

Working Group “Probabilistic and Interactive Machine Learning”

AAMAS’23: Learning Constraints From Human Stop-Feedback

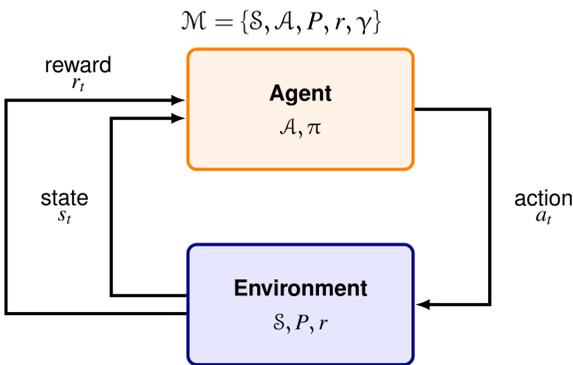
Sebastian Tschiatschek (with S. Poletti and A. Testolin)

Abstract

- **Safety**: central for the usage of intelligent agents in many domains
- **Lightbulb** In this paper: learning about dangerous behavior via **stop-feedback** in RL
 - **Probabilistic feedback model** inspired by how humans might provide feedback
 - **Bayesian inference** for inferring constraints
- **Experiments**:
 - ✓ Learning with our proposed feedback model is efficient
 - ✓ Human stop-feedback aligns reasonably well with our model

Setting

Markov Decision Processes



Classical Goal

- Cumulative rewards:

$$J(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \mid \pi \right],$$

- Find optimal policy:

$$\pi^* = \arg \max_{\pi} J(\pi).$$

Here: Learning With Constraints

- Rewards $r: \mathcal{S} \rightarrow \mathbb{R}_+$ known
- Constraints $c: \mathcal{S} \rightarrow \mathbb{R}_+$ unknown
- Dangerous states $\Leftrightarrow c(s) > 0$
- Agent should avoid dangerous states
- Find optimal policy for

$$J(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t) - c(s_t)) \mid \pi \right],$$

Stop-Feedback

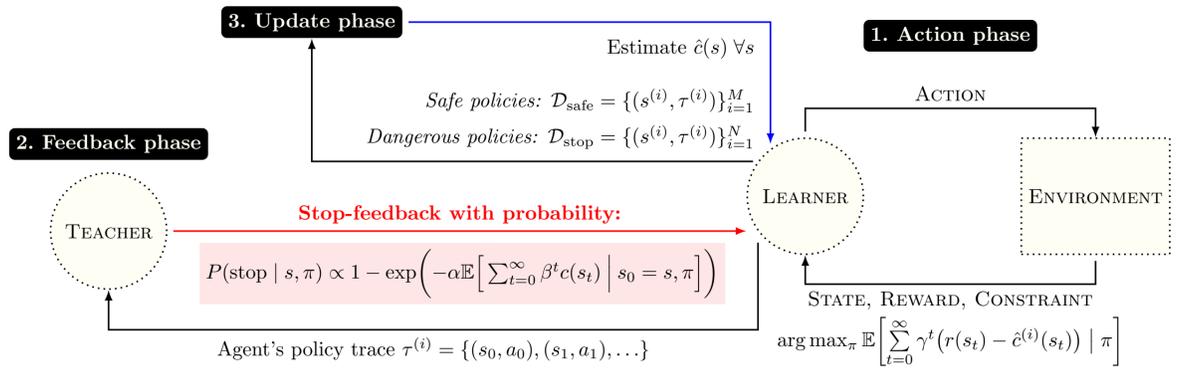
- Learn about constraints from stop-feedback (provided by a human supervisor)
- **Lightbulb** Model how a human supervisor might provide feedback

$$P(\text{stop} \mid s, \pi) \propto$$

$$1 - \exp \left(-\alpha \mathbb{E} \left[\sum_{t=0}^{\infty} \beta^t c(s_t) \mid s_0 = s, \pi \right] \right) \quad (1)$$

- β : Horizon for reasoning into the future
- α : “Worriedness” of supervisor
- Goal: sample-efficient learning from stop-feedback

Interaction of Agent and Supervisor



Our Approach

Central Ideas

- Estimate constraints from stop-feedback
- Be Bayesian:
 - Encourage exploration
 - Incorporate prior knowledge
- Featurize environment for scalability:

$$c(s) = \langle \phi(s), c^* \rangle$$

Major Steps

1 Action phase

2 Feedback phase

3 Update phase

- Compute $\hat{c}^{(i+1)}$ as posterior via MCMC
- Optimize policy:

$$\pi^{(i)} = \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t) - \hat{c}^{(i)}(s_t)) \mid \pi \right] \quad (2)$$

Input: maximum number of interactions K

Output: Final learner’s policy π^{K+1}

/ Initialization */*

- 1: $\mathcal{D}_{\text{stop}} \leftarrow \emptyset, \mathcal{D}_{\text{safe}} \leftarrow \emptyset$
- 2: $\hat{c}^{(1)}(s) \leftarrow 0 \quad \forall s \in \mathcal{S}$
- 3: $\pi^{(1)} \leftarrow$ (approx.) optimal policy for $r, \hat{c}^{(1)}$, cf. Eq. (2)
- /* Learner-teacher interaction */*
- 4: **for all** $i = 1, \dots, K$ **do**
- /* Action & feedback phase */*
- 5: $s \leftarrow s_0$
- 6: **for all** $t = 1, \dots, T$ **do**
- 7: $a_t \sim \pi^{(i)}(s), s_{t+1} \sim \mathcal{P}(\cdot \mid s_t, a_t), r_t \sim r_{a_t}(s_t)$
- 8: $f \leftarrow$ Teacher’s feedback according to Eq. (1)
- 9: **if** $f = \text{stop}$ **then**
- 10: $\mathcal{D}_{\text{stop}} \leftarrow \mathcal{D}_{\text{stop}} \cup \{(s_t, \pi^{(i)})\}$
- 11: **break**
- 12: **else**
- 13: $\mathcal{D}_{\text{safe}} \leftarrow \mathcal{D}_{\text{safe}} \cup \{(s_t, \pi^{(i)})\}$
- /* Update phase */*
- 14: Learner updates its estimate of the constraints to $\hat{c}^{(i+1)}$ based on the datasets $\mathcal{D}_{\text{stop}}, \mathcal{D}_{\text{safe}}$
- 15: $\pi^{(i+1)} \leftarrow$ (appr.) optimal policy for $r, \hat{c}^{(i+1)}$, cf. Eq. (2)
- 16: **return** Final learner’s policy π^{K+1}

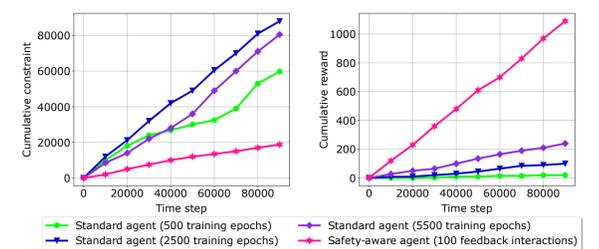
Experimental Setup

OpenAI Safety Gym

- 1 goal state
- 5 evenly distributed fixed hazards
- Agent can move in the 2d plane (turning & moving forward/backward)
- Environment is reset when reaching hazard
- PPO; 2nd layer of critic for feature extraction

Experimental Results

Learning from Synthetic Feedback



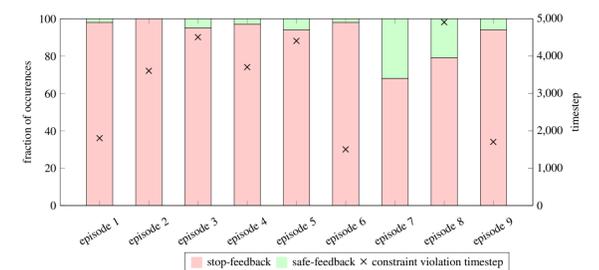
- Agents with constraint inference:

- Achieve higher cumulative rewards
- Violate fewer constraints

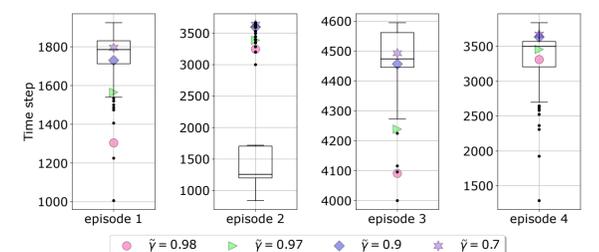
Human Feedback

- Survey with 100 volunteers (online)
- Evaluation of 9 videos:
 - Episodes with 5000 time steps
 - *Standard agent* colliding with ≤ 1 constraint
- Question:
 - Stop-feedback or not
 - Time step for feedback

Stop-feedback and collision times



Model-generated vs human stop-feedback



- Good alignment of feedback

Further Details

Paper link



Group link

