CS781: A Quick Primer on Abstract Interpretation for Neural Networks

Supratik Chakraborty IIT Bombay

Notion of State in Neural Network



State: $(x_1, x_2, ..., x_{18})$ in R¹⁸

State Change in Feed-Forward Neural Network



 $(x'_{1}, x'_{2}, ..., x'_{i-1}, x'_{i}) = f_{i}(x_{1}, x_{2}, ..., x_{i-1}), \text{ for i in } \{3, ..., 18 \}$

State Change in Feed-Forward NN as a sequence of instrns



$$(x'_{1}, x'_{2}, x'_{3}) = f_{3}(x_{1}, x_{2});$$

$$(x''_{1}, x''_{2}, x''_{3}, x''_{4}) = f_{4}(x'_{1}, x'_{2}, x'_{3});$$

NN computation: a sequence of state transitions caused by seq of instructions

Proving Property of a FF NN



{Pre-condition on (x1, x2)}

$$(x'_{1}, x'_{2}, x'_{3}) = f_{3}(x_{1}, x_{2});$$

$$(x''_{1}, x''_{2}, x''_{3}, x''_{4}) = f_{4}(x'_{1}, x'_{2}, x'_{3});$$

{Post-condition on (x17, x18) }

NN Computation as a State Transition System



{Pre-condition on (x1, x2)}

$$\begin{array}{ll} (x'_{1}, x'_{2}, x'_{3}) & = f_{3}(x_{1}, x_{2}); \\ (x''_{1}, x''_{2}, x''_{3}, x''_{4}) & = f_{4}(x'_{1}, x'_{2}, x'_{3}); \end{array}$$

{Post-condition on (x17, x18) }

Dealing with State Space Size

- Infinite state space
 - Difficult to represent using state transition diagram
 - · Can we still do some reasoning?
- Solution: Use of abstraction
 - Naive view
 - Bunch sets of states together "intelligently"
 - Don't talk of individual states, talk of a representation of a set of states
 - Transitions between state set representations
 - Granularity of reasoning shifted
 - Extremely powerful general technique
 - Allows reasoning about large/infinite state spaces

Concrete states

Abstract states

A Generic View of Abstraction



- Every subset of concrete states mapped to unique abstract state
- Desirable to capture containment relations
- > Transitions between state sets (abstract states)

The Game Plan



Abstract analysis engine

The Game Plan



How do we choose the right abstraction? Is there a method beyond domain expertise? Can we learn from errors in abstraction to build better (refined) abstractions? Can refinement be automated?

Abstract analysis engine

TPIE

The Game Plan



Abstract analysis engine

Desirable Properties of Abstraction



- $ightarrow \operatorname{Sup}S_1 \subseteq S_2$: subsets of concrete states
 - Any behaviour starting from can also happen starting from
 - If , we want this monotonicity in behaviour in abstr state space too
 - Need ordering of abstract states, similar in spirit to

Structure of Concrete State Space



Structure of Abstract State Space

- [≻] Abstract lattice $A = (A, \sqsubseteq, \sqcup, \sqcap, \top, \bot)$
- - Monotone: $S_1 \subseteq S_2 \Rightarrow \alpha(S_1) \sqsubseteq \alpha(S_2)$ for all $S_1, S_2 \subseteq S$

• $\alpha(S) = \top, \quad \alpha(\emptyset) = \bot$

- - Monotone: $a_1 \sqsubseteq a_2 \Rightarrow \gamma(a_1) \subseteq \gamma(a_2)$ for all $a_1, a_2 \in \mathcal{A}$

•
$$\gamma(\top) = S$$
, $\gamma(\bot) = \emptyset$

A Simple Abstract Domain

- Simplest domain for analyzing numerical programs
- Represent values of each variable separately using intervals
- Example:



Represent values of inputs by intervals,

Compute values of hidden layer nodes and outputs as intervals

- > Abstract states: intervals of values of x, (ignore values of y)
 - [-10, 7]: { (x, y) | -10 <= x <= 7 }
 - (-∞, 20]: { (x, y) | x <= 20 }
 - relation: Inclusion of intervals
 [-10, 7] [-20, 9]
 - □ and □: union and intersection of intervals
 [-10, 9] □ [-20, 7] = [-20, 9]
 [-10, 9] □ [-20, 7] = [-10, 7]
 - \perp is empty interval of x
 - \top is (- ∞ , + ∞)

- > Abstract states: intervals of values of x, (ignore values of y)
 - [-10, 7]: { (x, y) | -10 <= x <= 7 }
 - (-∞, 20]: { (x, y) | x <= 20 }
 - relation: Inclusion of intervals
 [-10, 7] [-20, 9]
 - □ and □: union and intersection
 [-10, 9] □ [-20, 7] = [-20, 9]
 [-10, 9] □ [-20, 7] = [-10, 7]
 - \perp is empty interval of x
 - \top is (- ∞ , + ∞)

 $\begin{array}{l} \alpha(\ \{(1,\ 3),\ (2,\ 4),\ (5,\ 7)\}\)=[1,\ 5]\\ \alpha(\ \{(5,\ 7),\ (7,\ 6),\ (9,\ 10)\}\)=[5,\ 9]\\ \alpha(\ \{(5,\ 7)\}\)=[5,\ 5] \end{array}$



- > Abstract states: pairs of intervals (one for x, y)
 - · ([-10, 7], (-1, 20])
 - └ relation: Inclusion of intervals
 ([-10, 7], (-1, 20]) └ ([-20, 9], (-1, +∞))
 - \square and \square : union and intersection of intervals
 - $([-10, 9], (-1, 20]) \sqcap ([-20, 7], [3, +\infty)) = ([-10, 7], [3, 20])$
 - ([-10, 9], (-1, 20]) \sqcup ([-20, 7], [3,+ ∞)) = ([-20, 9],(-1,+ ∞))
 - \perp is empty interval of x and y
 - \top is ((- ∞ , + ∞), (- ∞ , + ∞))

Desirable Properties of α and γ

For all $S_1 \subseteq \mathcal{C}$ $S_1 \subseteq \gamma(\alpha(S_1))$

.



Desirable Properties of α and γ

$$S_1 \subseteq \gamma(\alpha(S_1)) \quad \text{forall} \quad S_1 \subseteq \mathcal{C}$$
$$\alpha(\gamma(a_1)) \sqsubseteq a_1 \quad \text{forall} \quad a_1 \in \mathcal{A}$$



Set of abstract states



α and γ form a Galois connection

Desirable Properties of α and γ

- - · Second (equivalent) view:

 $\alpha(S_1) \sqsubseteq a_1 \Leftrightarrow S_1 \subseteq \gamma(a_1) \text{ for all } S_1 \subseteq S, a_1 \in \mathcal{A}$



Homework Problem

Let $\alpha : (S) \to \mathcal{A}$ and $\gamma : \mathcal{A} \to (S)$ be monotone maps

Show that

 $S_1 \subseteq \gamma(\alpha(S_1))$ for all $S_1 \in (S)$ and $\alpha(\gamma(a_1)) \sqsubseteq a_1$ for all $a_1 \in \mathcal{A}$ holds if and only if

 $\alpha(S_1) \sqsubseteq a_1 \Leftrightarrow S_1 \subseteq \gamma(a_1) \text{ for all } S_1 \in (S) \text{ and } a_1 \in \mathcal{A}$