

Safe Policy Improvement with Baseline Bootstrapping

Romain Laroche, Paul Trichelair, Rémi Tachet des Combes



Problem setting

Batch setting

- Fixed dataset, no direct interaction with the environment.
- Access to the behavioural policy, called baseline.
- Objective: improve the baseline with high probability.
- Commonly encountered in real world applications.

Problem setting

Batch setting

- Fixed dataset, no direct interaction with the environment.
- Access to the behavioural policy, called baseline.
- Objective: improve the baseline with high probability.
- Commonly encountered in real world applications.

Distributed systems



Problem setting

Batch setting

- Fixed dataset, no direct interaction with the environment.
- Access to the behavioural policy, called baseline.
- Objective: improve the baseline with high probability.
- Commonly encountered in real world applications.

Distributed systems



Long trajectories



Contributions

Novel batch RL algorithm: SPIBB

- SPIBB comes with reliability guarantees in finite MDPs.
- SPIBB is as computationally efficient as classic RL.

Contributions

Novel batch RL algorithm: SPIBB

- SPIBB comes with reliability guarantees in finite MDPs.
- SPIBB is as computationally efficient as classic RL.

Finite MDPs benchmark

- Extensive benchmark of existing algorithms.
- Empirical analysis on random MDPs and baselines.

Contributions

Novel batch RL algorithm: SPIBB

- SPIBB comes with reliability guarantees in finite MDPs.
- SPIBB is as computationally efficient as classic RL.

Finite MDPs benchmark

- Extensive benchmark of existing algorithms.
- Empirical analysis on random MDPs and baselines.

Infinite MDPs benchmark

- Model-free SPIBB for use with function approximation.
- First deep RL algorithm reliable in the batch setting.

Robust Markov Decision Processes

[Iyengar, 2005, Nilim and El Ghaoui, 2005]

- True environment $M^* = \langle \mathcal{X}, \mathcal{A}, P^*, R^*, \gamma \rangle$ is unknown.
- Maximum Likelihood Estimation (MLE) MDP built from counts: $\hat{M} = \langle \mathcal{X}, \mathcal{A}, \hat{P}, \hat{R}, \gamma \rangle$.
- Robust MDP set $\Xi(\hat{M}, e)$: $M^* \in \Xi(\hat{M}, e)$ with probability at least $1 - \delta$.
- Error function $e(x, a)$ derived from concentration bounds.

Existing algorithms

[Petrik et al., 2016]: SPI by robust baseline regret minimization

- Robust MDPs considers the maxmin of the value over Ξ ,
→ favors over-conservative policies.
- They also consider the maxmin of the value improvement,
→ NP-hard problem.
- RaMDP hacks the reward to account for uncertainty:

$$\tilde{R}(x, a) \leftarrow \hat{R}(x, a) - \frac{\kappa_{adj}}{\sqrt{N_{\mathcal{D}}(x, a)}},$$

→ not theoretically grounded.

Existing algorithms

[Petrik et al., 2016]: SPI by robust baseline regret minimization

- Robust MDPs considers the maxmin of the value over Ξ ,
→ favors over-conservative policies.
- They also consider the maxmin of the value improvement,
→ NP-hard problem.
- RaMDP hacks the reward to account for uncertainty:

$$\tilde{R}(x, a) \leftarrow \hat{R}(x, a) - \frac{\kappa_{adj}}{\sqrt{N_{\mathcal{D}}(x, a)}},$$

→ not theoretically grounded.

[Thomas, 2015]: High-Confidence Policy Improvement

- HCPI searches for the best regularization hyperparameter to allow safe policy improvement.

Safe Policy Improvement with Baseline Bootstrapping

Safe Policy Improvement with Baseline Bootstrapping (SPIBB)

- Tractable approximate solution to the robust policy improvement formulation.
- SPIBB allows policy update only with sufficient evidence.
- Sufficient evidence = state-action count that exceeds some threshold hyperparameter N_{\wedge} .

Safe Policy Improvement with Baseline Bootstrapping

Safe Policy Improvement with Baseline Bootstrapping (SPIBB)

- Tractable approximate solution to the robust policy improvement formulation.
- SPIBB allows policy update only with sufficient evidence.
- Sufficient evidence = state-action count that exceeds some threshold hyperparameter N_\wedge .

SPIBB algorithm

- Construction of the *bootstrapped set*:
$$\mathfrak{B} = \{(x, a) \in \mathcal{X} \times \mathcal{A}, N_D(x, a) < N_\wedge\}.$$
- Optimization over a constrained policy set:

$$\pi_{spibb}^\odot = \operatorname{argmax}_{\pi \in \Pi_b} \rho(\pi, \hat{M}),$$

$$\Pi_b = \{\pi, \text{ s.t. } \pi(a|x) = \pi_b(a|x) \text{ if } (x, a) \in \mathfrak{B}\}.$$

SPIBB policy iteration



SPIBB policy iteration



Policy improvement step example

Q-value	Baseline policy	Bootstrapping	SPIBB policy update
$Q_{\hat{M}}^{(l)}(x, a_1) = 1$	$\pi_b(a_1 x) = 0.1$	$(x, a_1) \in \mathfrak{B}$	
$Q_{\hat{M}}^{(l)}(x, a_2) = 2$	$\pi_b(a_2 x) = 0.4$	$(x, a_2) \notin \mathfrak{B}$	
$Q_{\hat{M}}^{(l)}(x, a_3) = 3$	$\pi_b(a_3 x) = 0.3$	$(x, a_3) \notin \mathfrak{B}$	
$Q_{\hat{M}}^{(l)}(x, a_4) = 4$	$\pi_b(a_4 x) = 0.2$	$(x, a_4) \in \mathfrak{B}$	

SPIBB policy iteration



Policy improvement step example

Q-value	Baseline policy	Bootstrapping	SPIBB policy update
$Q_{\hat{M}}^{(l)}(x, a_1) = 1$	$\pi_b(a_1 x) = 0.1$	$(x, a_1) \in \mathfrak{B}$	$\pi^{(i+1)}(a_1 x) = 0.1$
$Q_{\hat{M}}^{(l)}(x, a_2) = 2$	$\pi_b(a_2 x) = 0.4$	$(x, a_2) \notin \mathfrak{B}$	
$Q_{\hat{M}}^{(l)}(x, a_3) = 3$	$\pi_b(a_3 x) = 0.3$	$(x, a_3) \notin \mathfrak{B}$	
$Q_{\hat{M}}^{(l)}(x, a_4) = 4$	$\pi_b(a_4 x) = 0.2$	$(x, a_4) \in \mathfrak{B}$	$\pi^{(i+1)}(a_4 x) = 0.2$

SPIBB policy iteration



Policy improvement step example

Q-value	Baseline policy	Bootstrapping	SPIBB policy update
$Q_{\hat{M}}^{(i)}(x, a_1) = 1$	$\pi_b(a_1 x) = 0.1$	$(x, a_1) \in \mathfrak{B}$	$\pi^{(i+1)}(a_1 x) = 0.1$
$Q_{\hat{M}}^{(i)}(x, a_2) = 2$	$\pi_b(a_2 x) = 0.4$	$(x, a_2) \notin \mathfrak{B}$	$\pi^{(i+1)}(a_2 x) = 0.0$
$Q_{\hat{M}}^{(i)}(x, a_3) = 3$	$\pi_b(a_3 x) = 0.3$	$(x, a_3) \notin \mathfrak{B}$	$\pi^{(i+1)}(a_3 x) = 0.7$
$Q_{\hat{M}}^{(i)}(x, a_4) = 4$	$\pi_b(a_4 x) = 0.2$	$(x, a_4) \in \mathfrak{B}$	$\pi^{(i+1)}(a_4 x) = 0.2$

Theoretical analysis

Theorem (Convergence)

Policy iteration converges to a policy π_{spibb}^{\ominus} that is Π_b -optimal in the MLE MDP \hat{M} .

Theoretical analysis

Theorem (Convergence)

Policy iteration converges to a policy π_{spibb}^{\odot} that is Π_b -optimal in the MLE MDP \hat{M} .

Theorem (Safe policy improvement)

With high probability $1 - \delta$:

$$\rho(\pi_{spibb}^{\odot}, M^*) - \rho(\pi_b, M^*) \geq \rho(\pi_{spibb}^{\odot}, \hat{M}) - \rho(\pi_b, \hat{M}) - \frac{4V_{max}}{1-\gamma} \sqrt{\frac{2}{N_{\wedge}} \log \frac{2|\mathcal{X}||\mathcal{A}|2^{|\mathcal{X}|}}{\delta}}$$

Model-free formulation

SPIBB algorithm

- It may be formulated in a model-free manner by setting the targets:

$$y_j^{(i)} = r_j + \gamma \sum_{a' | (x'_j, a') \in \mathfrak{B}} \pi_b(a' | x'_j) Q^{(i)}(x'_j, a') \\ + \gamma \left(\sum_{a' | (x'_j, a') \notin \mathfrak{B}} \pi_b(a' | x'_j) \right) \max_{a' | (x'_j, a') \notin \mathfrak{B}} Q^{(i)}(x'_j, a').$$

Model-free formulation

SPIBB algorithm

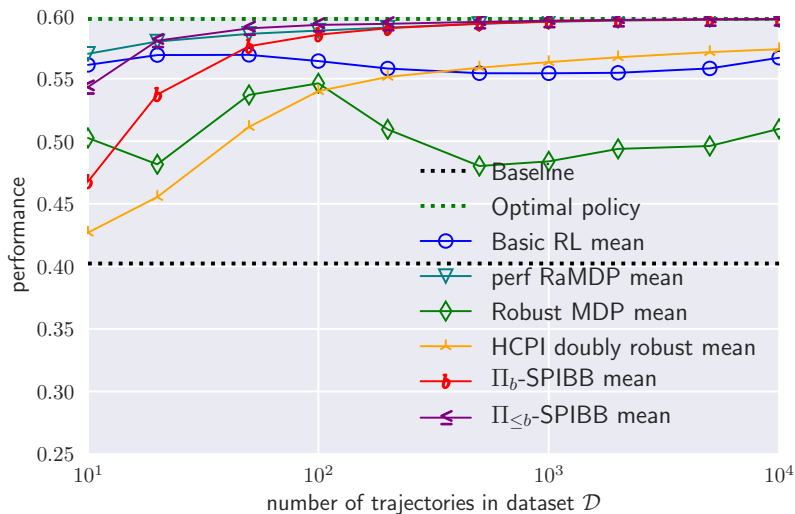
- It may be formulated in a model-free manner by setting the targets:

$$y_j^{(i)} = r_j + \gamma \sum_{a' | (x'_j, a') \in \mathfrak{B}} \pi_b(a' | x'_j) Q^{(i)}(x'_j, a') \\ + \gamma \left(\sum_{a' | (x'_j, a') \notin \mathfrak{B}} \pi_b(a' | x'_j) \right) \max_{a' | (x'_j, a') \notin \mathfrak{B}} Q^{(i)}(x'_j, a').$$

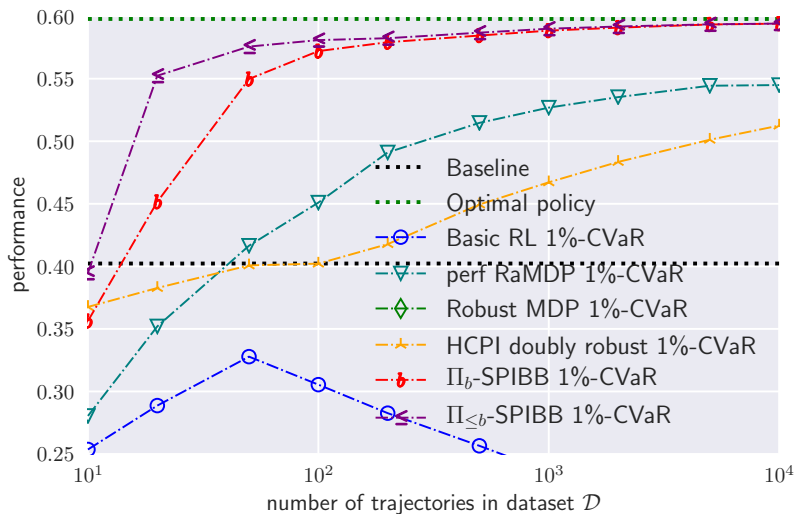
Theorem (Model-free formulation equivalence)

In finite MDPs, the model-free formulation admits a unique fixed point that coincides with the Q -value of π_{spibb}^\odot .

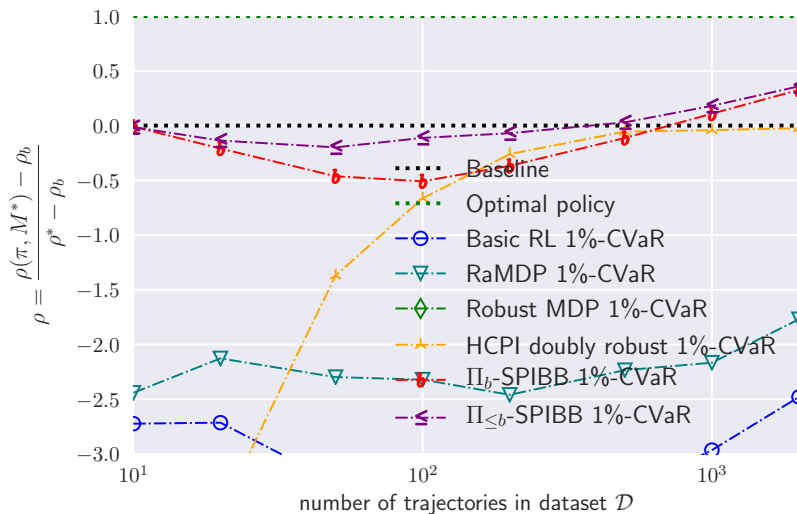
25-state stochastic gridworld – mean



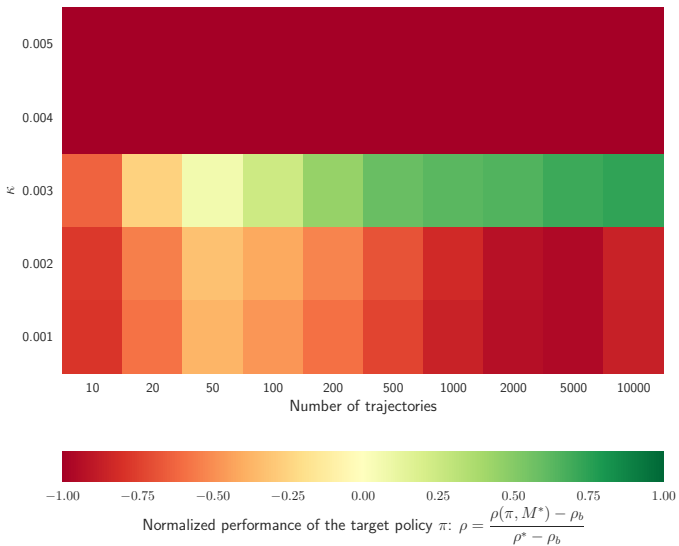
25-state stochastic gridworld – 1%-CVaR



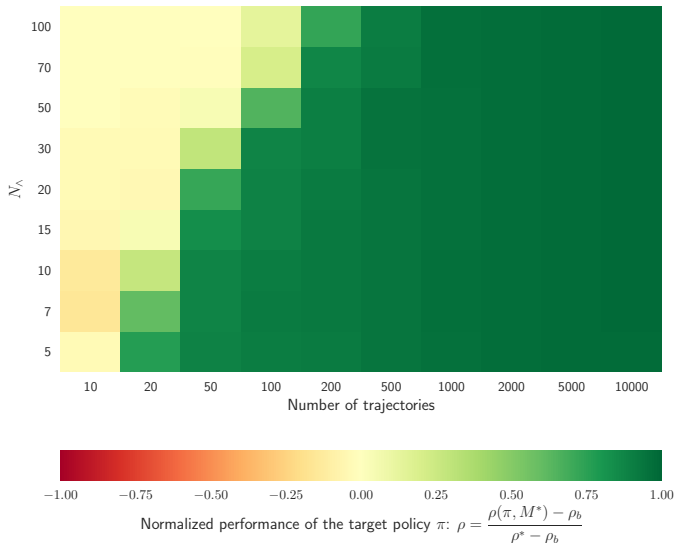
Random MDPs, random baseline – 1%-CVaR



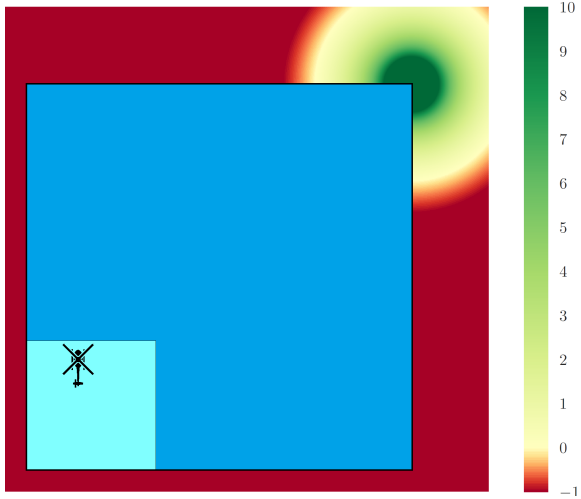
Gridworld – RaMDP hyperparameter sensitivity



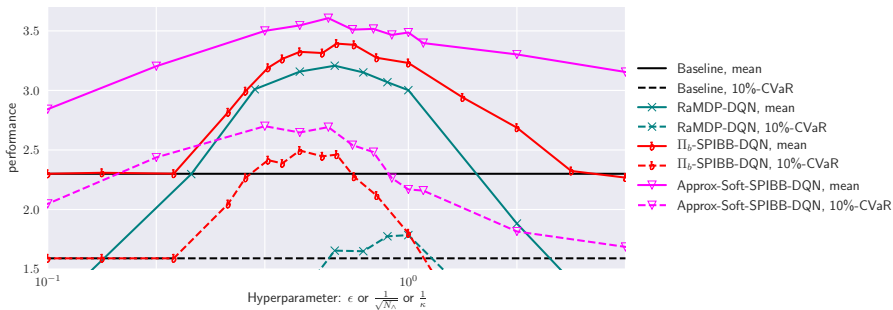
Gridworld – SPIBB hyperparameter sensitivity



Helicopter domain (continuous task)



Helicopter domain - benchmark (improved results)



Vanilla DQN is off the chart

- mean = 0.22,
- 10%-CVaR = -1 (minimal score).

Conclusion

SPIBB

- Assumes fixed dataset, and known behavioural policy.
- Tractable, provably reliable, sample-efficient algorithm.
- Successfully transferred to DQN architectures.

Conclusion

SPIBB

- Assumes fixed dataset, and known behavioural policy.
- Tractable, provably reliable, sample-efficient algorithm.
- Successfully transferred to DQN architectures.

Follow-up work

- Factored SPIBB [Simão and Spaan, 2019a].
- Structure learning coupled [Simão and Spaan, 2019b].
- Soft SPIBB [Nadjahi et al., 2019].

Conclusion

SPIBB

- Assumes fixed dataset, and known behavioural policy.
- Tractable, provably reliable, sample-efficient algorithm.
- Successfully transferred to DQN architectures.

Follow-up work

- Factored SPIBB [Simão and Spaan, 2019a].
- Structure learning coupled [Simão and Spaan, 2019b].
- Soft SPIBB [Nadjahi et al., 2019].

Still to do

- Improve the pseudo-count/error estimates.
- Investigate an online SPIBB inspired algorithm.

Thanks for your attention

(POSTER #101)



Iyengar, G. N. (2005).

Robust dynamic programming.

Mathematics of Operations Research.



Nadjahi, K., Laroche, R., and Tachet des Combes, R. (2019).

Safe policy improvement with soft baseline bootstrapping.

In Proceedings of the 17th European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD).



Nilim, A. and El Ghaoui, L. (2005).

Robust control of markov decision processes with uncertain transition matrices.

Operations Research.



Petrik, M., Ghavamzadeh, M., and Chow, Y. (2016).

Safe policy improvement by minimizing robust baseline regret.

In Proceedings of the 29th Advances in Neural Information Processing Systems (NIPS).



Simão, T. D. and Spaan, M. T. J. (2019a).

Safe policy improvement with baseline bootstrapping in factored environments.

In Proceedings of the 32nd AAAI Conference on Artificial Intelligence.



Simão, T. D. and Spaan, M. T. J. (2019b).

Structure learning for safe policy improvement.

In Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI).



Thomas, P. S. (2015).

Safe reinforcement learning.

PhD thesis, Stanford university.