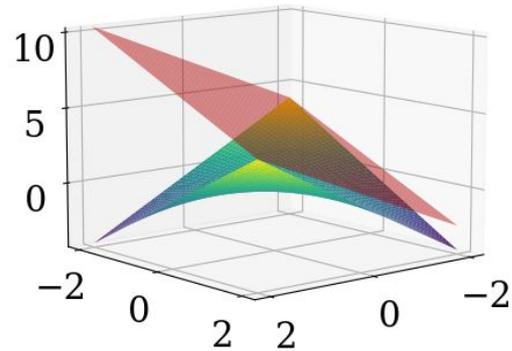
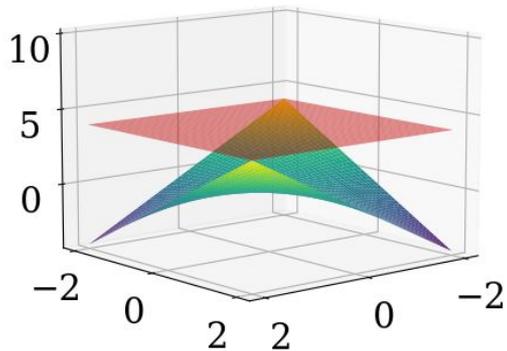
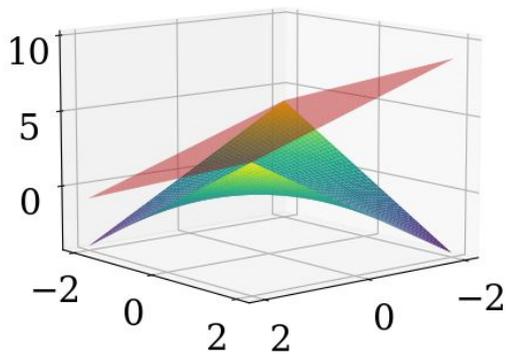
Bound propagation on computation graphs

For $\hat{z} = z \times v$, need intermediate layer bounds for z and v

$$\underline{a}^\top z + \underline{b}^\top v + \underline{c} \leq \sigma(z, v) \leq \bar{a}^\top z + \bar{b}^\top v + \bar{c}$$



These bounds can be optimized as well!

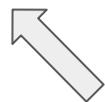


Bound propagation on computation graphs

$$\underline{a}^\top z + \underline{b}^\top v + \underline{c} \leq \sigma(z, v) \leq \bar{a}^\top z + \bar{b}^\top v + \bar{c}$$

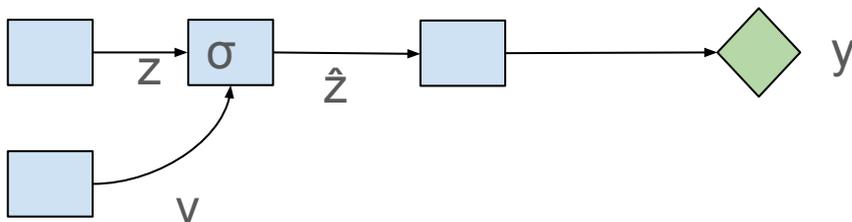


$$y \leq a_z^\top z + a_v^\top v + b'$$



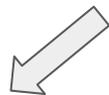
Choose lower or upper based on the sign of $a_{\hat{z}}$

$$y \leq a_{\hat{z}}^\top \hat{z} + b$$



Bound propagation on computation graphs

$$\underline{a}^\top z + \underline{b}^\top v + \underline{c} \leq \sigma(z, v) \leq \bar{a}^\top z + \bar{b}^\top v + \bar{c}$$



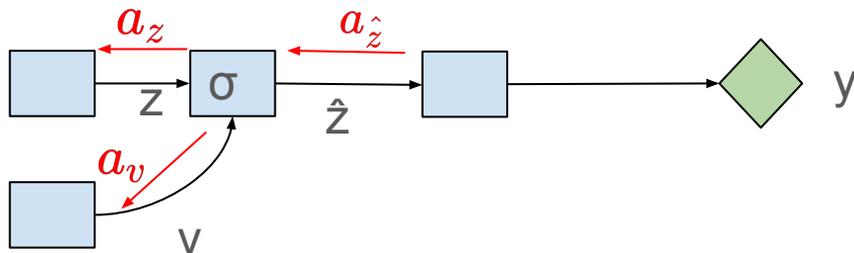
$$y \leq a_z^\top z + a_v^\top v + b'$$



Choose lower or upper based on the sign of $a_{\hat{z}}$

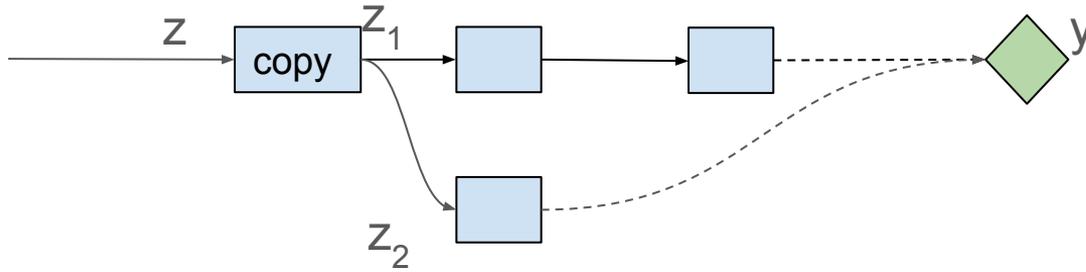
$$y \leq a_{\hat{z}}^\top \hat{z} + b$$

We essentially propagate those **coefficients** on the graph



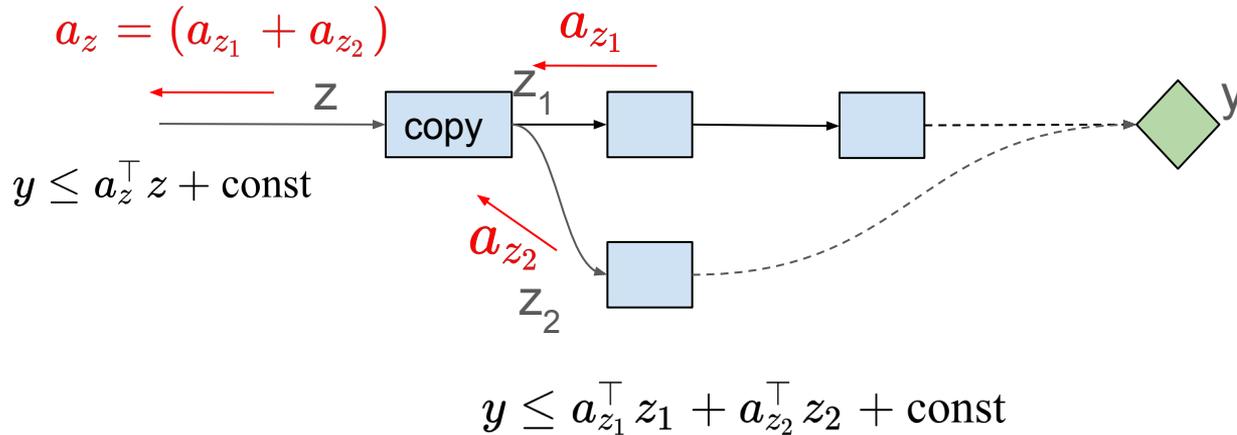
Bound propagation on computation graphs

One compute node's output can be used by multiple nodes



Bound propagation on computation graphs

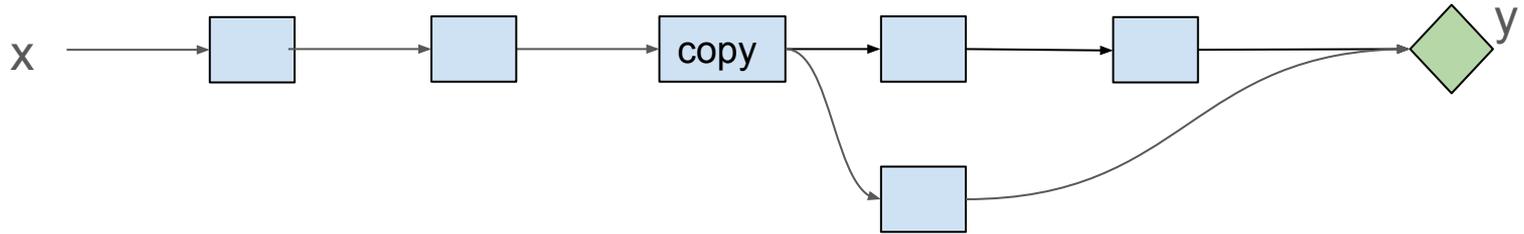
One compute node's output can be used by multiple nodes



Must wait until both coefficients become available

Bound propagation on computation graphs

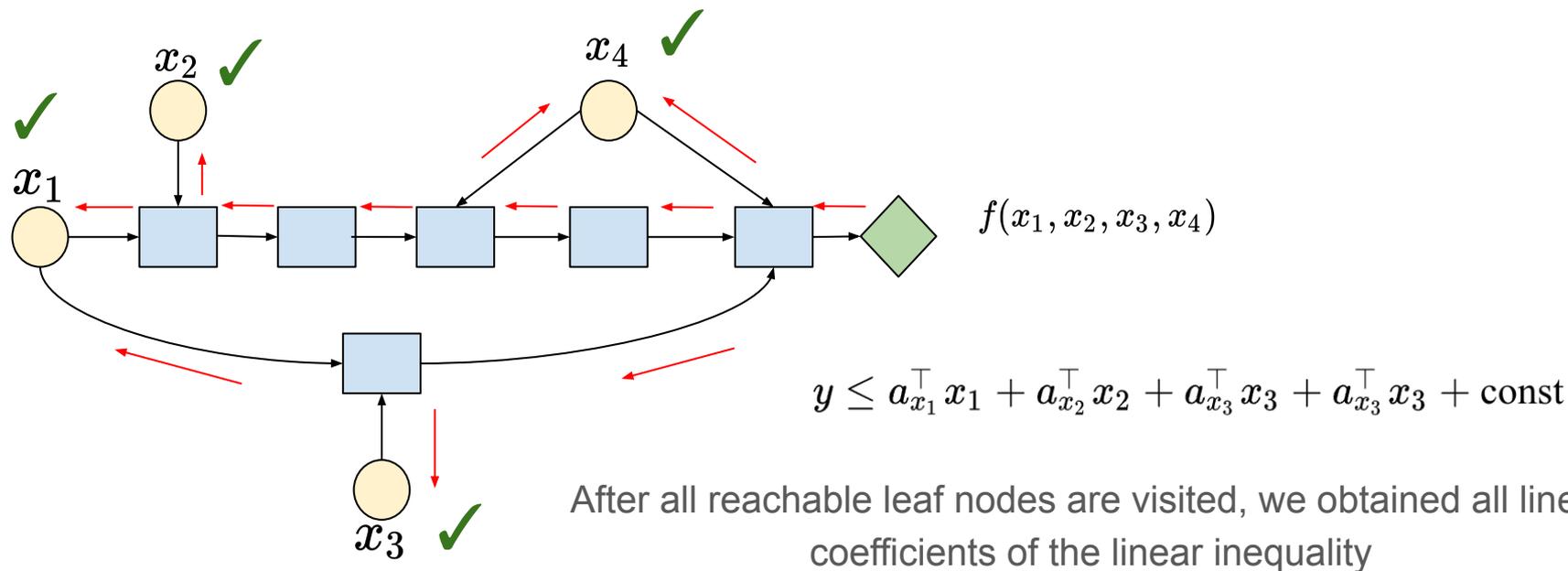
One compute node's output can be used by multiple nodes



Must wait until both coefficients become available

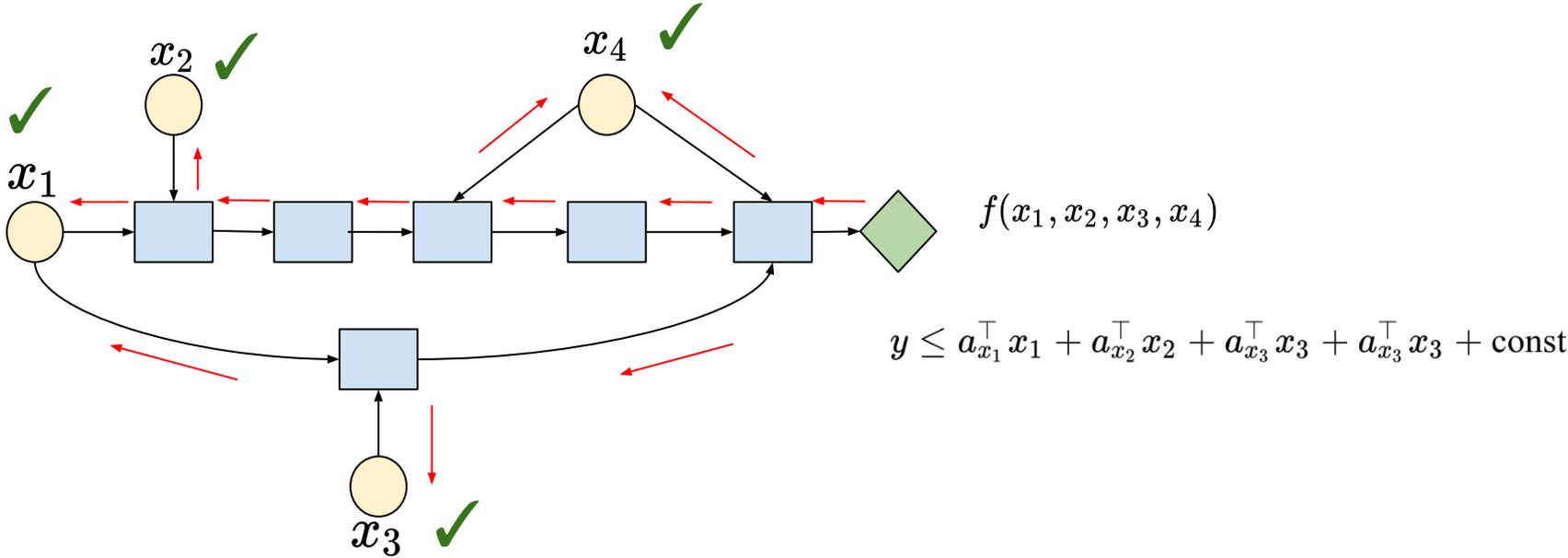
Bound propagation on computation graphs

When to stop? Reaching a leaf node.



Verification Beyond Neural Networks

Actually, CROWN can work on general computation graphs, not limited to neural networks!
Using CROWN for some novel applications that requires bounding is a great project idea!



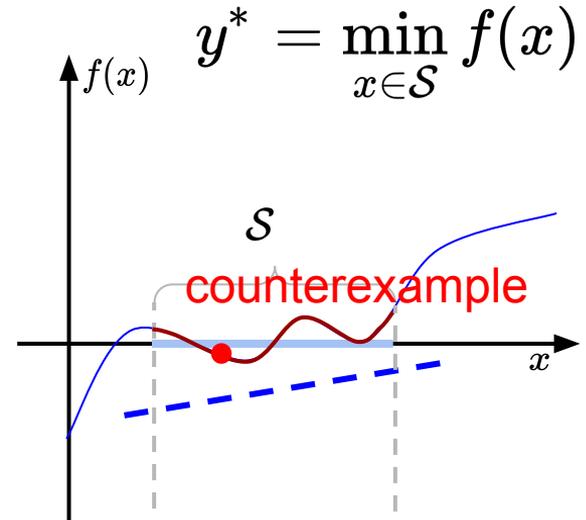
Falsification methods

Lower bounds are good for verification, but are not very helpful for finding a counterexample (feasible solution) so far. Back to our original problem:

$$\exists x \in S \wedge y \leq 0 \wedge y = f(x)$$

In some case, we want to find some x such that $y \leq 0$

- Use these counterexamples to fix model bugs
Similar to failed test cases in software engineering
- For large models, lower bound can be loose.
Verification is challenging and falsification has more hope
- Also called “adversarial attacks” in ML literature
Counterexample == adversarial example

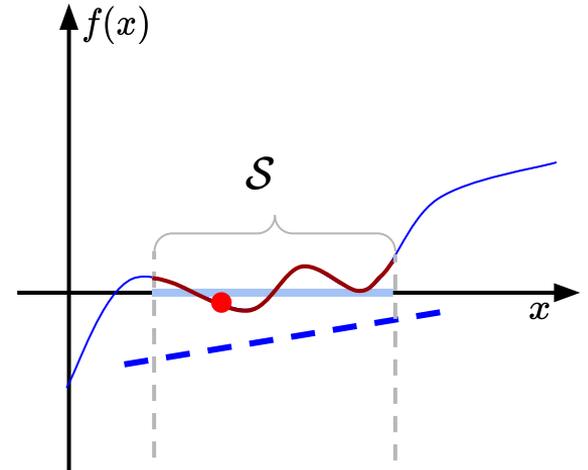


Falsification methods

To find a counterexample, we minimize the objective function using any method

- Randomly sample some $x \in S$ and check $f(x)$?
- **Gradient-based method**

$$y^* = \min_{x \in S} f(x)$$

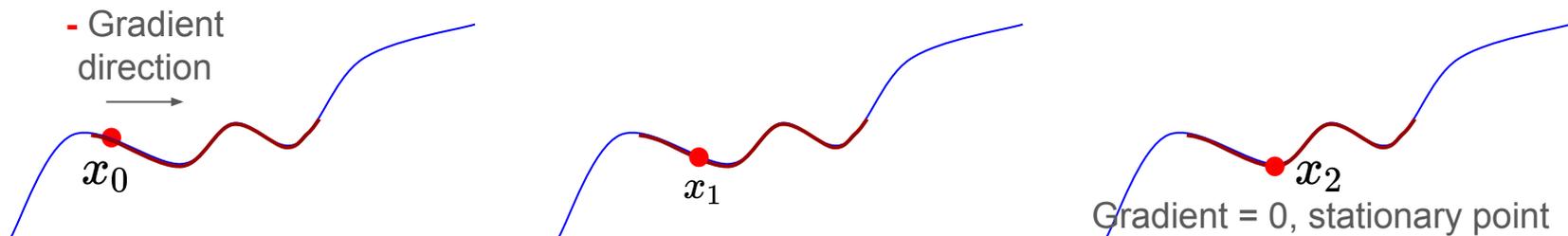


Projected gradient descent

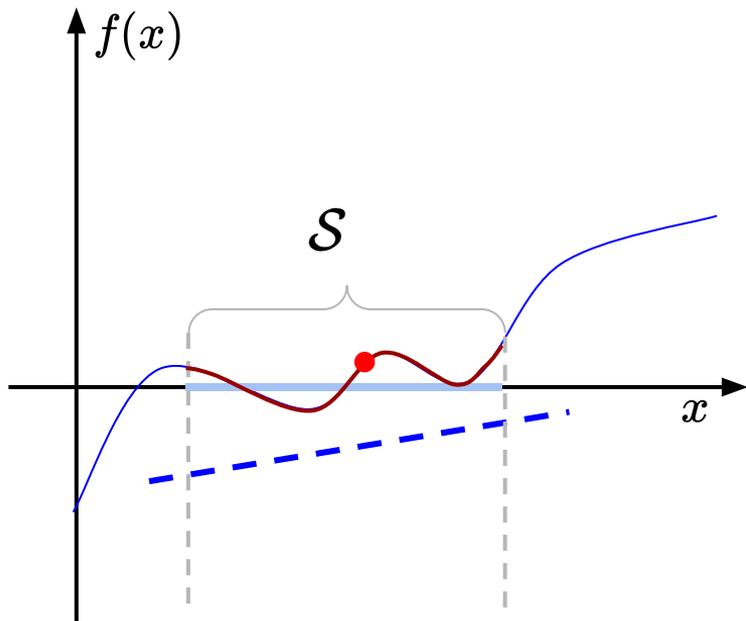
$$x_{t+1} \leftarrow \text{Proj}_{\mathcal{S}}(x_t - \eta \nabla f(x_t))$$

Follow the “downhill” to decrease $f(x)$

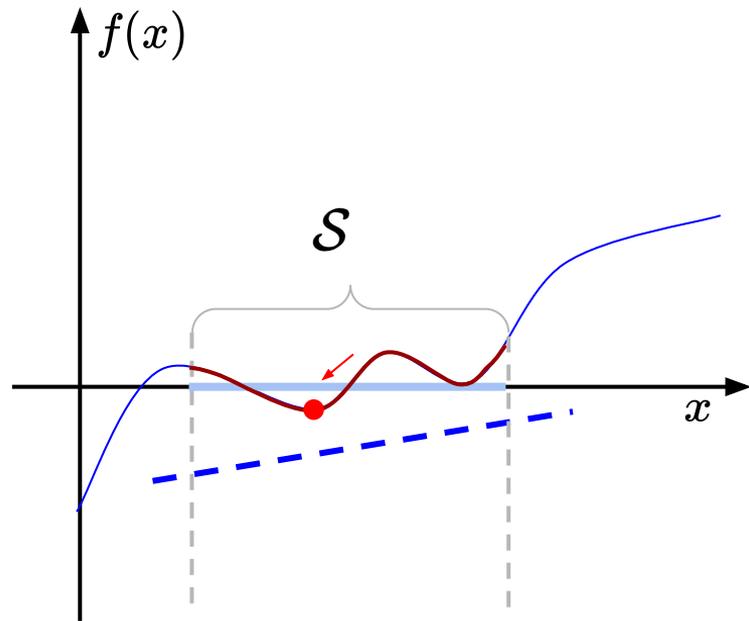
In the meanwhile, do not go outside of the set \mathcal{S}



Projected gradient descent

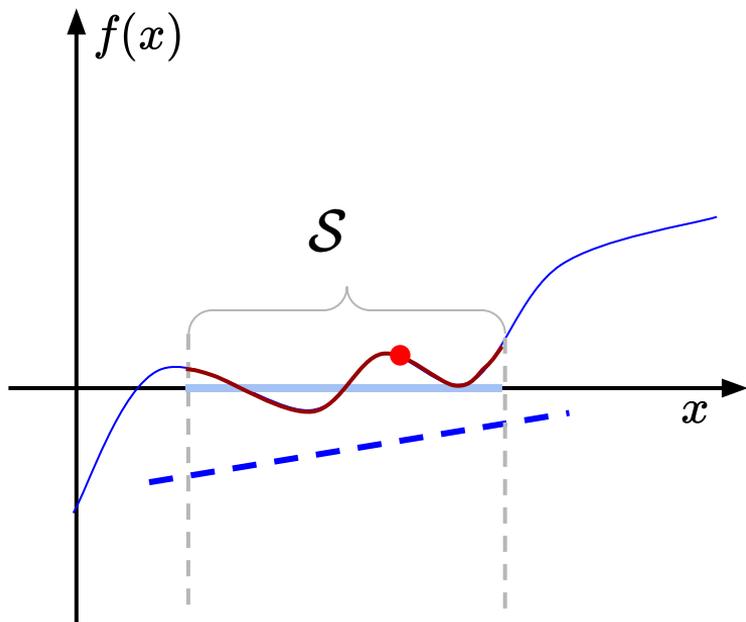


After a few iterations

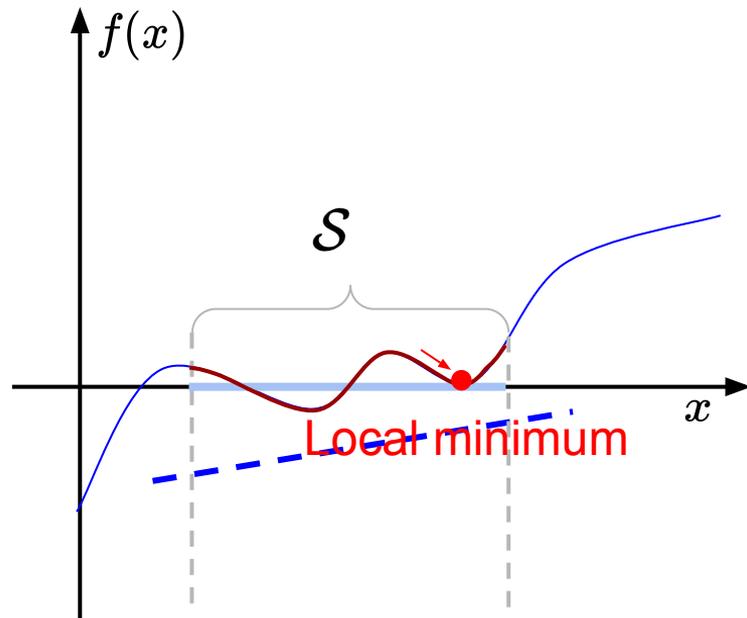
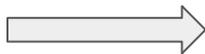


$f(x) < 0$, counterexample found!

Projected gradient descent may fail

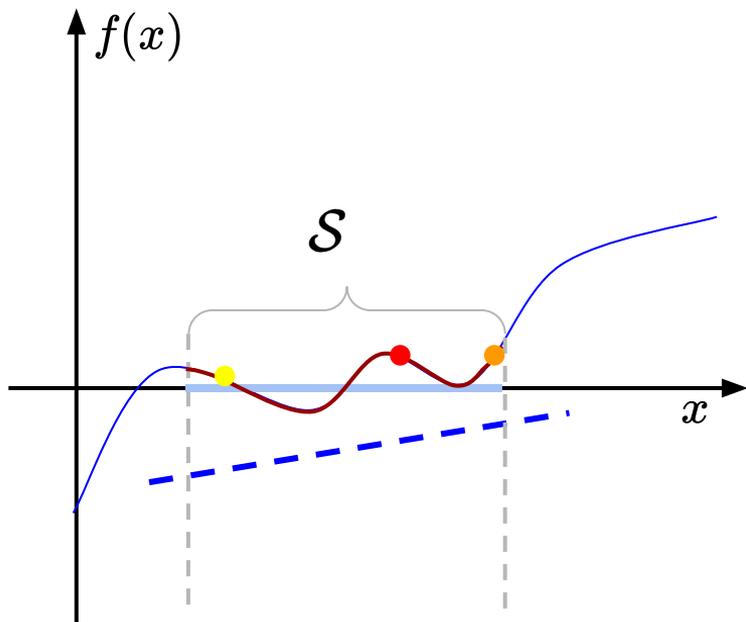


After a few iterations

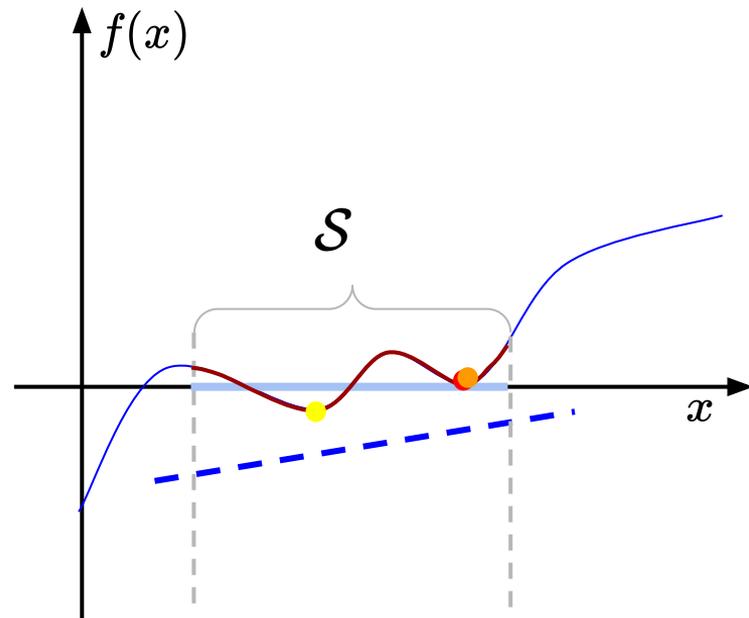


Get stuck! No further improvements possible

Projected gradient descent with random restarts



After a few iterations



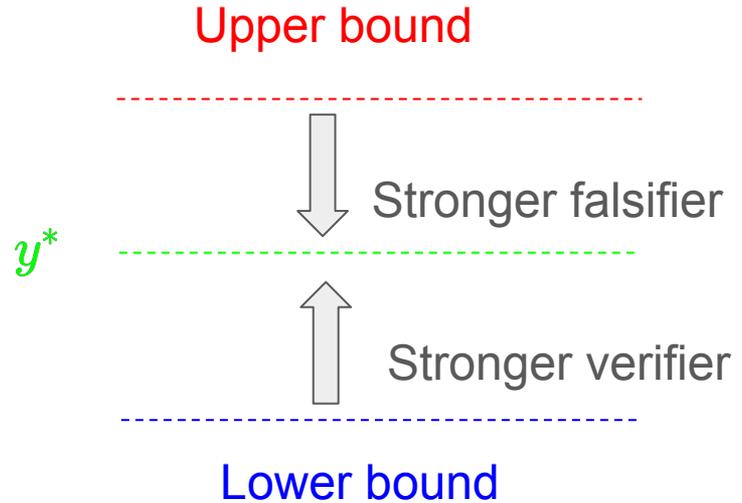
Try random starting points to increase chance

Falsification gives an upper bound of the verification problem

$$y^* = \min_{x \in \mathcal{S}} f(x)$$

Verification gives a lower bound of y^* :
Lower bound $\geq 0 \Rightarrow$ verified

Falsification gives an upper bound of y^* :
Upper bound $\leq 0 \Rightarrow$ falsified



Adversarial training

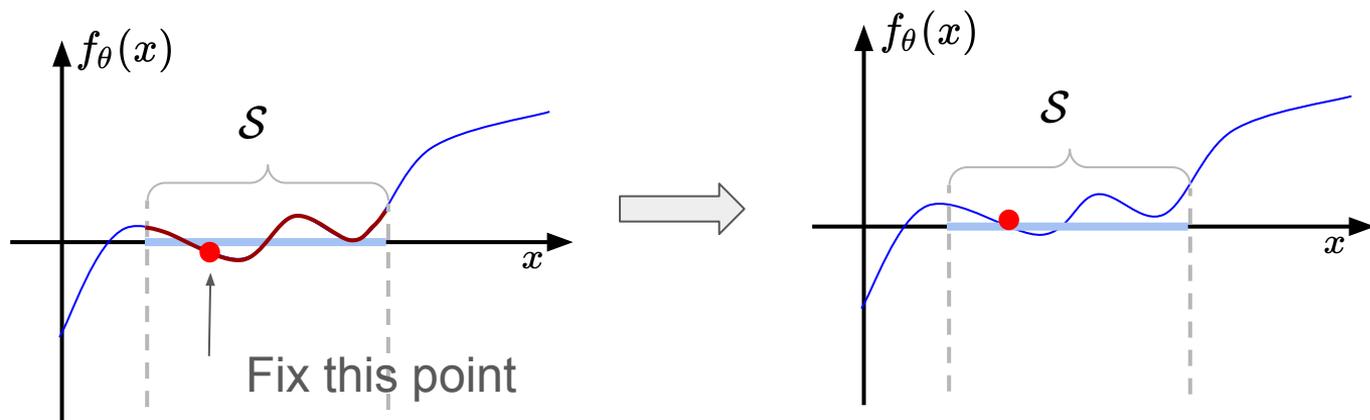
While not converged:

$x \leftarrow$ sample training data

$S \leftarrow$ a small neighborhood around x

$x_{\text{adv}} \leftarrow$ counterexample: $x_{\text{adv}} \in S \wedge f_{\theta}(x_{\text{adv}}) \leq 0$

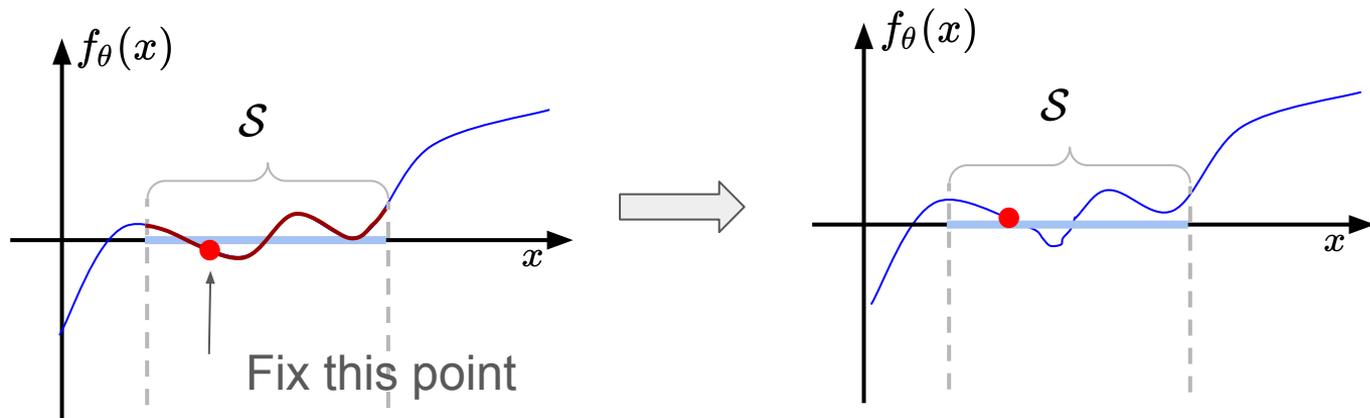
Update model parameter θ using gradient ascent to maximize $f_{\theta}(x_{\text{adv}})$



Adversarial training

Pros: (relatively) efficient to find counterexamples, less impact on model performance

Cons: just fixing a finite number of counterexamples may not lead to a truly verified model; usually challenging for verifiers



Verification-guided training

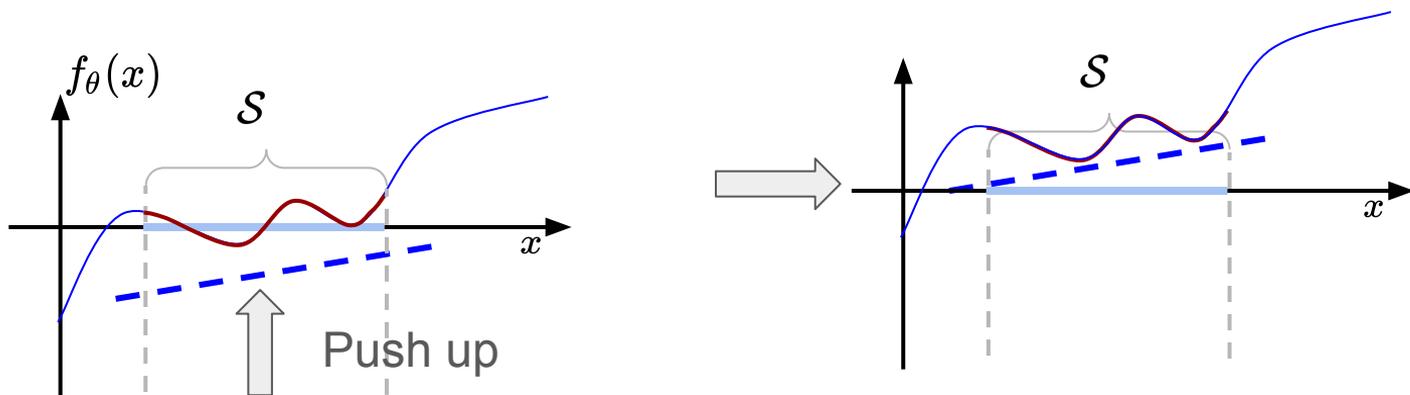
While not converged:

$x \leftarrow$ sample training data

$S \leftarrow$ a small neighborhood around x

$LB_{\theta}(S) \leftarrow$ lower bound of $f_{\theta}(x)$ in S

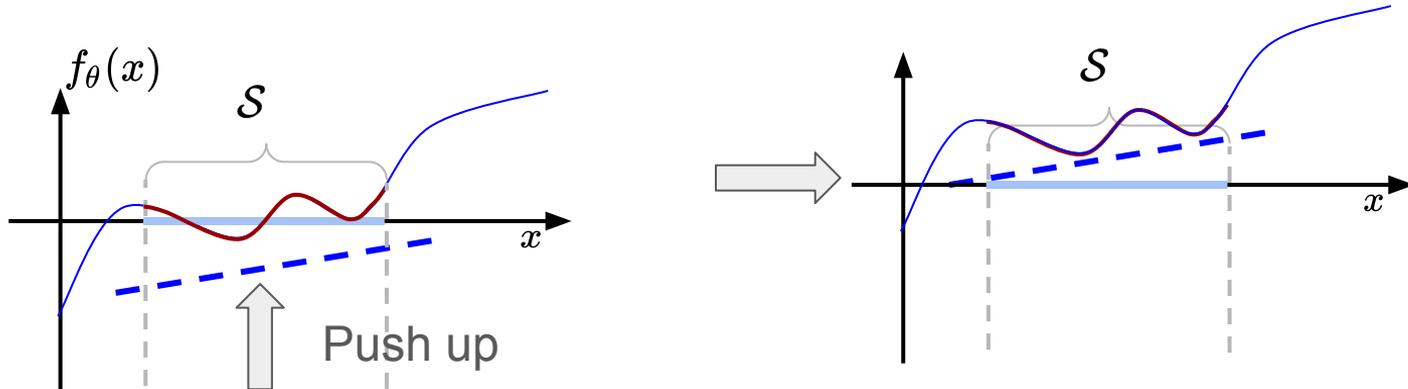
Update model parameter θ using gradient ascent to maximize $LB_{\theta}(S)$



Verification-guided training

Pros: Maximizing the lower bound guarantees $f(x) > 0$; models trained in this way are typically very easy to verify (verification friendly)

Cons: Calculating lower bounds is expensive; lower bounds can be too conservative, leading to poor model performance



Adversarial training vs Verification-guided training

CIFAR10: pixel-wise perturbation 8/255

Adversarial training: 90% clean accuracy, 70% accuracy under adversarial attack, ~0% verified accuracy

Verification-guided training: 55% clean accuracy, ~40% accuracy under adversarial attack, ~35% verified accuracy