

Inductive synthesis for variable RL

Monday, 9 November 2020 23:10

- Rather than enforcing safety by examining and altering the structure of complex NN implementations → use black-box methods to synthesize deterministic programs, simpler, more interpretable approximations of NN.
- Forms the problem of NN verification in terms of:
 - ① counterexample and ② syntax-guided inductive synthesis procedure.
- The procedure searches for ① a deterministic program and ② an inductive invariant over an infinite state transition system that represents a specification of an application's control logic.
- The synthesis procedure treats the NN as an oracle extracting a deterministic program P intended to approximate the policy actions implemented by the network P also serves as a safety shield that operates in tandem with the network
- Pros: generalizes to infinite state systems with a continuous state space.
- A verification toolchain for ensuring that the control policies learned by RL-based neural network are safe.
- Verification Procedure:
 - To verify NN control policy π_w w.r.t environment context C we have to synthesize a deterministic policy program P .
 - P should:
 1. broadly resembles its neural oracle
 2. satisfies a desired safety property when deployed in C .
 - Safety proof of $C[P]$ is easier than of $C[\pi_w]$

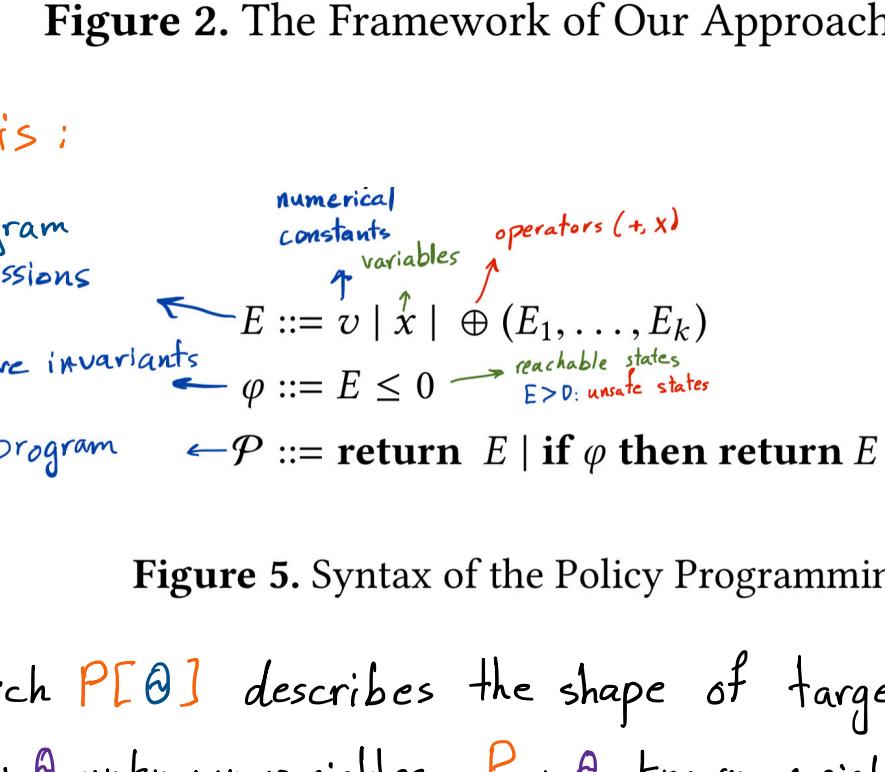


Figure 2. The Framework of Our Approach.

1. Synthesis:

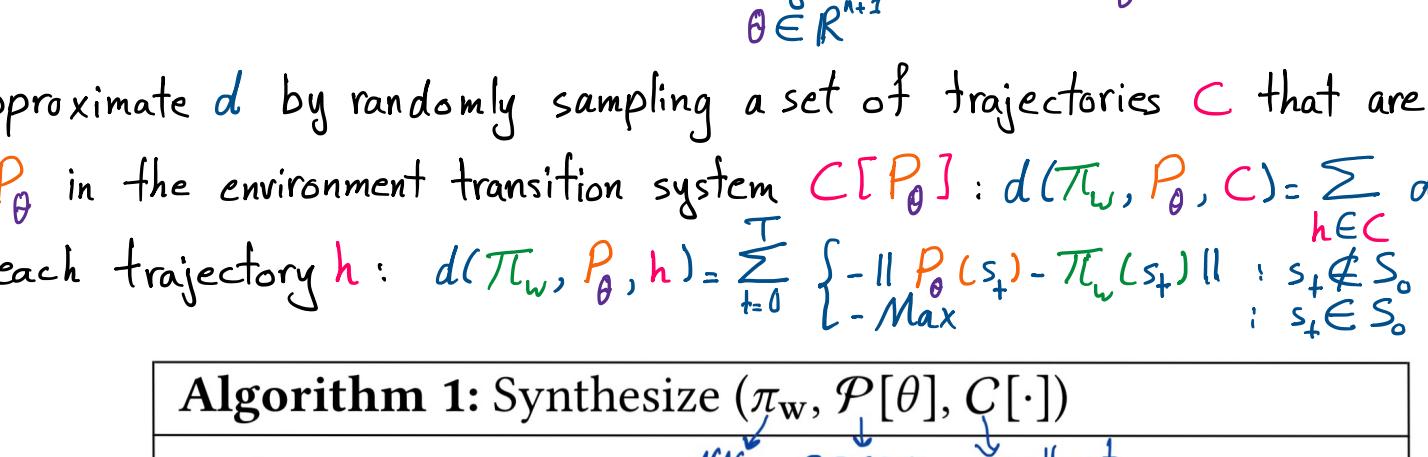


Figure 5. Syntax of the Policy Programming Language.

- The sketch $P[\theta]$ describes the shape of target policy using the grammar. $P[\theta] : \theta$ unknown variables, $P_\theta : \theta$ known variables (P_θ is the synthesized program). ex: $P[\theta](x) := \text{return } \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n + \theta_{n+1}$
- the goal of the algorithm is to find θ that maximize the likelihood that P_θ resembles the neural oracle π_w , while still being safe w.r.t. C :

$$\theta^* = \arg \max_{\theta \in \mathbb{R}^{n+1}} d(\pi_w, P_\theta, C)$$

- we approximate d by randomly sampling a set of trajectories C that are encountered by P_θ in the environment transition system $C[P_\theta]$: $d(\pi_w, P_\theta, C) = \sum_{h \in C} d(\pi_w, P_\theta, h)$
- for each trajectory h : $d(\pi_w, P_\theta, h) = \sum_{t=0}^T \frac{\|P_\theta(s_t) - \pi_w(s_t)\|}{L - \text{Max}} : s_t \notin S_0 \quad ; \quad s_t \in S_0$

Algorithm 1: Synthesize $(\pi_w, P[\theta], C[\cdot])$

1	$\theta \leftarrow 0;$	\downarrow	\downarrow	\downarrow
2	do	\downarrow	\downarrow	\downarrow
3	Sample δ from a zero mean Gaussian vector;	\downarrow	\downarrow	\downarrow
4	Sample a set of trajectories C_1 using $C[P_{\theta+\nu\delta}]$;	\downarrow	\downarrow	\downarrow
5	Sample a set of trajectories C_2 using $C[P_{\theta-\nu\delta}]$;	\downarrow	\downarrow	\downarrow
6	$\theta \leftarrow \theta + \alpha \frac{d(\pi_w, P_{\theta+\nu\delta}, C_1) - d(\pi_w, P_{\theta-\nu\delta}, C_2)}{\nu} \delta$;	\downarrow	\downarrow	\downarrow
7	while θ is not converged;	\downarrow	\downarrow	\downarrow
8	return P_θ	\downarrow	\downarrow	\downarrow

δ : Gaussian noise
 ν : small real number
 P_θ : perturbed policy program

2. Verification:

Verification algorithm learns an inductive variant φ over the transition relation $T_w[P]$ formally proving that all possible system behaviors are encapsulated in φ and φ has to be disjoint with all unsafe states S_u .

by defining a function $E: \mathbb{R}^n \rightarrow \mathbb{R}$ that serves as a barrier between reachable states $E \leq 0$ and unsafe states $E > 0$

$$\varphi[c] := E[c](X) \leq 0 ; c \text{ parameters to be learned}$$

E is found using SMT-solvers: $E[c](X) = \sum_{i=0}^P c_i b_i(X)$ one term

- The algorithm Counter Example-Guided Inductive Synthesis (CEGIS):

Algorithm 2: CEGIS $(\pi_w, P[\theta], C[\cdot])$

1	policies $\leftarrow \emptyset$;	\downarrow	\downarrow	\downarrow
2	covers $\leftarrow \emptyset$;	\downarrow	\downarrow	\downarrow
3	while $C.S_0 \not\subseteq \text{covers}$ do	\downarrow	\downarrow	\downarrow
4	search s_0 such that $s_0 \in C.S_0 \wedge s_0 \notin \text{covers}$;	\downarrow	\downarrow	\downarrow
5	$r^* \leftarrow \text{Diameter}(C.S_0)$;	\downarrow	\downarrow	\downarrow
6	do	\downarrow	\downarrow	\downarrow
7	$\phi_{\text{bound}} \leftarrow \{s \mid s \in [s_0 - r^*, s_0 + r^*]\}$;	\downarrow	\downarrow	\downarrow
8	$\tilde{C} \leftarrow C \text{ where } \tilde{C}.S_0 = (C.S_0 \cap \phi_{\text{bound}})$;	\downarrow	\downarrow	\downarrow
9	$\theta \leftarrow \text{Synthesize}(\pi_w, P[\theta], \tilde{C}[\cdot])$;	\downarrow	\downarrow	\downarrow
10	$\varphi \leftarrow \text{Verify}(\tilde{C}[P_\theta])$;	\downarrow	\downarrow	\downarrow
11	if φ is False then	\downarrow	\downarrow	\downarrow
12	$r^* \leftarrow r^*/2$;	\downarrow	\downarrow	\downarrow
13	else	\downarrow	\downarrow	\downarrow
14	covers $\leftarrow \text{covers} \cup \{s \mid \varphi(s)\}$;	\downarrow	\downarrow	\downarrow
15	policies $\leftarrow \text{policies} \cup (P_\theta, \varphi)$;	\downarrow	\downarrow	\downarrow
16	break ;	\downarrow	\downarrow	\downarrow
17	while True;	\downarrow	\downarrow	\downarrow
18	end	\downarrow	\downarrow	\downarrow
19	return policies	\downarrow	\downarrow	\downarrow

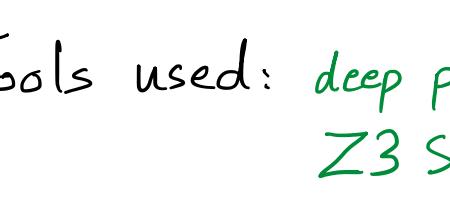
δ : Gaussian noise
 ν : small real number
 P_θ : perturbed policy program

the algorithm is sound but ① it may not terminate as r^* can become small
 ② no restriction on the size of potential counter example

3. Shielding:

Problem: the network may exhibit behaviors not captured by the simpler deterministic program.

Solution: using the synthesized policy program and its inductive invariant as a shield.



Algorithm 3: Shield $(s, C[\pi_w], P, \varphi)$

1 Predict s' such that $(s, s') \in C[\pi_w].T_w(s)$; → predict the next state

2 if $\varphi(s')$ then return $\pi_w(s)$;

3 else return $P(s)$; → shield to ensure safety

($P(s)$ is proved safe in Algorithm 2)

Figure 4. The Framework of Neural Network Shielding.

- Tools used: deep policy gradient for RL training

Z3 SMT-solver to check convergence of CEGIS loop.

Mosek constraint solver to generate inductive invariants of synthesized program from a sketch.

- Synthesizing the shield is better (faster) than retraining NN for a new environment.