# Microsoft<sup>®</sup> Research

# ABSTRACT

# Problem setting

- Batch setting: fixed set of trajectories.
- Access to the behavioral policy called baseline.
- Objective: Improve the baseline with high probability.

# **Reliability issues in Reinforcement Learning**

Deep Reinforcement Learning is unreliable.

- No convergence/optimality proof,  $\rightarrow$  benchmark performance is the standard evaluation.
- Many, poorly understood hyperparameters,  $\rightarrow$  not reproducible from one task to another.
- Very unstable learning process [1],  $\rightarrow$  different seeds may yield very different results.

Batch Reinforcement Learning [2] is unreliable.

- Classic algorithms trained on a fixed dataset consider it as the environment (either explicitly or implicitly).
- This may be statistically insufficient, due to environmental stochasticity or function approximation.
- This leads to overestimates in the trained values.
- Classic algorithms perform planning, which turns out to be over-reliant on the overestimated states.

# Safe Policy Improvement with **Baseline Bootstrapping (SPIBB, [3])**

- Recent computationally efficient and provably-safe methodology for batch RL.
- SPIBB updates the policy for frequent state-action pairs in the dataset only.
- SPIBB relies on a binary decision-making and may be too conservative.

# Our contributions

- Novel batch RL method: Soft-SPIBB, a reformulation of SPIBB objective which makes it more flexible.
- Soft-SPIBB has proofs of safety guarantees and of computational efficiency in finite MDPs.
- Model-free Soft-SPIBB for function approximation.
- Empirical validation (performance, safety) on two domains.

# SAFE POLICY IMPROVEMENT WITH SOFT BASELINE BOOTSTRAPPING

Kimia Nadjahi<sup>\*</sup>, Romain Laroche<sup>\*</sup>, Rémi Tachet des Combes

kimia.nadjahi@telecom-paris.fr, {romain.laroche, remi.tachet}@microsoft.com

### THEORY

## Safe Policy Improvement with Soft Baseline Bootstrapping (Soft-SPIBB)

- ► True environment  $M^* = \langle \mathcal{X}, \mathcal{A}, P^*, R^*, \gamma \rangle$  is unknown,
- ► Maximum Likelihood Estimation (MLE) MDP built from counts:  $\widehat{M} = \langle \mathcal{X}, \mathcal{A}, \widehat{P}, \widehat{R}, \gamma \rangle$ .
- Error function e is derived from concentrations inequalities to bound the difference between parameters of M and  $M^*$ :  $e_Q$  between Q-functions,  $e_P$  between probabilities.
- Soft-SPIBB relies on a softer mechanism where, for a given error function, a local error budget is allocated for policy changes in each state.
- A policy  $\pi$  is  $(\pi_b, e, \epsilon)$ -constrained with hyper-parameter  $\epsilon$  if, for each state  $x \in \mathcal{X}$ ,

$$\sum_{a \in \mathcal{A}} e(x, a) |\pi(a|x) - \pi_b(a|)$$

A policy  $\pi$  is  $\pi_b$ -advantageous in  $\widehat{M}$  if, for all  $x \in \mathcal{X}$ :

$$\sum_{a \in \mathcal{A}} \left( \mathcal{Q}_{\widehat{M}}^{\pi_b}(x, a) - V_{\widehat{M}}^{\pi_b}(x) \right) \pi(a)$$

- Fixed Two algorithms performing policy iteration in the space of  $(\pi_b, e, \epsilon)$ -constrained policies: *Exact-Soft-SPIBB* (exact solution) and *Approx-Soft-SPIBB* (tractable approximate solution).
- $\blacktriangleright$  Model-free formulation, which fits the Q-function to the following targets:  $y_j^{(i+1)} = r_j + \gamma \sum \pi^{(i+1)}(a'|x_j)Q^{(i)}(x_j, a').$

### Theorems

### Safe policy improvement bounds.

► Assume  $\pi$  is  $(\pi_b, e_Q, \epsilon)$ -constrained and  $\pi_b$ -advantageous in  $\widehat{M}$ . For each state x, with high probability  $1 - \delta$ :

$$V_{M^*}^{\pi}(x) - V_{M^*}^{\pi_b}(x) \geq -\frac{\epsilon}{1}$$

Assume  $\pi$  is  $(\pi_b, e_P, \epsilon)$ -constrained. Suppose there exists  $\kappa < \frac{1}{2}$  such that,

$$orall (x,a) \in \mathcal{X} imes \mathcal{A}, \ \sum_{x',a'} e_P(x',a') \pi_b(a'|x') P$$

We denote by  $d_M(\cdot | x)$  the  $\gamma$ -discounted future state distribution starting from x when following  $\pi_b$  in M. Then, for each state x, with high probability  $1 - \delta$ :

$$V_{M^*}^{\pi}(x) - V_{M^*}^{\pi_b}(x) \geq V_{\widehat{M}}^{\pi}(x) - V_{\widehat{M}}^{\pi_b}(x) - 2 \left\| d_{M^*}(\cdot | x) - d_{\widehat{M}}(\cdot | x) 
ight\|_1 V_{max} \ - rac{1+\gamma}{\left(1-\gamma
ight)^2 \left(1-\kappa\gamma
ight)} \, \, \epsilon V_{max} \, .$$

Analysis of Approx-Soft-SPIBB. The policy improvement step of Approx-Soft-SPIBB generates  $(\pi_b, e, \epsilon)$ -constrained policies, and has a complexity of  $\mathcal{O}(|\mathcal{X}||\mathcal{A}|^2)$ .

Model-free formulation equivalence. In finite MDPs, the model-free policy iteration of (Exact or Approx)-Soft-SPIBB coincides with the model-based counterparts.

 $|\mathbf{X}| \leq \epsilon.$ 

 $(a|x) \geq 0.$ 

/ max  $\mathcal{P}^*(x'|x,a) \leq \kappa e_P(x,a)$ .





[1] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger. Deep RL that matters. In AAAI, 2018. [2] S. Lange, T. Gabel, and M. Riedmiller. Batch Reinforcement Learning. In *Reinforcement Learning*. 2012. [3] R. Laroche, P. Trichelair, and R. Tachet des Combes. Safe policy improvement with baseline bootstrapping. In ICML, 2019. [4] M. Petrik, M. Ghavamzadeh, and Y. Chow. Safe policy improvement by minimizing robust baseline regret. In NIPS, 2016.



 $\epsilon = 2$ , influence of  $\eta$ 

Figures (a-b) show the mean and CVaR benchmark. Figures (c-d) (resp. (e-f)) study the sensitivity to  $\eta$  (resp. to hyperparameters) for RaMDP [4] and Soft-SPIBB.

# REFERENCES