



TRANSFER OF VALUE FUNCTIONS VIA VARIATIONAL METHODS

ANDREA TIRINZONI, RAFAEL RODRIGUEZ, AND MARCELLO RESTELLI
{andrea.tirinzi, marcello.restelli}@polimi.it, rafaelalberto.rodriguez@mail.polimi.it



PROBLEM

- The agent has solved a finite set of **source tasks** $\mathcal{M}_{\tau_1}, \mathcal{M}_{\tau_2}, \dots, \mathcal{M}_{\tau_M}$ sampled from some **distribution** \mathcal{D}
- Each task is an MDP $\mathcal{M}_\tau = \langle S, \mathcal{A}, \mathcal{P}_\tau, \mathcal{R}_\tau, p_0 \rangle$
- A parametric approximation to their **optimal value functions** is available
 $\mathcal{W}_s = \{w_1, w_2, \dots, w_M\}$ s.t. $Q_{w_j} \simeq Q_{\tau_j}^*$
- Assumption**: all tasks share **similarities** in their optimal value functions [4]
- Goal**: use this knowledge to speed-up the learning process of a new **target task** \mathcal{M}_τ sampled from \mathcal{D}

MOTIVATION

- Reinforcement learning** algorithms have enjoyed many success stories in complicated tasks
- High **sample complexity** remains a major issue
- Must **adapt** to changing environments and goals
- Prior knowledge** from related tasks is often available in practice \rightarrow **Transfer learning** [6]
- Need for transfer algorithms that are **general** and **widely applicable**

CONTRIBUTIONS

- Algorithmic**. We propose a general **framework** for transferring value functions in RL and two **practical algorithms**
 - We learn a **prior** distribution over optimal Q -functions using the given source tasks
 - Variational** approximation of the corresponding posterior for a new target task
 - Efficient **exploration** via posterior sampling
 - Any** differentiable Q -function approximator and prior/posterior models could be used
- Theoretical**. We provide a theoretical analysis of our practical algorithms offering a better insight into their behavior
- Empirical**. We empirically evaluate our algorithms on four different domains with increasing level of difficulty

VARIATIONAL TRANSFER FRAMEWORK

IDEA: use the source weights \mathcal{W}_s to estimate the distribution $p(w)$ over optimal Q -functions induced by \mathcal{D}

- How to characterize $p(w|D) \propto p(D|w)p(w)$ given a dataset D of N samples from the target task?
- PAC-Bayes argument** [3]: the likelihood $p(D|w)$ decays exponentially as the TD error of Q_w on D increases

$$p(w|D) \simeq \frac{e^{-\Lambda \|B_w\|_D^2} p(w)}{\int e^{-\Lambda \|B_{w'}\|_D^2} p(dw')}$$

- Problem**: computing the Gibbs posterior is often intractable \rightarrow **Variational approximation** [1]

$$\min_{\xi \in \Xi} \mathcal{L}(\xi) = \mathbb{E}_{w \sim q_\xi} [\|B_w\|_D^2] + \frac{\lambda}{N} KL(q_\xi(w) \| p(w))$$

MAIN PROPERTIES

- Prior estimation**: summarize the information to transfer into a single distribution and use it to guide the learning process of the target task
- Exploration via posterior sampling** [5, 2]: at each time, the agent guesses the solution of the target task according to the current posterior and acts accordingly
- Black-box optimization**: minimizing the variational objective requires only differentiability of the models involved

Algorithm Variational Transfer

Input: Target task \mathcal{M}_τ , source weights \mathcal{W}_s

Estimate prior $p(w)$ from \mathcal{W}_s

$\xi \leftarrow \operatorname{argmin}_\xi KL(q_\xi \| p), D \leftarrow \emptyset$

repeat

Sample initial state: $s_0 \sim p_0$

while s_h is not terminal **do**

$a_h = \operatorname{argmax}_a Q_w(s_h, a)$ for $w \sim q_\xi(w)$

$s_{h+1} \sim \mathcal{P}_\tau(\cdot | s_h, a_h), r_{h+1} = \mathcal{R}_\tau(s_h, a_h)$

$D \leftarrow D \cup \langle s_h, a_h, r_{h+1}, s_{h+1} \rangle$

$\xi \leftarrow \operatorname{optimizer}(\xi, \nabla_\xi \mathcal{L}(\xi))$

end while

until forever

PRACTICAL ALGORITHMS

GAUSSIAN VARIATIONAL TRANSFER (GVT)

- Prior: $p(w) = \mathcal{N}(\mu_p, \Sigma_p)$
- Posterior: $q_\xi(w) = \mathcal{N}(\mu, \Sigma)$

MIXTURE OF GAUSSIAN VARIATIONAL TRANSFER (MGVT)

- Prior: $p(w) = |\mathcal{W}_s|^{-1} \sum_{w_s \in \mathcal{W}_s} \mathcal{N}(w | w_s, \sigma_p^2 I)$
- Posterior: $q_\xi(w) = C^{-1} \sum_{i=1, \dots, C} \mathcal{N}(w | \mu_i, \Sigma_i)$

FINITE-SAMPLE ANALYSIS

Bound the **expected Bellman error** under the optimal variational distribution for a dataset of N samples

$$\mathbb{E}_{q_\xi} [\| \tilde{B}_w \|^2] \leq 2 \| \tilde{B}_{w^*} \|^2 + v(w^*) + c_1 \sqrt{\frac{\log \frac{2}{\delta}}{N}} + \frac{c_2 + \lambda d \log N + \lambda \varphi(\mathcal{W}_s)}{N} + \frac{c_3}{N^2}$$

- Approximation error** due to the limited hypothesis space
- Variance** due to a biased estimation of the Bellman error
- Variance** due to the finite samples
- Likelihood** of the optimal target weights under the prior

GVT: Distance to the prior mean

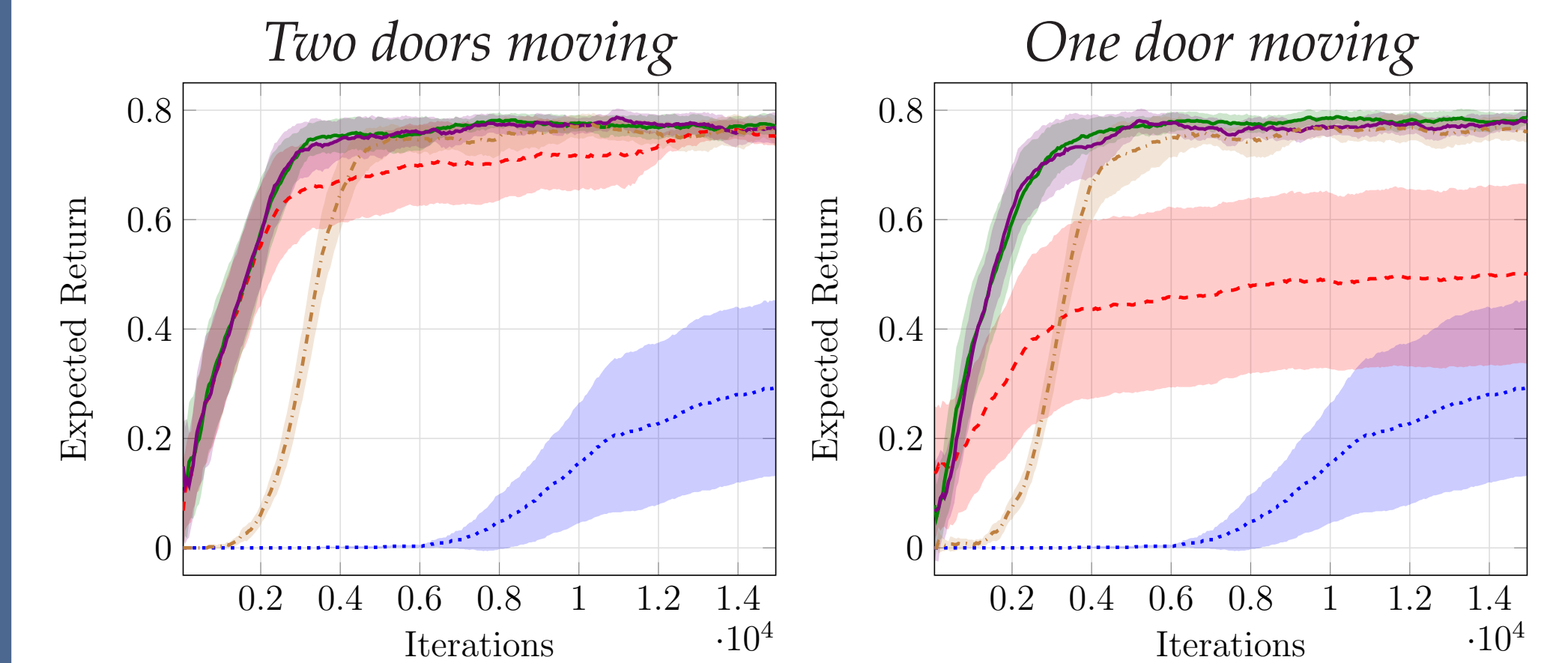
$$\varphi(\mathcal{W}_s) = \|w^* - \mu_p\|_{\Sigma_p^{-1}}$$

MGVT: **Softmin** distance to the sources

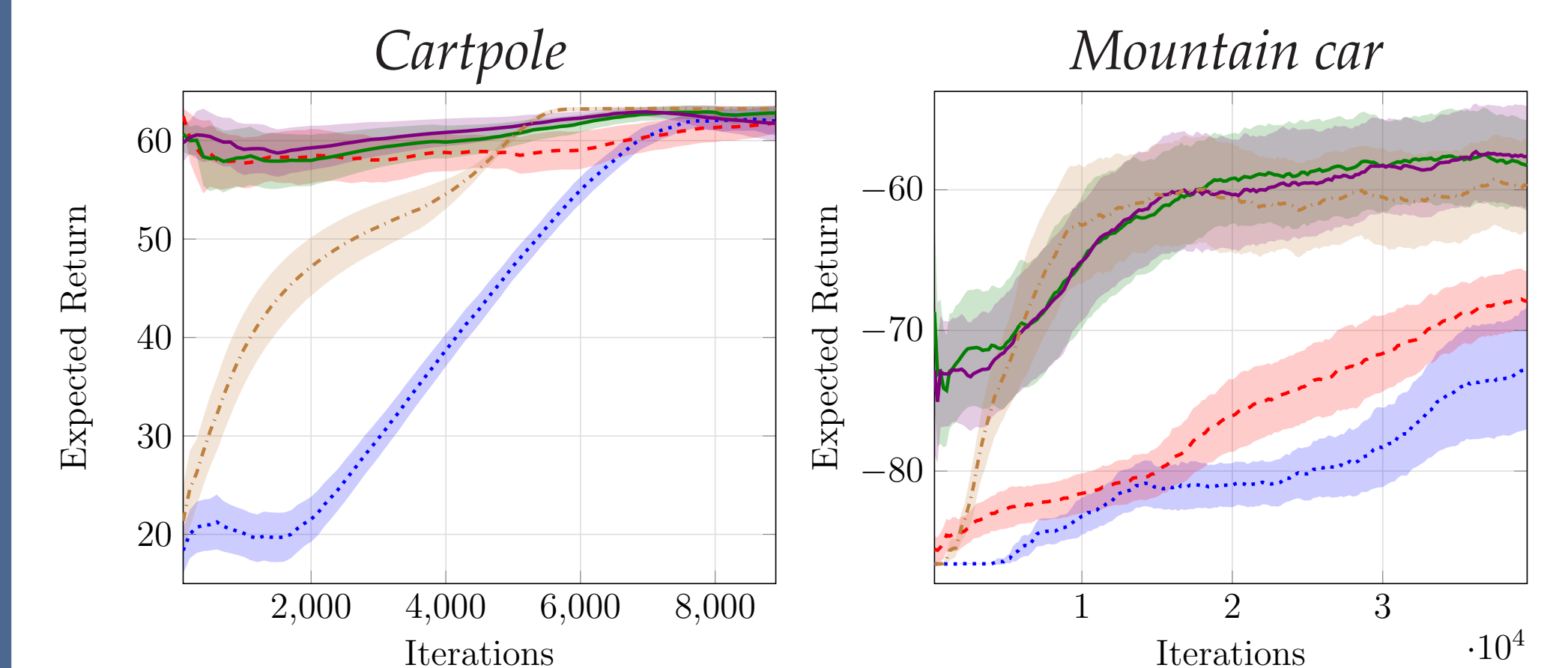
$$\varphi(\mathcal{W}_s) = \operatorname{softmin}_{w \in \mathcal{W}_s} (\|w^* - w\|)$$

EMPIRICAL RESULTS

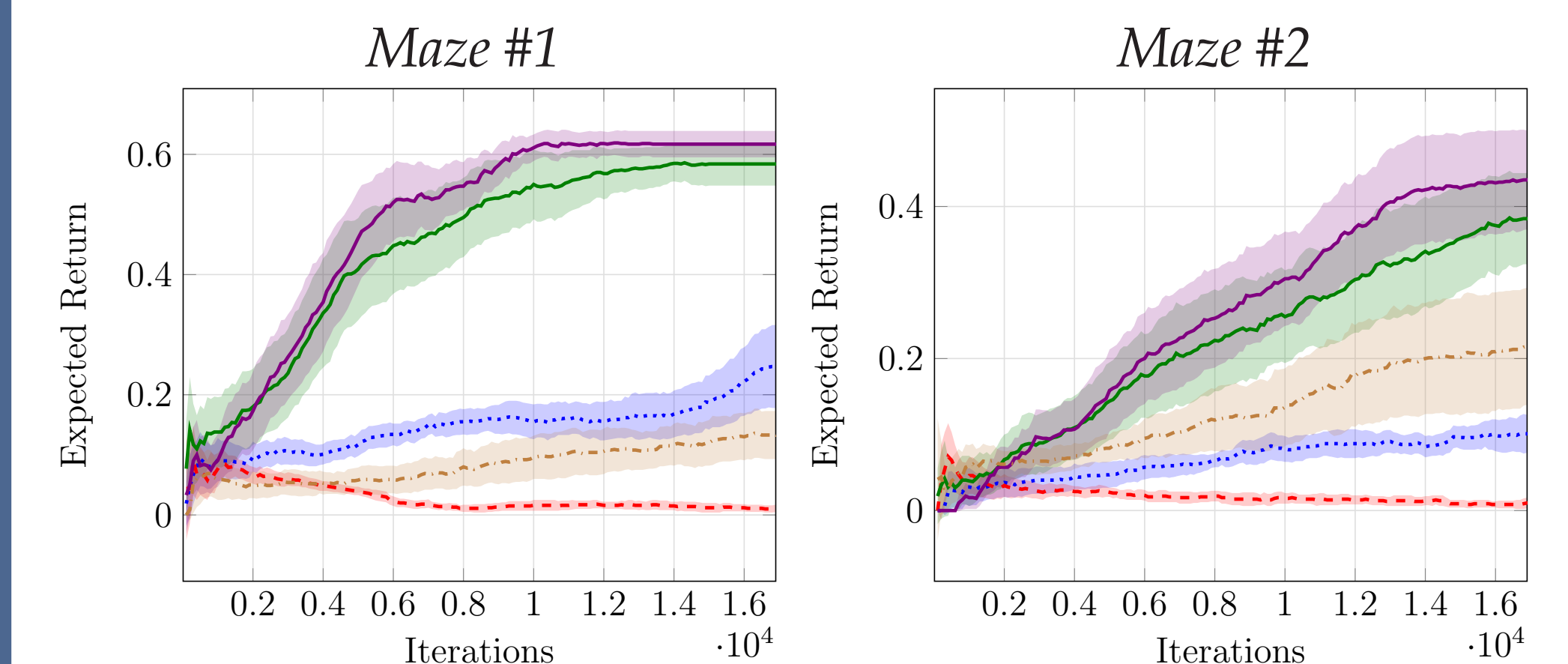
THE ROOMS PROBLEM



CLASSIC CONTROL



MAZE NAVIGATION



- No transfer (DDQN)
- - - - Fine-tuning from random source task
- - - - GVT
- MGVT with 1 component
- MGVT with 3 components

REFERENCES

- Pierre Alquier, James Ridgway, and Nicolas Chopin. On the properties of variational approximations of gibbs posteriors. *Journal of Machine Learning Research*, 17(239):1–41, 2016.
- Kamyar Aizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient exploration through bayesian deep q-networks. *arXiv preprint arXiv:1802.04412*, 2018.
- Olivier Catoni. Pac-bayesian supervised classification: the thermodynamics of statistical learning. *arXiv preprint arXiv:0712.0248*, 2007.
- Alessandro Lazaric and Mohammad Ghavamzadeh. Bayesian multi-task reinforcement learning. In *ICML-27th International Conference on Machine Learning*, pages 599–606. Omnipress, 2010.
- Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. *arXiv preprint arXiv:1402.0635*, 2014.
- Matthew E Taylor and Peter Stone. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 2009.