Cicaps 34th International Conference on Automated Planning and Scheduling

Epistemic Exploration for Generalizable Planning and Learning in Non-Stationary Settings

Rushang Karia^{*}, <u>Pulkit Verma</u>^{*}, Alberto Speranzon, Siddharth Srivastava



Rushang



Pulkit

Alberto



Siddharth



Motivation: Learn Environment Dynamics Efficiently

- Stream of tasks not known in advance
- Unknown, non-stationary environment dynamics (action models)
- Limited budget per task

Actively Conduct Experiments to Learn Environment Dynamics Efficiently



Problem Setting

Setting

- A stream of tasks M_1, \ldots, M_n with
 - different initial states and goals
 - different state/action spaces
- A simulator whose transition function changes in an arbitrary fashion at unknown intervals.
- Reward for reaching a goal is +1 and is 0 otherwise.

Objective

- Maximize #tasks completed within a fixed budget.
- Need to adapt fast (minimize adaptive delay), and compute good solutions (minimize regret).



Image Source: J. Balloch et al., NovGrid: A Flexible Grid World for Evaluating Agent Response to Novelty, AAAI'22 Spring Symposium on Designing AI for Open Worlds

Research Questions

Q1

Can we develop a method to conduct experiments in the environment actively?

Q1

Can we develop a method to conduct experiments in the environment actively? **Q2**

Can the exploration be more deliberative to learn more about the environment?

Q1

Can we develop a method to conduct experiments in the environment actively? Can the exploration be more deliberative to learn more about the environment?

Q2

Q3

Would it be useful to learn a STRIPSlike model to direct exploration?

Is it even feasible?

Our Solution: Continual Learning and Planning (CLaP)

Challenges in the Learning and Planning Paradigm

- 1. How do we generate sample-efficient strategies for learning?
 - Need to explore the state space
 - Need to be systematic

- 2. Is learning a full model worth it?
 - Might learn irrelevant actions wrt the current task.
 - Non-stationarity might waster a lot of the computational effort.

Continual Learning and Planning (CLaP)

Algorithm 1: Continual Learning and Planning **Input** : RMDP M, Simulator Δ , Simulator Budget $\Delta_{\mathcal{S}}$, Learned Model \mathcal{M}^{\uparrow} , Horizon *H*, Sampling Count η , Threshold θ , Failure Threshold β **Output:** \mathcal{M}^{\uparrow} 1 $s \leftarrow s_0; h \leftarrow 0; f \leftarrow 0$ 2 $\pi \leftarrow \text{modelBasedSolver}(S, A, s_0, g, \mathcal{M}^{\uparrow}, R, \gamma, H)$ 3 while $|\Delta| < \Delta_S$ do if $f > \beta$ or unreachableGoal $(s_0, g, \mathcal{M}^{\uparrow}, \pi)$ then explore($\mathcal{M}^{\uparrow}, \Delta$) 5 if *needsLearning*(\mathcal{M}) then 6 $\mathcal{M}^{\uparrow} \leftarrow \text{learnModel}(\Delta, \mathcal{M}^{\uparrow})$ 7 $\pi \leftarrow \text{modelBasedSolver}(S, A, s_0, g, \mathcal{M}^{\uparrow}, R, \gamma, H)$ 8 $a \leftarrow \pi(s)$ 9 $s' \leftarrow \Delta(s, a)$ 10 $h \leftarrow h + 1$ 11 if $(s, a, s') \rightleftharpoons \mathcal{M}^{\uparrow}$ then 12 $\mathcal{M} \leftarrow \text{goodnessOfFitTest}(s, a, s', \Delta, \mathcal{M}^{\uparrow}, \theta, Freq)$ 13 else 14 $\mathcal{M}^{\uparrow} \leftarrow \text{addInconsistentPredicates}(s, a, s', \mathcal{M}^{\uparrow})$ 15 if $s \models q$ or h > H then 16 $s \leftarrow s_0; f \leftarrow f + 1 \text{ iff } s \not\models g$ 17 else 18 $s \leftarrow s'$ 19 20 return \mathcal{M}^{\uparrow}

10

Continual Learning and Planning (CLaP)



Experiment Generation using Active Querying

[Verma, Karia, Srivastava; NeurlPS '23]

• Automated query generation to help fix each $(+/-/\emptyset)$ to a correct form.



Policy: Generated Autonomously by Reduction to Non-Deterministic Planning

This process monotonically improves the learned PPDDL model

```
(:action open-door
:parameters (?l1)
:precondition (and
  (+/-/\emptyset) (has_key)
  (+/-/\emptyset) (door_open)
  (+/-/∅) (door_adjacent ?l1)
  (+/-/Ø)(player_at ?l1))
:effect (probabilistic
  0.xx (and
     (+/-/\emptyset) (has_key)
     (+/-/\emptyset) (door_open)
     (+/-/∅) (door_adjacent ?l1)
     (+/-/\emptyset) (player_at ?l1))
  0.yy (and
     (+/-/\emptyset) (has_key)
     (+/-/\emptyset) (door_open)
     (+/-/∅) (door_adjacent ?l1)
     (+/-/Ø)(player_at ?l1)))
```

Continual Learning and Planning (CLaP)



Empirical Evaluation

Setup

Baselines

- 4 benchmark domains with 5 tasks each
 - Tireworld
 - ExplodingBlocks
 - FirstResponders
 - Elevators

- Changing action transitions, initial state, and goal for each task
- Limited simulator budget per task

Oracle

Knows precisely when each novelty is introduced.

Q-Learning

Vanilla Q-Learning that has no information about novelty.

Adaptive-QACE

Modification of QACE to incorporate differential learning. [Verma, Karia, Srivastava; NeurIPS '23]

CLaP's Performance Comparative to Oracle



CLaP Few Shot Transfers in Non-Stationary Settings



Tasks

Same Trends Across All Domains



Theoretical Guarantee: Locally Convergent Learning

Theorem 1. Let M be an RMDP with a series of transition system changes $\delta_1, \ldots, \delta_n$ at timesteps t_1, \ldots, t_n implemented using a simulator Δ , then during each stationary epoch between t_i and t_{i+1} CLaP performs locally convergent model learning.

- Between each consecutive stage, ClaP's model learning is *locally convergent*.
- Locally Convergent: The model keeps getting better.

Conclusions and Future Work

- CLaP is a sample-efficient method for solving tasks under non-stationarity
- If we use some resources for experiment design, that makes up for the costs by increasing efficiency

Future Directions

- Include information about the future goal(s)
- Include priors on the transition function change when available

