

---

# What the hell is the title of this paper?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

The abstract paragraph should be indented  $\frac{1}{2}$  inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

## 1 Introduction

Recent advances have allowed reinforcement learning (RL) [] to achieve impressive results in a wide variety of complex tasks, ranging from Atari [], to the game of Go [], to the control of sophisticated robotics systems []. The main limitation is that RL algorithms still require an enormous amount of experience samples before successfully learning such complicated tasks. One of the most promising solutions is transfer learning, which focuses on reusing past knowledge available to the agent in order to reduce the sample-complexity for learning new tasks. In the typical settings of transfer in RL [], the agent is assumed to have already solved a set of *source tasks* generated from some unknown distribution. Then, given a *target task* drawn from the same distribution, or a slightly different one, the agent can rely on knowledge from the source tasks to speed-up the learning process. This constitutes a significant advantage over plain RL, where the agent learns each new task from scratch independently of previously learned tasks. Several algorithms have been proposed in the literature to transfer experience samples [], policies/options [], rewards [], value functions [], features [], and so on. We refer the reader to [] for a thorough survey on transfer in RL.

One of the most relevant problems in this context is how to efficiently explore the target task based on knowledge from the source tasks. Intuitively, assuming the tasks under consideration share some similarities due to the common distribution, much better exploration strategies than uninformed ones (e.g.,  $\epsilon$ -greedy) can be adopted for quickly learning the target task. Among the appealing approaches for this problem we find Bayesian methods (e.g., []), which are able to model the uncertainty over the current task based on previous knowledge and drive exploration so that this uncertainty is reduced as quickly as possible. Similarly, model-based algorithms (e.g., []) typically transfer samples to improve their estimates of the task model and adopt classic count-based exploration to drive the agent towards regions where such estimates are more uncertain. However, all these approaches either require strong assumptions (for example, on the distribution involved in Bayesian methods) or do not scale well to large problems. This greatly limits their practical applicability.

In this work, we tackle such limitations by considering a more general approach. Similarly to [], we assume tasks to share similarities in their value functions and use the given source tasks to learn the distribution over such functions. Then, we use this distribution as a prior for learning the target task and we propose an efficient variational approximation of the corresponding posterior. Leveraging on recent ideas from randomized value functions ([1]), we design a Thompson sampling-based algorithm which efficiently explores the target task by repeatedly sampling from the posterior and acting greedily w.r.t. (with respect to) the sampled value function. We show that our approach is very general, in the sense that it does not require any specific choice of function approximator or prior/posterior distribution models.

Maybe mention something about theory/experiments here

The rest of this document is organized as follows...

Complete this when the structure is defined or simply remove it to save space.

Transfer approaches to mention:

- Taylor 2009: survey
- Konidaris and Barto transfer shaped reward functions PROs: this allows the agent to have a more goal-directed behavior, thus limiting unnecessary random exploration CONs: applicable only to simple problems, does not scale
- Wilson 2007 propose a hierarchical Bayesian model for the distribution over tasks PROs it explicitly models the uncertainty over which tasks are being solved, thus quickly adapting to new ones and allowing informed exploration decisions to be taken
- Lazaric 2008 transfer samples PROs select good samples, can scale to continuous domains CONs eps greedy is used
- Taylor 2008 transfer samples for model-based RL PROs good exploration via model-based RL CONs does not scale, can negatively transfer when tasks are very different
- Lazaric 2010 propose a hierarchical bayesian model for the distribution over value functions PROs quickly adapts to new tasks, non-parameteric (GP) models CONs scaling, strong assumptions (GPTD)
- Fernandez and veloso 2006 propose an exploration strategy based on probabilistic policy reuse PROs good exploration (?) CONs Do not try to figure out which policies are good or bad
- Brunskill propose a method to transfer in model based RL (E3) - PROs theory, exploration CONs scaling
- Barreto use successor features PROs simple, theory CONs eps-greedy exploration

Exploration approaches to mention:

- Osband 2014 propose a method to efficiently explore via randomized value functions
- Osband 2016 adapt such algorithm to DQNs
- Houthoofd 2016 use variational inference to approximate the posterior distribution over parameters of the dynamics, and use that to drive exploration
- Azizzadenesheli 2018 extend Osband 2016 to use Bayesian DQN instead. Still makes Gaussian assumptions though

## 2 Preliminaries

We define a Markov decision process (MDP) as a tuple  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, p_0, \gamma \rangle$ , where  $\mathcal{S}$  is the state-space,  $\mathcal{A}$  is a finite set of actions,  $\mathcal{P}(\cdot|s, a)$  is the distribution of the next state  $s'$  given that action  $a$  is taken in state  $s$ ,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function,  $p_0$  is the initial-state distribution, and  $\gamma \in [0, 1)$  is the discount factor. We assume the reward function to be uniformly bounded by a constant  $R_{max} > 0$ . A deterministic policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  is a mapping from states to actions. At the beginning of each episode of interaction, the initial state  $s_0$  is drawn from  $p_0$ . Then, the agent takes the action  $a_0 = \pi(s_0)$ , receives a reward  $\mathcal{R}(s_0, a_0)$ , transitions to the next state  $s_1 \sim \mathcal{P}(\cdot|s_0, a_0)$ , and the process is repeated. The goal is to find the policy maximizing the long-term return over a possibly infinite horizon:  $\max_{\pi} J(\pi) \triangleq \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | \mathcal{M}, \pi]$ . To this end, we define the optimal value function  $Q^*(s, a)$  as the expected return obtained by taking action  $a$  in state  $s$  and following an optimal policy thereafter. Then, an optimal policy  $\pi^*$  is a policy that is greedy with respect to the optimal value function, i.e.,  $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$  for all states  $s$ . It can be shown (e.g., [1]) that  $Q^*$  is the unique fixed-point of the optimal Bellman operator  $T$  defined by  $TQ(s, a) = \mathcal{R}(s, a) + \gamma \mathbb{E}_{\mathcal{P}}[\max_{a'} Q(s', a')]$  for any value function  $Q$ . From now on, we

85 adopt the term  $Q$ -function to denote any plausible value function, i.e., any function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$   
 86 uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ .

87 When learning the optimal value function, a quantity of interest is how close a given  $Q$ -function  
 88 is to the fixed-point of the Bellman operator. This is given by its Bellman residual, defined by  
 89  $B(Q) \triangleq TQ - Q$ . Notice that  $Q$  is optimal if, and only if,  $B(Q)(s, a) = 0$  for all  $s, a$ . Furthermore,  
 90 if we assume the existence of a distribution  $\nu$  over  $\mathcal{S} \times \mathcal{A}$ , the squared Bellman error of  $Q$  is  
 91 defined as the expected squared Bellman residual of  $Q$  under  $\nu$ ,  $\|B(Q)\|_\nu^2 = \mathbb{E}_\nu [B^2(Q)]$ . Although  
 92 minimizing the empirical Bellman error is an appealing objective, it is well-known that an unbiased  
 93 estimator requires two independent samples of the next state  $s'$  of each  $s, a$  (e.g., []). In practice,  
 94 the empirical Bellman error is typically replaced by the TD error, which approximates the former  
 95 using a single transition sample. Given a dataset of  $N$  samples, the TD error is computed as  
 96  $\|B(Q)\|_D^2 = \frac{1}{N} \sum_{i=1}^N (r_i + \gamma \max_{a'} Q(s'_i, a') - Q(s_i, a_i))^2$ .

cite Maillard

### 97 3 Variational Transfer Learning

98 In this section, we describe our variational approach to transfer in RL. In Section 3.1, we start  
 99 by introducing our algorithm from a high-level perspective, in such a way that it can be used  
 100 for any choice of prior and posterior distributions. Then, in Sections 3.2 and 3.3, we propose  
 101 practical implementations based on Gaussian prior/posterior and mixture of Gaussian prior/posterior,  
 102 respectively.

#### 103 3.1 Algorithm

104 We begin with a simple consideration: the distribution  $\mathcal{D}$  over tasks clearly induces a distribution over  
 105 optimal  $Q$ -functions. Since, for any MDP, learning its optimal  $Q$ -function is sufficient for solving the  
 106 problem, one can safely replace the distribution over tasks with the distribution over their optimal  
 107 value functions. Furthermore, assume we know such distribution and we are given a new task  $\tau$  to  
 108 solve. Our goal is to design an algorithm that efficiently explores  $\tau$  so as to quickly adapt the prior  
 109 distribution in a Bayesian fashion to put all probability mass over the optimal  $Q$ -function of  $\tau$ .

110 We consider a parametric family of  $Q$ -functions,  $\mathcal{Q} = \{Q_{\mathbf{w}} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \mid \mathbf{w} \in \mathbb{R}^K\}$ . For simplic-  
 111 ity, we assume each function in  $\mathcal{Q}$  to be uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ . Then, we can reduce our prior  
 112 distribution over  $Q$ -functions to a prior distribution over weights  $p(\mathbf{w})$ . Assume that we are given  
 113 a dataset  $D = \{(s_i, a_i, s'_i, r_i) \mid i = 1, 2, \dots, N\}$  of samples from some task  $\tau$  that we want to solve.  
 114 Then, the posterior distribution over weights given such dataset can be computed by applying Bayes  
 115 theorem as  $p(\mathbf{w}|D) \propto p(D|\mathbf{w})p(\mathbf{w})$ . Unfortunately, this cannot be directly used in practice since  
 116 we do not have a model of the likelihood  $p(D|\mathbf{w})$ . In such case, it is very common to make strong  
 117 assumptions on the MDPs or the  $Q$ -functions so as to get tractable posteriors []. However, in our  
 118 transfer settings all distributions involved depend on the family of tasks under consideration, and  
 119 making such assumptions is likely to limit the applicability of the approach. Thus, we take a different  
 120 approach to derive a more general and meaningful solution. Recall that our final goal is move all  
 121 probability mass over the weights minimizing some empirical loss measure, which in our case is the  
 122 TD error  $\|B(\mathbf{w})\|_D^2$ . Then, given a prior  $p(\mathbf{w})$ , we know from PAC-Bayesian theory that the optimal  
 123 Gibbs posterior takes the form []:

Cite some-  
body

Cite Catoni  
2007

$$q(\mathbf{w}) = \frac{e^{-\Lambda \|B(\mathbf{w})\|_D^2} p(\mathbf{w})}{\int e^{-\Lambda \|B(\mathbf{w}')\|_D^2} p(d\mathbf{w}')} \quad (1)$$

124 for some parameter  $\Lambda > 0$ . Since  $\Lambda$  is typically chosen to increase with the number of samples  
 125  $N$ , in the remaining we set it to  $\lambda^{-1}N$ , for some constant  $\lambda > 0$ . Notice that, whenever the term  
 126  $e^{-\Lambda \|B(\mathbf{w})\|_D^2}$  can be interpreted as the actual likelihood of  $D$ ,  $q$  becomes a classic Bayesian posterior.  
 127 Although we now have an appealing distribution, the integral at the denominator of (1) is intractable  
 128 to compute even for simple  $Q$ -function models. Thus, we propose a variational approximation  $q_\xi$  by  
 129 considering a simpler family of distributions parameterized by  $\xi \in \Xi$ . Then, our problem reduces to  
 130 finding the variational parameters  $\xi$  such that  $q_\xi$  minimizes the Kullback-Leibler (KL) divergence  
 131 w.r.t. the Gibbs posterior  $q$ . From the theory of variational inference (e.g., []), this can be shown to be

cite

<sup>1</sup>In practice, this is easily achieved by truncation.

---

**Algorithm 1** Variational Transfer

---

**Require:** Target task  $\tau$ , source  $Q$ -function weights  $\mathcal{W}_s$ , batch sizes  $M_D$  and  $M_{\mathcal{W}}$ , prior weight  $\lambda$

---

```
1: Estimate prior  $p(\mathbf{w})$  from  $\mathcal{W}_s$ 
2: Initialize variational parameters:  $\xi \leftarrow \arg\min_{\xi} KL(q_{\xi} || p)$ 
3: Initialize replay buffer:  $D = \emptyset$ 
4: repeat
5:   Sample initial state:  $s_0 \sim p_0^{(\tau)}$ 
6:   while  $s_h$  is not terminal do
7:     Sample weights:  $\mathbf{w} \sim q_{\xi}(\mathbf{w})$ 
8:     Take action  $a_h = \arg\max_a Q_{\mathbf{w}}(s_h, a)$ 
9:     Observe transition  $s_{h+1} \sim \mathcal{P}^{(\tau)}(\cdot | s_h, a_h)$ 
10:    Collect reward  $r_h = \mathcal{R}^{(\tau)}(s_h, a_h)$ 
11:    Add sample to the replay buffer:  $D \leftarrow D \cup \langle s_h, a_h, r_h, s_{h+1} \rangle$ 
12:    Sample batch  $D' = \langle s_i, a_i, r_i, s'_i \rangle_{i=1}^{M_D}$  from  $D$  and  $\mathcal{W} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{M_{\mathcal{W}}}\}$  from  $q_{\xi}$ 
13:    Approximate objective:  $\mathcal{L}(\xi) = \frac{1}{M_{\mathcal{W}}} \sum_{\mathbf{w} \in \mathcal{W}} \|B(\mathbf{w})\|_{D'}^2 - \frac{\lambda}{N} KL(q_{\xi} || p)$ 
14:    Compute the gradient  $\nabla_{\xi} \mathcal{L}(\xi)$ 
15:    Update  $\xi$  in the direction of  $\nabla_{\xi} \mathcal{L}(\xi)$  using any stochastic optimizer (e.g., ADAM)
16:  end while
17: until forever
```

---

132 equivalent to minimizing the well-known (negative) *evidence lower bound* (ELBO):

$$\min_{\xi \in \Xi} \mathcal{L}(\xi) = \mathbb{E}_{\mathbf{w} \sim q_{\xi}} \left[ \|B(\mathbf{w})\|_D^2 \right] - \frac{\lambda}{N} KL(q_{\xi}(\mathbf{w}) || p(\mathbf{w})) \quad (2)$$

133 Intuitively, the approximate posterior trades-off between placing probability mass over those weights  
134  $\mathbf{w}$  that have low TD error (first term), and staying close to the prior distribution (second term).  
135 Assuming that we are able to compute the gradients of (2) w.r.t. the variational parameters  $\xi$ , our  
136 objective can be easily optimized with any stochastic optimization algorithm.

137 We now highlight our general transfer procedure in Alg. 1, while deferring a description of practical  
138 implementations with specific choices for the distributions involved to the next two sections. Given a  
139 set of weights  $\mathcal{W}_s$  from the source tasks, we start by estimating the prior distribution (line 1) and we  
140 initialize the variational parameters by minimizing the KL divergence w.r.t. such distribution<sup>2</sup> (line  
141 2). Then, at each time step of interaction, we re-sample the weights from the current approximate  
142 posterior and act greedily w.r.t. the corresponding  $Q$ -function (lines 7,8). This resembles the well-  
143 known Thompson sampling approach adopted in multi-armed bandits [1] and allows our algorithm to  
144 efficiently explore the target task. In some sense, at each time we guess what is the task we are trying  
145 to solve based on our current belief and we act as if such guess were actually true. After collecting  
146 and storing the new experience (lines 9-11), we draw a batch of samples from the replay buffer and  
147 a batch of weights from the posterior (line 12). We use these to approximate the negative ELBO,  
148 compute its gradient, and finally update the variational parameters (lines 13-15).

149 The main advantage of our approach is that it exploits knowledge from the source tasks to perform an  
150 efficient adaptive exploration. Intuitively, during the first steps of interaction, our algorithm has no  
151 idea about what is the current task. However, it can rely on the learned prior to take early informed  
152 decisions. As the learning process goes on, it will quickly figure out which task is being solved, thus  
153 moving all probability mass over the weights minimizing the TD error. From that point, sampling  
154 from the posterior is approximately equivalent to deterministically taking such weights, and no more  
155 exploration will be performed. Finally, notice the generality of the proposed approach: as far as  
156 the objective  $\mathcal{L}$  is differentiable in the variational parameters  $\xi$ , and its gradients can be efficiently  
157 computed, any function approximator for the  $Q$ -functions and any family for the prior and posterior  
158 distributions can be adopted. For the latter, we describe two practical choices in the next two sections.

---

<sup>2</sup>If the prior and approximate posterior were in the same family of distributions we could simply set  $\xi$  to the prior parameters, however this does not always hold in practice.

Cite

### 3.2 Gaussian Variational Transfer

We now restrict ourselves to a specific choice of the prior and posterior families that makes our algorithm very efficient and easy to implement. We assume that optimal  $Q$ -functions according to our task distribution (or better, their weights) follow a multivariate Gaussian law. That is, we model the prior as  $p(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$  and we learn its parameters from the set of source weights using, e.g., maximum likelihood estimation (with small regularization to make sure the covariance is positive definite). Then, our variational family is the set of all well-defined Gaussian distributions, i.e., the variational parameters are  $\Xi = \{(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mid \boldsymbol{\mu} \in \mathbb{R}^K, \boldsymbol{\Sigma} \in \mathbb{R}^{K \times K}, \boldsymbol{\Sigma} \succ 0\}$ . To prevent the covariance from going not positive definite, we consider its Cholesky decomposition  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$  and learn the lower-triangular Cholesky factor  $\mathbf{L}$  instead. Under Gaussian distributions, all quantity of interest for using Alg. 1 can be computed very easily. The KL divergence between the prior and approximate posterior can be computed in closed-form as:

$$KL(q_{\xi}(\mathbf{w}) \parallel p(\mathbf{w})) = \frac{1}{2} \left( \log \frac{|\boldsymbol{\Sigma}_p|}{|\boldsymbol{\Sigma}|} + \text{Tr}(\boldsymbol{\Sigma}_p^{-1}\boldsymbol{\Sigma}) + (\boldsymbol{\mu} - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_p^{-1}(\boldsymbol{\mu} - \boldsymbol{\mu}_p) - K \right) \quad (3)$$

for  $\xi = (\boldsymbol{\mu}, \mathbf{L})$  and  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$ . Its gradients with respect to the variational parameters are []:

$$\nabla_{\boldsymbol{\mu}} KL(q_{\xi}(\mathbf{w}) \parallel p(\mathbf{w})) = \boldsymbol{\Sigma}_p^{-1}(\boldsymbol{\mu} - \boldsymbol{\mu}_p) \quad (4)$$

Cite matrix cookbook

$$\nabla_{\mathbf{L}} KL(q_{\xi}(\mathbf{w}) \parallel p(\mathbf{w})) = \boldsymbol{\Sigma}_p^{-1}\mathbf{L} - (\mathbf{L}^{-1})^T \quad (5)$$

Finally, the gradients w.r.t. the expected likelihood term of the variational objective (2) can be computed using the reparameterization trick (e.g., []):

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^T)} [\|B(\mathbf{w})\|_D^2] = \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\nabla_{\mathbf{w}} \|B(\mathbf{w})\|_D^2] \quad \text{for } \mathbf{w} = \mathbf{L}\mathbf{v} + \boldsymbol{\mu} \quad (6)$$

Cite deepmind and another

$$\nabla_{\mathbf{L}} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^T)} [\|B(\mathbf{w})\|_D^2] = \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\nabla_{\mathbf{w}} \|B(\mathbf{w})\|_D^2 \cdot \mathbf{v}^T] \quad \text{for } \mathbf{w} = \mathbf{L}\mathbf{v} + \boldsymbol{\mu} \quad (7)$$

### 3.3 Mixture of Gaussian Variational Transfer

Although the Gaussian assumption of the previous section is very appealing as it allows for a simple and efficient way of computing the variational objective and its gradients, we believe that such assumption almost never holds in practice. In fact, even for families of tasks in which the reward and transition models follow a Gaussian law, the  $Q$ -values might be far from it. Depending on the family of tasks under consideration and, since we are learning a distribution over weights, on the chosen function approximator, the prior might have arbitrarily complex shapes. When the information loss due to the Gaussian approximation becomes too severe, the algorithm is likely to fail at transferring knowledge, thus reducing to almost random exploration. We now propose a variant to successfully solve this problem, while keeping the algorithm simple and efficient enough to be applied in practice. In order to capture arbitrarily complex distributions, we use a kernel estimator []for learning our prior.

Assume we are given a set  $\mathcal{W}_s$  of weights from the source tasks. Then, our estimated prior places a single isotropic Gaussian over each weight:  $p(\mathbf{w}) = \frac{1}{|\mathcal{W}_s|} \sum_{\mathbf{w}_s \in \mathcal{W}_s} \mathcal{N}(\mathbf{w} | \mathbf{w}_s, \sigma_p^2 \mathbf{I})^3$ . This takes the form of a mixture of Gaussians with equally weighted components. Consistently with the prior, we model our approximate posterior as a mixture of Gaussians. However, we allow a different number of components (typically much less than the prior's) and we adopt full covariances instead of only diagonals, so that our posterior has the potential to match complex distributions with less components. Using  $C$  components, our posterior is  $q_{\xi}(\mathbf{w}) = \frac{1}{C} \sum_{i=1}^C \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ , with variational parameters  $\xi = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_C, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_C)$ . Once again, we learn Cholesky factors instead of full covariances.

Although this new model has the potential to capture much more complex distributions, it poses a major complication: the KL divergence between two mixture of Gaussians is well-known to have no closed-form equation. To solve this issue, we can rely on an upper bound to such quantity, so that negative ELBO we are optimizing still represents an upper bound on the KL between the approximate and true posterior. However, this turns out to be non-trivial as well. In fact, it is very easy to bound the KL between two mixtures with the KLs between each couple of components. However, the loss of information is such that minimizing the upper bound via gradient methods converges to a local optimum in which all components tend to go to the same point, thus almost reducing to the single

203 Gaussian case. To solve this issue, we adopt the variational upper bound proposed in [], which we  
 204 found to be able to preserve the needed information. We report it here for the sake of completeness.  
 205 See the original paper for the proof.

Cite UB

206 **Theorem 1.** Let  $p = \sum_i c_i^{(p)} f_i^{(p)}$  and  $q = \sum_j c_j^{(q)} f_j^{(q)}$  be two mixture of Gaussian distributions,  
 207 where  $f_i^{(p)} = \mathcal{N}(\mu_i^{(p)}, \Sigma_i^{(p)})$  denotes the  $i$ -th component of  $p$ ,  $c_i^{(p)}$  denotes its weight, and similarly  
 208 for  $q$ . Introduce two vectors  $\chi^{(1)}$  and  $\chi^{(2)}$  such that  $c_i^{(p)} = \sum_j \chi_{j,i}^{(2)}$  and  $c_j^{(q)} = \sum_i \chi_{i,j}^{(1)}$ . Then:

$$KL(p||q) \leq KL(\chi^{(2)}||\chi^{(1)}) + \sum_{i,j} \chi_{j,i}^{(2)} KL(f_i^{(p)}||f_j^{(q)}) \quad (8)$$

209 Our new algorithm replaces the KL with the above-mentioned upper bound. Each time we require its  
 210 value, we have to recompute the parameters  $\chi^{(1)}$  and  $\chi^{(2)}$  that tighten the bound. As shown in [],  
 211 this can be achieved by a simple fixed-point procedure. Furthermore, both terms in the approximate  
 212 negative ELBO are now linear combinations of functions of the variational parameters for different  
 213 components, thus their gradients can be straightforwardly derived from the ones of the Gaussian case.

cite them  
again

### 214 3.4 Optimizing the TD error

215 From Sections 3.2 and 3.3, we know that differentiating the negative ELBO  $\mathcal{L}$  w.r.t.  $\xi$  requires  
 216 differentiating  $\|B(\mathbf{w})\|_D^2$  w.r.t.  $\mathbf{w}$ . Unfortunately, the TD error is well-known to be non-differentiable  
 217 due to the presence of the max operator. This rarely represents a problem since typical value-based  
 218 algorithms are actually semi-gradient methods, i.e., they do not differentiate the targets (see, e.g.,  
 219 Chapter 11 of []). However, our transfer settings are rather different than common RL. In fact, our  
 220 algorithm is likely to always start from  $Q$ -functions that are very close to the optimum, and the  
 221 only thing that needs to be done is to adapt the weights in a direction of lower error (i.e., higher  
 222 likelihood) so as to quickly converge to the solution of the task that is being solved. Unfortunately,  
 223 this property cannot be guaranteed for most semi-gradient algorithms. Even worse, many online  
 224 RL algorithms combined with complex function approximators (e.g., DQNs) are well-known to be  
 225 unstable, especially when approaching the optimum, and require a lot of tuning and tricks to work  
 226 well. This is obviously an undesirable property in our case, as we only aim at adapting already good  
 227 solutions. Thus, we consider using a residual gradient algorithm (after []). In order to differentiate  
 228 the targets, we replace the optimal Bellman operator with the mellow Bellman operator introduced in  
 229 [], which adopts a softened version of max called *mellowmax*:

Cite Sutton

Baird

Cite MM

$$\text{mm}_a Q_{\mathbf{w}}(s, a) = \frac{1}{\kappa} \log \frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \quad (9)$$

230 where  $\kappa$  is a hyperparameter and  $|\mathcal{A}|$  is the number of actions. The mellow Bellman operator, which  
 231 we denote as  $\tilde{T}$ , has several appealing properties that make it suitable for our settings: (i) it converges  
 232 to the maximum as  $\kappa \rightarrow \infty$ , (ii) it has a unique fixed point, and (iii) it is *differentiable*. Denoting  
 233 by  $\tilde{B}(\mathbf{w}) = \tilde{T}Q_{\mathbf{w}} - Q_{\mathbf{w}}$  the Bellman residual w.r.t. the mellow Bellman operator  $\tilde{T}$ , we have  
 234 that the corresponding TD error,  $\|\tilde{B}(\mathbf{w})\|_D^2$ , is now differentiable with respect to  $\mathbf{w}$ . Although  
 235 residual algorithms have the sought guaranteed convergence, they are typically much slower than  
 236 their semi-gradient counterpart. This problem was addressed in [], where the authors proposed a  
 237 simple remedy consisting in projecting the gradient in a direction that still guarantees convergence but  
 238 is also closer to the one of the semi-gradient, thus achieving higher learning speed. This can be easily  
 239 done by including a parameter  $\psi \in [0, 1]$  in the TD error gradient such that:

Baird again

$$\nabla_{\mathbf{w}} \|\tilde{B}(\mathbf{w})\|_D^2 = \frac{2}{N} \sum_{i=1}^N b_i(\mathbf{w}) \left( \gamma \psi \nabla_{\mathbf{w}} \text{mm}_{a'} Q_{\mathbf{w}}(s'_i, a') - \nabla_{\mathbf{w}} Q_{\mathbf{w}}(s_i, a_i) \right) \quad (10)$$

240 where  $b_i(\mathbf{w}) = r_i + \gamma \text{mm}_{a'} Q_{\mathbf{w}}(s'_i, a') - Q_{\mathbf{w}}(s_i, a_i)$ . Notice that  $\psi$  trades-off between the semi-  
 241 gradient ( $\psi = 0$ ) and the full residual gradient ( $\psi = 1$ ). A good criterion for choosing such  
 242 parameter is to start with values close to zero (to have faster learning) and move to higher values  
 243 when approaching the optimum (to guarantee converge). Since in our case we are likely to start from  
 244 this latter point, we consider only higher values (e.g., above 0.5).

<sup>3</sup>Notice that this is slightly different than the typical kernel estimator (e.g., [])



## 4 Theoretical Analysis

In this section, we theoretically analyze our variational transfer algorithm...

A first important question that we need to answer is whether replacing max with mellow-max in the Bellman operator constitutes a strong approximation or not. It has been proved [] that the mellow Bellman operator is a contraction under the  $L_\infty$ -norm and, thus, has a unique fixed-point. However, how such fixed-point differs from the one of the optimal Bellman operator remains an open question. Since mellow-max monotonically converges to max as  $\kappa \rightarrow \infty$ , it would be desirable if the corresponding operator also monotonically converged to the optimal one. We confirm that this property actually holds in the following theorem.

Cite MM

**Theorem 2.** *Let  $V$  be the fixed-point of the optimal Bellman operator  $T$ , and  $Q$  the corresponding action-value function. Define the action-gap function  $g(s)$  as the difference between the value of the best action and the second best action at each state  $s$ . Let  $\tilde{V}$  be the fixed-point of the mellow Bellman operator  $\tilde{T}$  with parameter  $\kappa > 0$  and denote by  $\beta > 0$  the inverse temperature of the induced Boltzmann distribution (as in []). Let  $\nu$  be a probability measure over the state-space. Then, for any  $p \geq 1$ :*

Cite MM

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{2R_{max}}{(1-\gamma)^2} \left\| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g}} \right\|_{\nu,p}^p \quad (11)$$

## 5 Related Works

## 6 Experiments

### 6.1 Gridworld

### 6.2 Classic Control

### 6.3 Maze Navigation

## 7 Conclusion

## References

- [1] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1994.

## 269 A Proofs

270 **Theorem 2.** Let  $V$  be the fixed-point of the optimal Bellman operator  $T$ , and  $Q$  the corresponding  
 271 action-value function. Define the action-gap function  $g(s)$  as the difference between the value of  
 272 the best action and the second best action at each state  $s$ . Let  $\tilde{V}$  be the fixed-point of the mellow  
 273 Bellman operator  $\tilde{T}$  with parameter  $\kappa > 0$  and denote by  $\beta > 0$  the inverse temperature of the  
 274 induced Boltzmann distribution (as in []). Let  $\nu$  be a probability measure over the state-space. Then,  
 275 for any  $p \geq 1$ :

Cite MM

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{2R_{max}}{(1-\gamma)^2} \left\| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g}} \right\|_{\nu,p}^p \quad (11)$$

276 *Proof.* We begin by noticing that:

$$\begin{aligned} \|V - \tilde{V}\|_{\nu,p}^p &= \|TV - \tilde{T}\tilde{V}\|_{\nu,p}^p \\ &= \|TV - \tilde{T}V + \tilde{T}V - \tilde{T}\tilde{V}\|_{\nu,p}^p \\ &\leq \|TV - \tilde{T}V\|_{\nu,p}^p + \|\tilde{T}V - \tilde{T}\tilde{V}\|_{\nu,p}^p \\ &\leq \|TV - \tilde{T}V\|_{\nu,p}^p + \gamma \|V - \tilde{V}\|_{\nu,p}^p \end{aligned}$$

277 where the first inequality follows from Minkowsky's inequality and the second one from the contrac-  
 278 tion property of the mellow Bellman operator. This implies that:

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{1}{1-\gamma} \|TV - \tilde{T}V\|_{\nu,p}^p \quad (12)$$

279 Let us bound the norm on the right-hand side separately. In order to do that, we will bound the  
 280 function  $|TV(s) - \tilde{T}V(s)|$  point-wisely for any state  $s$ . By applying the definition of the optimal  
 281 and mellow Bellman operators, we obtain:

$$\begin{aligned} |TV(s) - \tilde{T}V(s)| &= \left| \max_a \{R(s, a) + \gamma \mathbb{E}[V(s')]\} - \min_a \{R(s, a) + \gamma \mathbb{E}[V(s')]\} \right| \\ &= \left| \max_a Q(s, a) - \min_a Q(s, a) \right| \end{aligned}$$

282 Recall that applying the mellow-max is equivalent to computing an expectation under a Boltzmann  
 283 distribution with inverse temperature  $\beta$  induced by  $\kappa$  []. Thus, we can write:

Cite MM

$$\begin{aligned} \left| \max_a Q(s, a) - \min_a Q(s, a) \right| &= \left| \sum_a \pi^*(a|s) Q(s, a) - \sum_a \pi_\beta(a|s) Q(s, a) \right| \\ &= \left| \sum_a Q(s, a) (\pi^*(a|s) - \pi_\beta(a|s)) \right| \\ &\leq \sum_a |Q(s, a)| |\pi^*(a|s) - \pi_\beta(a|s)| \\ &\leq \frac{R_{max}}{1-\gamma} \sum_a |\pi^*(a|s) - \pi_\beta(a|s)| \end{aligned} \quad (13)$$

284 where  $\pi^*$  is the optimal (deterministic) policy w.r.t.  $Q$  and  $\pi_\beta$  is the Boltzmann distribution induced  
 285 by  $Q$  with inverse temperature  $\beta$ :

$$\pi_\beta(a|s) = \frac{e^{\beta Q(s,a)}}{\sum_{a'} e^{\beta Q(s,a')}}$$



286 Denote by  $a_1(s)$  the optimal action for state  $s$  under  $Q$ . We can then write:

$$\begin{aligned}
\sum_a |\pi^*(a|s) - \pi_\beta(a|s)| &= |\pi^*(a_1(s)|s) - \pi_\beta(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi^*(a|s) - \pi_\beta(a|s)| \\
&= |1 - \pi_\beta(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi_\beta(a|s)| \\
&= 2 |1 - \pi_\beta(a_1(s)|s)|
\end{aligned} \tag{14}$$

287 Finally, let us bound this last term:

$$\begin{aligned}
|1 - \pi_\beta(a_1(s)|s)| &= \left| 1 - \frac{e^{\beta Q(s, a_1(s))}}{\sum_{a'} e^{\beta Q(s, a')}} \right| \\
&= \left| 1 - \frac{e^{\beta(Q(s, a_1(s)) - Q(s, a_2(s)))}}{\sum_{a'} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&= \left| 1 - \frac{e^{\beta g(s)}}{\sum_{a'} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&= \left| 1 - \frac{e^{\beta g(s)}}{e^{\beta g(s)} + \sum_{a' \neq a_1(s)} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&\leq \left| 1 - \frac{e^{\beta g(s)}}{e^{\beta g(s)} + |\mathcal{A}|} \right| \\
&= \left| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g(s)}} \right|
\end{aligned} \tag{15}$$

288 Combining Eq. (13), (14), and (15), we obtain:

$$\left| \max_a Q(s, a) - \min_a Q(s, a) \right| \leq \frac{2R_{max}}{1 - \gamma} \left| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g(s)}} \right|$$

289 Taking the norm and plugging this into Eq. (12) concludes the proof.  $\square$

290 **Lemma 1.** Let  $p$  and  $\nu$  denote probability measures over  $Q$ -functions and state-action pairs, respectively. Assume  $Q^*$  is the unique fixed-point of the optimal Bellman operator  $T$ . Then, for any  $\delta > 0$ ,  
291 with probability at least  $1 - \delta$  over the choice of a  $Q$ -function  $Q$ , the following holds:  
292

$$\|Q - Q^*\|_\nu^2 \leq \frac{\mathbb{E}_p \left[ \|B(Q)\|_\nu^2 \right]}{(1 - \gamma)\delta} \tag{16}$$

293 *Proof.* First notice that:

$$\begin{aligned}
\|Q - Q^*\| &= \|Q + TQ - TQ - TQ^*\| \\
&\leq \|Q - TQ\| + \|TQ - TQ^*\| \\
&\leq \|Q - TQ\| + \gamma \|Q - Q^*\| \\
&= \|B(Q)\| + \gamma \|Q - Q^*\|
\end{aligned}$$

294 which implies that:

$$\|Q - Q^*\| \leq \frac{1}{1 - \gamma} \|B(Q)\|$$

295 Then we can write:

$$P(\|Q - Q^*\| > \epsilon) \leq P(\|B(Q)\| > \epsilon(1 - \gamma)) \leq \frac{\mathbb{E}_p \left[ \|B(Q)\|_\nu^2 \right]}{(1 - \gamma)\epsilon}$$

296 Settings the right-hand side equal to  $\delta$  and solving for  $\epsilon$  concludes the proof.  $\square$

297 **Corollary 1.** Let  $p$  and  $\nu$  denote probability measures over  $Q$ -functions and state-action pairs,  
 298 respectively. Assume  $\tilde{Q}$  is the unique fixed-point of the mellow Bellman operator  $\tilde{T}$ . Then, for any  
 299  $\delta > 0$ , with probability at least  $1 - \delta$  over the choice of a  $Q$ -function  $Q$ , the following holds:

$$\|Q - \tilde{Q}\|_\nu^2 \leq \frac{\mathbb{E}_p \left[ \|\tilde{B}(Q)\|_\nu^2 \right]}{(1 - \gamma)\delta} \quad (17)$$

300 **Lemma 2.** Assume  $Q$ -functions belong to a parametric space of functions bounded by  $\frac{R_{max}}{1-\gamma}$ . Let  $p$   
 301 and  $q$  be arbitrary distributions over the parameter space  $\mathcal{W}$ , and  $\nu$  be a probability measure over  
 302  $\mathcal{S} \times \mathcal{A}$ . Consider a dataset  $D$  of  $N$  samples and define  $v(\mathbf{w}) \triangleq \mathbb{E}_\nu [\text{Var}_{\mathcal{P}} [b(\mathbf{w})]]$ . Then, for any  
 303  $\delta > 0$ , with probability at least  $1 - \delta$ , the following two inequalities hold simultaneously:

$$\mathbb{E}_q \left[ \|B(\mathbf{w})\|_\nu^2 \right] \leq \mathbb{E}_q \left[ \|B(\mathbf{w})\|_D^2 \right] - \mathbb{E}_q [v(\mathbf{w})] + \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (18)$$

304

$$\mathbb{E}_q \left[ \|B(\mathbf{w})\|_D^2 \right] \leq \mathbb{E}_q \left[ \|B(\mathbf{w})\|_\nu^2 \right] + \mathbb{E}_q [v(\mathbf{w})] + \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (19)$$

305 *Proof.* From Hoeffding's inequality we have:

$$P \left( \left| \mathbb{E}_{\nu, \mathcal{P}} \left[ \|B(\mathbf{w})\|_D^2 \right] - \|B(\mathbf{w})\|_D^2 \right| > \epsilon \right) \leq 2 \exp \left( - \frac{2N\epsilon^2}{\left( 2 \frac{R_{max}}{1-\gamma} \right)^4} \right)$$

306 which implies that, for any  $\delta > 0$ , with probability at least  $1 - \delta$ :

$$\left| \mathbb{E}_{\nu, \mathcal{P}} \left[ \|B(\mathbf{w})\|_D^2 \right] - \|B(\mathbf{w})\|_D^2 \right| \leq 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}}$$

307 Under independence assumptions, the expected TD error can be re-written as:

$$\begin{aligned} \mathbb{E}_{\nu, \mathcal{P}} \left[ \|B(\mathbf{w})\|_D^2 \right] &= \mathbb{E}_{\nu, \mathcal{P}} \left[ \frac{1}{N} \sum_{i=1}^N (r_i + \gamma \min_{a'} Q_{\mathbf{w}}(s'_i, a') - Q_{\mathbf{w}}(s_i, a_i))^2 \right] \\ &= \mathbb{E}_{\nu, \mathcal{P}} \left[ (R(s, a) + \gamma \min_{a'} Q_{\mathbf{w}}(s', a') - Q_{\mathbf{w}}(s, a))^2 \right] \\ &= \mathbb{E}_\nu \left[ \mathbb{E}_{\mathcal{P}} [b(\mathbf{w})^2] \right] \\ &= \mathbb{E}_\nu \left[ \text{Var}_{\mathcal{P}} [b(\mathbf{w})] + \mathbb{E}_{\mathcal{P}} [b(\mathbf{w})]^2 \right] \\ &= v(\mathbf{w}) + \|B(\mathbf{w})\|_\nu^2 \end{aligned}$$

308 where  $v(\mathbf{w}) \triangleq \mathbb{E}_\nu [\text{Var}_{\mathcal{P}} [b(\mathbf{w})]]$ . Thus:

$$\left| \|B(\mathbf{w})\|_\nu^2 + v(\mathbf{w}) - \|B(\mathbf{w})\|_D^2 \right| \leq 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (20)$$

309 From the change of measure inequality [], we have that, for any measurable function  $f(\mathbf{w})$  and any  
 310 two probability measures  $p$  and  $q$ :

$$\log \mathbb{E}_p \left[ e^{f(\mathbf{w})} \right] \geq \mathbb{E}_q [f(\mathbf{w})] - KL(q||p)$$

311 Thus, multiplying both sides of (20) by  $\lambda^{-1}N$  and applying the change of measure inequality with  
 312  $f(\mathbf{w}) = \lambda^{-1}N \left| \|B(\mathbf{w})\|_\nu^2 + v(\mathbf{w}) - \|B(\mathbf{w})\|_D^2 \right|$ , we obtain:

$$\mathbb{E}_q [f(\mathbf{w})] - KL(q||p) \leq \log \mathbb{E}_p \left[ e^{f(\mathbf{w})} \right] \leq 4 \frac{R_{max}^2 \lambda^{-1}N}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}}$$

Find a ref-  
erence for  
this

where the second inequality holds since the right-hand side of (20) does not depend on  $\mathbf{w}$ . Finally, we can explicitly write:

$$\mathbb{E}_q \left[ \left\| B(\mathbf{w}) \right\|_\nu^2 + v(\mathbf{w}) - \left\| B(\mathbf{w}) \right\|_D^2 \right] \leq \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}}$$

from which the lemma follows straightforwardly.  $\square$

**Lemma 3.** Let  $p$  be a prior distribution over the parameter space  $\mathcal{W}$ , and  $\nu$  be a probability measure over  $\mathcal{S} \times \mathcal{A}$ . Assume  $\hat{\xi}$  is the minimizer of  $ELBO(\xi) = \mathbb{E}_{q_\xi} \left[ \left\| B(\mathbf{w}) \right\|_D^2 \right] + \frac{\lambda}{N} KL(q_\xi||p)$  for a dataset  $D$  of  $N$  samples. Define  $v(\mathbf{w}) \triangleq \mathbb{E}_\nu [Var_{\mathcal{P}} [b(\mathbf{w})]]$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ :

$$\mathbb{E}_{q_{\hat{\xi}}} \left[ \left\| B(\mathbf{w}) \right\|_\nu^2 \right] \leq \inf_{\xi \in \Xi} \left\{ \mathbb{E}_{q_\xi} \left[ \left\| B(\mathbf{w}) \right\|_\nu^2 \right] + \mathbb{E}_{q_\xi} [v(\mathbf{w})] + 2 \frac{\lambda}{N} KL(q_\xi||p) \right\} + 2 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{N}}$$

*Proof.* Let us use Lemma 2 for the specific choice  $q = q_{\hat{\xi}}$ . From Eq. (18), we have:

$$\begin{aligned} \mathbb{E}_{q_{\hat{\xi}}} \left[ \left\| B(\mathbf{w}) \right\|_\nu^2 \right] &\leq \mathbb{E}_{q_{\hat{\xi}}} \left[ \left\| B(\mathbf{w}) \right\|_D^2 \right] - \mathbb{E}_{q_{\hat{\xi}}} [v(\mathbf{w})] + \frac{\lambda}{N} KL(q_{\hat{\xi}}||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \\ &\leq \mathbb{E}_{q_{\hat{\xi}}} \left[ \left\| B(\mathbf{w}) \right\|_D^2 \right] + \frac{\lambda}{N} KL(q_{\hat{\xi}}||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \\ &= \inf_{\xi \in \Xi} \left\{ \mathbb{E}_{q_\xi} \left[ \left\| B(\mathbf{w}) \right\|_D^2 \right] + \frac{\lambda}{N} KL(q_\xi||p) \right\} + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \end{aligned}$$

where the second inequality holds since  $v(\mathbf{w}) > 0$ , while the equality holds from the definition of  $\hat{\xi}$ . We can now use Eq. (19) to bound  $\mathbb{E}_{q_{\hat{\xi}}} \left[ \left\| B(\mathbf{w}) \right\|_D^2 \right]$ , thus obtaining:

$$\mathbb{E}_{q_{\hat{\xi}}} \left[ \left\| B(\mathbf{w}) \right\|_\nu^2 \right] \leq \inf_{\xi \in \Xi} \left\{ \mathbb{E}_{q_\xi} \left[ \left\| B(\mathbf{w}) \right\|_\nu^2 \right] + \mathbb{E}_{q_\xi} [v(\mathbf{w})] + 2 \frac{\lambda}{N} KL(q_\xi||p) \right\} + 2 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{N}}$$

This concludes the proof.  $\square$