

Programmatically interpretable RL

Wednesday, 11 November 2020 22:48

- PIRL: is based on the idea of learning policies that are represented in a human-readable language.
- The problem: traditional RL + a (policy) sketch that syntactically defines a set of programmatic policies in PIRL language
- key benefits:
 1. Interpretability.
 2. Implicitly encode the learner's inductive bias, which will be used for generalization.
 3. Effective pruning of undesired policies, to make the search for a good policy more efficient.
 4. Allow to use symbolic program verification techniques.

- key challenge: the space of policies permitted in an instance can be vast and non-smooth; making optimization extremely challenging.

- Solution: Neurally Directed Program Synthesis (NDPS)

uses DRL to train NN → use the NN to direct a local search

over programmatic policies then;

interesting inputs → minimization the distance between

① the program output and ② the output of the oracle (NN policy).

- Partially observable Markov decision process (POMDP):

$$M = \{S, A, O, T, Z(\cdot|s), r, \text{Init}, \gamma\}$$

$Z(\cdot|s)$ - the probability that the agent makes an observation o at state s .

Init - the initial distribution over environment states.

A history $h = o_0, a_0, o_1, a_1, \dots, o_{k-1}, a_k$

\mathcal{H}_M a set of histories.

- A programmatic language for policies:

- atoms: observations and actions as well as auxiliary integers and reals generated during computation.

- Two types of data: ① atoms and ② sequences of atoms (e.g. histories).

$$E ::= c \mid x \mid \oplus(E_1, \dots, E_k) \mid \text{peek}(x, i) \mid \text{fold}((\lambda x_1, x_2. E_1), \alpha)$$

most recent observation $\text{peek}(x, -1)$

$$\text{fold}(f, [e_1, \dots, e_k], e) = f(e_k, f(e_{k-1}, \dots, f(e_1, e)))$$

the system identifies a set of expressions whose

① inputs are histories ② outputs are actions

this expressions are known as programs (programmatic policies).

- sketch: is a grammar of expressions over atoms and sequences of atoms,

places restrictions on the kind of basic operators that

will be considered during policy search.

the set of programs permitted by a sketch S is denoted $\llbracket S \rrbracket$

- The PIRL problem: find a program $e^* \in \llbracket S \rrbracket$ with optimal reward

$$e^* = \underset{e \in \llbracket S \rrbracket}{\operatorname{argmax}} R(e)$$

example ; PID controller.

sketch:

$$\begin{aligned} P &::= \text{peek}((\epsilon - h_i), -1) \\ I &::= \text{fold}(+, \epsilon - h_i) \\ D &::= \text{peek}(h_i, -2) - \text{peek}(h_i, -1) \\ C &::= c_1 * P + c_2 * I + c_3 * D \rightarrow \text{PID controllers} \\ B &::= c_0 + c_1 * \text{peek}(h_1, -1) + \dots \rightarrow \text{boolean} \\ &\quad \dots + c_k * \text{peek}(h_m, -1) > 0 \mid \text{operators} \\ E &::= C \mid \text{if } B \text{ then } E_1 \text{ else } E_2 \leftarrow \text{programs permitted by the sketch} \end{aligned}$$

- Neurally Directed Program search (NDPS):

Algorithm 1 Neurally Directed Program Search

Input: POMDP M , neural policy e_N , sketch S

$\mathcal{H} \leftarrow \text{create_histories}(e_N, M)$

$e \leftarrow \text{initialize}(e_N, \mathcal{H}, M, S)$

$R \leftarrow \text{collect_reward}(e, M)$

repeat

$(e', R') \leftarrow (e, R)$

$\mathcal{H} \leftarrow \text{update_histories}(e, e_N, M, \mathcal{H})$

$E \leftarrow \text{neighborhood_pool}(e)$

$e \leftarrow \arg \min_{e' \in E} \sum_{h \in \mathcal{H}} \|e'(h) - e_N(h)\|$

$R \leftarrow \text{collect_reward}(e, M)$

until $R' \geq R$

Output: e'

the main intuition is that, the distance from e_N is

a simpler objective than reward function.

Problem: the agent may encounter histories that are not possible under e_N

or any of the programs encountered in previous iterations.

Solution: Input augmentation or periodic updates to the set \mathcal{H} .

→ every while we sample a set of additional histories

by simulating the current programmatic policy, and

add these samples to \mathcal{H} .

- The optimization: two steps:

1. Enumerating a set of program templates (numerically parametrized programs) that are ① strictly similar to e and ② permitted by the sketch S , giving priority to shorter templates.

2. Use Bayesian optimization to find the optimal parameters for the enumerated templates.

→ Bayesian optimization is better than SMT solving here.