



TRANSFERRING VALUE FUNCTIONS VIA VARIATIONAL METHODS

Andrea Tirinzoni, Rafael A. Rodriguez, and Marcello Restelli

14th European Workshop on Reinforcement Learning, Lille, France



POLITECNICO
MILANO 1863

Why do we need transfer?

Reinforcement Learning (RL) has been successfully applied to many **complex tasks**



[Heess et al., 2017]



[OpenAI, 2018]



[Vinyals et al., 2017]

- High **sample complexity** remains a major limitation
- Both humans and artificial agents repeatedly face several related tasks
 - Changing environments
 - New goals

} **Transfer**

Our Settings

- The agent has solved a *finite* set of **source tasks** sampled from a **distribution** \mathcal{D}

$$\mathcal{M}_{\tau_1}, \mathcal{M}_{\tau_2}, \dots, \mathcal{M}_{\tau_M} \text{ s.t. } \mathcal{M}_{\tau} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}_{\tau}, \mathcal{R}_{\tau}, p_0 \rangle \sim \mathcal{D}$$

- A parametric approximation to their **optimal value functions** is available

$$\mathcal{W}_s = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\} \text{ s.t. } Q_{\mathbf{w}_j} \simeq Q_{\tau_j}^*$$

- **Goal:** use this knowledge to speed-up the learning process of a new **target task** \mathcal{M}_{τ} sampled from \mathcal{D}

Related Works

- Use the source Q -functions as **initializers**
[Tanaka and Yamamura, 2003, Taylor and Stone, 2009, Abel et al., 2018]
- **Bayesian** methods
[Lazaric and Ghavamzadeh, 2010]
- Transfer via **successor features**
[Barreto et al., 2017, Barreto et al., 2018]
- **Fast adaptation** / Meta-learning
[Finn et al., 2017, Grant et al., 2018, Amit and Meir, 2018]

What we would like to have

Our transfer algorithm should

- *not* make strong **assumptions** that limit its applicability
- dynamically use information from the source tasks during the learning process
- drive the **exploration** of the target task based on transferred knowledge

Transfer via Variational Methods

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

- How to characterize $p(\mathbf{w}|D) \propto p(D|\mathbf{w}) p(\mathbf{w})$ given a dataset D of N samples from the target task?
- PAC-Bayes argument** [Catoni, 2007]: the likelihood $p(D|\mathbf{w})$ decays exponentially as the *TD error* $\|B_{\mathbf{w}}\|_D^2$ of $Q_{\mathbf{w}}$ on D increases

$$p(\mathbf{w}|D) \simeq \frac{e^{-\Lambda\|B_{\mathbf{w}}\|_D^2} p(\mathbf{w})}{\underbrace{\int e^{-\Lambda\|B_{\mathbf{w}'}\|_D^2} p(d\mathbf{w}')}_{\text{Gibbs posterior}}}$$

Transfer via Variational Methods

Problem: computing the Gibbs posterior $q(\mathbf{w})$ is often intractable

- **Variational approximation** [Alquier et al., 2016] $\rightarrow \operatorname{argmin}_{\xi} KL(q_{\xi}(\mathbf{w}) \parallel q(\mathbf{w}))$

$$\underbrace{\min_{\xi \in \Xi} \mathcal{L}(\xi)}_{\text{Variational objective}} = \underbrace{\mathbb{E}_{\mathbf{w} \sim q_{\xi}} \left[\|B_{\mathbf{w}}\|_D^2 \right]}_{\text{Expected TD error}} + \frac{\lambda}{N} \underbrace{KL(q_{\xi}(\mathbf{w}) \parallel p(\mathbf{w}))}_{\text{Divergence w.r.t. the prior}}$$

Algorithm 1 Variational Transfer

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

Estimate prior $p(\mathbf{w})$ 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**

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

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

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

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

end while

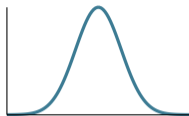
until forever

- Summarize transferred information into the **prior** distribution
- Exploration via **posterior sampling** [Osband et al., 2014]
- Requires only **differentiable** models

GAUSSIAN VARIATIONAL TRANSFER (GVT)

- Prior/posterior models are multivariate **Gaussians**

$$p(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p), \quad q_{\xi}(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$



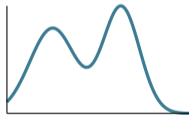
MIXTURE OF GAUSSIAN VARIATIONAL TRANSFER (MGVT)

- **Kernel density estimator** for the prior

$$p(\mathbf{w}) = \frac{1}{|\mathcal{W}_s|} \sum_{\mathbf{w}_s \in \mathcal{W}_s} \mathcal{N}(\mathbf{w} | \mathbf{w}_s, \sigma_p^2 \mathbf{I})$$

- C -component **mixture of Gaussian** model for the posterior

$$q_{\xi}(\mathbf{w}) = \frac{1}{C} \sum_{i=1}^C \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

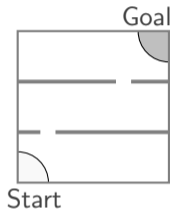


Theoretical Properties

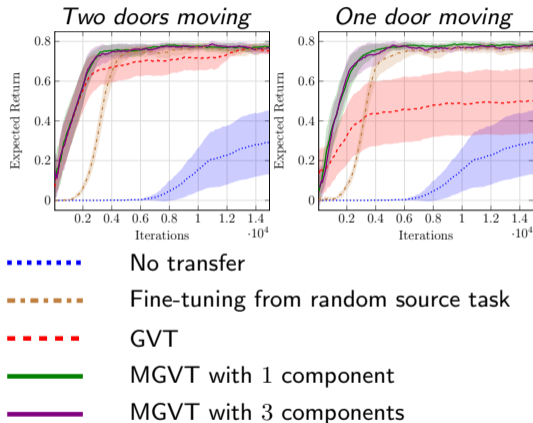
We provide a **finite-sample** analysis of both our practical algorithms

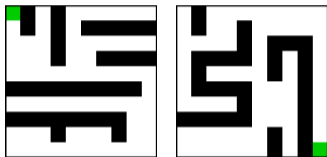
$$\underbrace{\mathbb{E}_{q_{\hat{\xi}}} \left[\left\| \tilde{B} \mathbf{w} \right\|_{\nu}^2 \right]}_{\text{Expected Bellman error}} \leq \underbrace{2 \left\| \tilde{B} \mathbf{w}^* \right\|_{\nu}^2}_{\text{Approximation error}} + \underbrace{v(\mathbf{w}^*)}_{\text{Variance}} + \underbrace{\frac{\lambda}{N} \varphi(\mathcal{W}_s)}_{\text{Distance to the prior}} + \mathcal{O}(1/N)$$

- Bound the expected Bellman error under the *optimal variational distribution* $q_{\hat{\xi}}$
- Same bounds, different distances to the prior
 - GVT: distance to the prior *mean* $\varphi(\mathcal{W}_s) = \|\mathbf{w}^* - \boldsymbol{\mu}_p\|_{\Sigma_p^{-1}}$
 - MGVT: **softmin** distance to the source weights $\varphi(\mathcal{W}_s) = \text{softmin}_{\mathbf{w} \in \mathcal{W}_s} (\|\mathbf{w}^* - \mathbf{w}\|)$

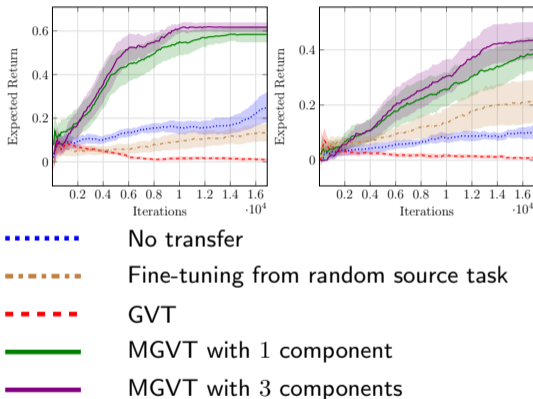


- Linear value functions with RBFs
- Different door positions
- $M = 10$ source tasks





- NN value function approximators
- 20 different mazes
- $M = 5$ source tasks
- Target *not* in the source tasks



Conclusion

We presented a general approach for transferring value functions in RL

- *No strong assumptions* on the approximators/distributions involved
- Exploration of the target task via *posterior sampling*
- Two practical and efficient algorithms

Future works

- transfer parameterized *policies*
- *active* exploration

Contacts







andrea.tirinzoni@polimi.it








<https://github.com/AndreaTirinzoni/>






References

- 
- Abel, D., Jinnai, Y., Guo, S. Y., Konidaris, G., and Littman, M. (2018).
Policy and value transfer in lifelong reinforcement learning.
In *International Conference on Machine Learning*, pages 20–29.
- 
- Alquier, P., Ridgway, J., and Chopin, N. (2016).
On the properties of variational approximations of gibbs posteriors.
Journal of Machine Learning Research, 17(239):1–41.
- 
- Amit, R. and Meir, R. (2018).
Meta-learning by adjusting priors based on extended PAC-Bayes theory.
In *Proceedings of the 35th International Conference on Machine Learning*.
- 
- Barreto, A., Borsa, D., Quan, J., Schaul, T., Silver, D., Hessel, M., Mankowitz, D., Zidek, A., and Munos, R. (2018).
Transfer in deep reinforcement learning using successor features and generalised policy improvement.
In Dy, J. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 501–510, Stockholm, Sweden. PMLR.


References (cont.)

- 
- Barreto, A., Dabney, W., Munos, R., Hunt, J. J., Schaul, T., van Hasselt, H. P., and Silver, D. (2017). Successor features for transfer in reinforcement learning. In *Advances in neural information processing systems*, pages 4055–4065.
- 
- Catoni, O. (2007). Pac-bayesian supervised classification: the thermodynamics of statistical learning. *arXiv preprint arXiv:0712.0248*.
- 
- Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*.
- 
- Grant, E., Finn, C., Levine, S., Darrell, T., and Griffiths, T. (2018). Recasting gradient-based meta-learning as hierarchical bayes. *arXiv preprint arXiv:1801.08930*.
- 
- Heess, N., Sriram, S., Lemmon, J., Merel, J., Wayne, G., Tassa, Y., Erez, T., Wang, Z., Eslami, A., Riedmiller, M., et al. (2017). Emergence of locomotion behaviours in rich environments. *arXiv preprint arXiv:1707.02286*.

References (cont.)

-  Lazaric, A. and Ghavamzadeh, M. (2010).
Bayesian multi-task reinforcement learning.
In ICML-27th International Conference on Machine Learning, pages 599–606. Omnipress.
-  OpenAI (2018).
Learning dexterous in-hand manipulation.
CoRR, abs/1808.00177.
-  Osband, I., Van Roy, B., and Wen, Z. (2014).
Generalization and exploration via randomized value functions.
arXiv preprint arXiv:1402.0635.
-  Tanaka, F. and Yamamura, M. (2003).
Multitask reinforcement learning on the distribution of mdps.
In Computational Intelligence in Robotics and Automation, 2003. Proceedings. 2003 IEEE International Symposium on, volume 3, pages 1108–1113. IEEE.
-  Taylor, M. E. and Stone, P. (2009).
Transfer learning for reinforcement learning domains: A survey.
Journal of Machine Learning Research, 10(Jul):1633–1685.

References (cont.)

- 
- Vinyals, O., Ewalds, T., Bartunov, S., Georgiev, P., Vezhnevets, A. S., Yeo, M., Makhzani, A., Küttler, H., Agapiou, J., Schrittwieser, J., et al. (2017).
Starcraft ii: A new challenge for reinforcement learning.
arXiv preprint arXiv:1708.04782.