

---

# Formatting instructions for NIPS 2018

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

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

## 5 1 Introduction

## 6 2 Preliminaries

### 7 2.1 Markov Decision Processes

8 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  
9 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  
10 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  
11 distribution, and  $\gamma \in [0, 1)$  is the discount factor. We assume the reward function to be uniformly  
12 bounded by a constant  $R_{max} > 0$ . A deterministic policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  is a mapping from states  
13 to actions. At the beginning of each episode of interaction, the initial state  $s_0$  is drawn from  $p_0$ .  
14 Then, the agent takes the action  $a_0 = \pi(s_0)$ , receives a reward  $\mathcal{R}(s_0, a_0)$ , transitions to the next  
15 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  
16 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,  
17 we define the optimal value function  $Q^*(s, a)$  as the expected return obtained by taking action  $a$   
18 in state  $s$  and following an optimal policy thereafter. Then, an optimal policy  $\pi^*$  is a policy that  
19 is greedy with respect to the optimal value function, i.e.,  $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$  for all states  
20  $s$ . It can be shown (e.g., [1]) that  $Q^*$  is the unique fixed-point of the optimal Bellman operator  $T$   
21 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  
22 adopt the term  $Q$ -function to denote any plausible value function, i.e., any function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$   
23 uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ .

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

cite Maillard

## 34 2.2 Variational Inference

35 When working with Bayesian approaches, the posterior distribution of hidden variables  $\mathbf{w} \in \mathbb{R}^K$   
 36 given data  $D$ ,

$$p(\mathbf{w}|D) = \frac{p(D|\mathbf{w})p(\mathbf{w})}{p(D)} = \frac{p(D|\mathbf{w})p(\mathbf{w})}{\int_{\mathbf{w}} p(D|\mathbf{w})p(\mathbf{w})}, \quad (1)$$

37 is typically intractable for many models of interest (e.g., when working with deep neural networks)  
 38 due to difficulties in computing the integral of Eq. (1). The main intuition behind variational inference  
 39 [] is to approximate the intractable posterior  $p(\mathbf{w}|D)$  with a simpler distribution  $q_{\xi}(\mathbf{w})$ . The latter is  
 40 chosen in a parametric family, with variational parameters  $\xi$ , as the minimizer of the Kullback-Leibler  
 41 (KL) divergence w.r.t.  $p$ :

$$\min_{\xi} KL(q_{\xi}(\mathbf{w}) || p(\mathbf{w} | D)) \quad (2)$$

42 It is well-known that minimizing the KL divergence is equivalent to maximizing the so-called *evidence*  
 43 *lower bound* (ELBO), which is defined as:

$$\text{ELBO}(\xi) = \mathbb{E}_{\mathbf{w} \sim q_{\xi}} [\log p(D|\mathbf{w})] - KL(q_{\xi}(\mathbf{w}) || p(\mathbf{w})) \quad (3)$$

44 Intuitively, the best approximation is the one that maximizes the expected log-likelihood of the data,  
 45 while minimizing the KL divergence w.r.t. the prior  $p(\mathbf{w})$ .

## 46 3 Variational Transfer Learning

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

### 52 3.1 Algorithm

53 We begin with a simple consideration: the distribution  $\mathcal{D}$  over tasks clearly induces a distribution over  
 54 optimal  $Q$ -functions. Since, for any MDP, learning its optimal  $Q$ -function is sufficient for solving the  
 55 problem, one can safely replace the distribution over tasks with the distribution over their optimal  
 56 value functions. Furthermore, assume we know such distribution and we are given a new task  $\tau$  to  
 57 solve. Then, our main intuition is that it is possible to design an algorithm that efficiently explores  $\tau$   
 58 so as to quickly adapt the prior distribution in a Bayesian fashion to put all probability mass over the  
 59 optimal  $Q$ -function of  $\tau$ .

60 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 simplicity,  
 61 we assume each function in  $\mathcal{Q}$  to be uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ <sup>1</sup>. Then, we can reduce our prior  
 62 distribution over  $Q$ -functions to a prior distribution over weights  $p(\mathbf{w})$ . Assume that we are given a  
 63 dataset  $D = \{(s_i, a_i, s'_i, r_i) \mid i = 1, 2, \dots, N\}$  of samples from some task  $\tau$  we want to solve. Then,  
 64 the posterior distribution over weights given such dataset can be computed by applying Bayes theorem  
 65 as in Eq. 1. Unfortunately, this cannot be directly used in practice since we do not have a model of  
 66 the likelihood  $p(D|\mathbf{w})$ . In such case, it is very common to make strong assumptions on the MDPs  
 67 or the  $Q$ -functions so as to get a tractable posterior []. On the other hand, we take a PAC-Bayesian  
 68 approach to derive a more general and meaningful posterior form. Recall that our final goal is move  
 69 all probability mass over the weights minimizing some empirical loss measure, which in our case  
 70 is the TD error  $\|B(\mathbf{w})\|_D^2$ . Then, given a prior  $p(\mathbf{w})$  we know from PAC-Bayesian theory that the  
 71 optimal Gibbs posterior takes the form []:

$$q(\mathbf{w}) = \frac{e^{-\Lambda \|B(\mathbf{w})\|_D^2} p(\mathbf{w})}{\int e^{-\Lambda \|B(\mathbf{w}')\|_D^2} p(d\mathbf{w}')} \quad (4)$$

72 for some parameter  $\Lambda > 0$ . Since  $\Lambda$  is typically chosen to increase with the number of samples  
 73  $N$ , we set it to  $\lambda^{-1}N$ , for some constant  $\lambda > 0$ . Notice that, whenever the term  $e^{-\Lambda \|B(\mathbf{w})\|_D^2}$  can

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

CITE

Cite some-  
body

Cite Catoni  
2007

74 be interpreted as the actual likelihood,  $q$  becomes a classic Bayesian posterior. Unfortunately, the  
 75 integral at the denominator of  $q$  is still intractable to compute even for simple  $Q$ -function models.  
 76 Thus, we propose a variational approximation  $q_\xi$  in a simpler family of distributions parameterized  
 77 by  $\xi \in \Xi$ . Then, our problem reduces to finding the variational parameters  $\xi$  such that  $q_\xi$  minimizes  
 78 the KL divergence w.r.t.  $q$ :

$$\min_{\xi \in \Xi} KL(q_\xi(\mathbf{w}) \parallel q(\mathbf{w})) = \min_{\xi \in \Xi} \mathbb{E}_{\mathbf{w} \sim q_\xi} \left[ \|B(\mathbf{w})\|_D^2 \right] - \frac{\lambda}{N} KL(q_\xi(\mathbf{w}) \parallel p(\mathbf{w})) \quad (5)$$

79 where the last objective is the well-known *evidence lower bound* (ELBO) []. Intuitively, the approxi-  
 80 mate posterior trades-off between placing probability mass over those weights  $\mathbf{w}$  that have low TD  
 81 error (first term), and staying close to the prior distribution (second term). Assuming that we are  
 82 able to compute the gradients of (5) w.r.t. the variational parameters, our objective can be easily  
 83 optimized with any stochastic optimization algorithm. Notice, however, that differentiating w.r.t.  $\xi$   
 84 typically requires differentiating  $\|B(\mathbf{w})\|_D^2$  w.r.t.  $\mathbf{w}$  (e.g., when using the reparameterization trick []).  
 85 Unfortunately, the TD error is well-known to be non-differentiable due to the presence of the max  
 86 operator. In iterative approaches (e.g., DQNs), this is typically solved by keeping the targets fixed  
 87 at each iteration, thus without differentiating them. However, our algorithm ... . To solve this issue,  
 88 we replace the optimal Bellman operator with the mellow Bellman operator introduced in [], which  
 89 adopts a softened version of max called *mellowmax*:

$$\text{mm}_a Q_{\mathbf{w}}(s, a) = \frac{1}{\kappa} \log \frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \quad (6)$$

90 where  $\kappa$  is a hyperparameter and  $|\mathcal{A}|$  is the number of actions. The mellow Bellman operator, which  
 91 we denote as  $\tilde{T}$ , has several appealing properties that make it suitable for our settings: (i) it converges  
 92 to the maximum as  $\kappa \rightarrow \infty$ , (ii) it has a unique fixed point, and (iii) it is *differentiable*. Denoting by  
 93  $\tilde{B}(\mathbf{w}) = \tilde{T}Q_{\mathbf{w}} - Q_{\mathbf{w}}$  the Bellman residual w.r.t. the mellow Bellman operator  $\tilde{T}$ , we have that the  
 94 corresponding TD error,  $\|\tilde{B}(\mathbf{w})\|_D^2$ , is now differentiable with respect to  $\mathbf{w}$ .

## 95 3.2 Gaussian Variational Transfer

## 96 3.3 Mixture of Gaussian Variational Transfer

## 97 4 Theoretical Analysis

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

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

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

$$\|V - \tilde{V}\|_{\nu, p}^p \leq \frac{2R_{max}}{(1 - \gamma)^2} \left\| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g}} \right\|_{\nu, p}^p \quad (7)$$

112 **5 Related Works**

113 **6 Experiments**

114 **6.1 Gridworld**

115 **6.2 Classic Control**

116 **6.3 Maze Navigation**

117 **7 Conclusion**

118 **References**

- 119 [1] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley  
120 & Sons, Inc., New York, NY, USA, 1994.

## 121 A Proofs

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

Cite MM

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{2R_{max}}{(1-\gamma)^2} \left\| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g}} \right\|_{\nu,p}^p \quad (7)$$

128 *Proof.* We begin by noticing that:

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

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

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{1}{1-\gamma} \|TV - \tilde{T}V\|_{\nu,p}^p \quad (8)$$

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

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

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

Cite MM

$$\begin{aligned} \left| \max_a Q(s, a) - \min_a Q(s, a) \right| &= \left| \sum_a \pi^*(a|s) Q(s, a) - \sum_a \pi_\beta(a|s) Q(s, a) \right| \\ &= \left| \sum_a Q(s, a) (\pi^*(a|s) - \pi_\beta(a|s)) \right| \\ &\leq \sum_a |Q(s, a)| |\pi^*(a|s) - \pi_\beta(a|s)| \\ &\leq \frac{R_{max}}{1-\gamma} \sum_a |\pi^*(a|s) - \pi_\beta(a|s)| \end{aligned} \quad (9)$$

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

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

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

$$\begin{aligned}
\sum_a |\pi^*(a|s) - \pi_\beta(a|s)| &= |\pi^*(a_1(s)|s) - \pi_\beta(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi^*(a|s) - \pi_\beta(a|s)| \\
&= |1 - \pi_\beta(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi_\beta(a|s)| \\
&= 2 |1 - \pi_\beta(a_1(s)|s)|
\end{aligned} \tag{10}$$

139 Finally, let us bound this last term:

$$\begin{aligned}
|1 - \pi_\beta(a_1(s)|s)| &= \left| 1 - \frac{e^{\beta Q(s, a_1(s))}}{\sum_{a'} e^{\beta Q(s, a')}} \right| \\
&= \left| 1 - \frac{e^{\beta(Q(s, a_1(s)) - Q(s, a_2(s)))}}{\sum_{a'} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&= \left| 1 - \frac{e^{\beta g(s)}}{\sum_{a'} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&= \left| 1 - \frac{e^{\beta g(s)}}{e^{\beta g(s)} + \sum_{a' \neq a_1(s)} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&\leq \left| 1 - \frac{e^{\beta g(s)}}{e^{\beta g(s)} + |\mathcal{A}|} \right| \\
&= \left| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g(s)}} \right|
\end{aligned} \tag{11}$$

140 Combining Eq. (9), (10), and (11), we obtain:

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

141 Taking the norm and plugging this into Eq. (8) concludes the proof.  $\square$

142 **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  $\delta > 0$ ,  
143 with probability at least  $1 - \delta$  over the choice of a  $Q$ -function  $Q$ , the following holds:  
144

$$\|Q - Q^*\|_\nu^2 \leq \frac{\mathbb{E}_p \left[ \|B(Q)\|_\nu^2 \right]}{(1 - \gamma)\delta} \tag{12}$$

145 *Proof.* First notice that:

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

146 which implies that:

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

147 Then we can write:

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

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

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

$$\|Q - \tilde{Q}\|_\nu^2 \leq \frac{\mathbb{E}_p \left[ \|\tilde{B}(Q)\|_\nu^2 \right]}{(1 - \gamma)\delta} \quad (13)$$

152 **Lemma 2.** Assume  $Q$ -functions belong to a parametric space of functions bounded by  $\frac{R_{max}}{1-\gamma}$ . Let  $p$   
 153 and  $q$  be arbitrary distributions over the parameter space  $\mathcal{W}$ , and  $\nu$  be a probability measure over  
 154  $\mathcal{S} \times \mathcal{A}$ . 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  
 155  $\delta > 0$ , with probability at least  $1 - \delta$ , the following two inequalities hold simultaneously:

$$\mathbb{E}_q \left[ \|B(\mathbf{w})\|_\nu^2 \right] \leq \mathbb{E}_q \left[ \|B(\mathbf{w})\|_D^2 \right] - \mathbb{E}_q [v(\mathbf{w})] + \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1 - \gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (14)$$

156

$$\mathbb{E}_q \left[ \|B(\mathbf{w})\|_D^2 \right] \leq \mathbb{E}_q \left[ \|B(\mathbf{w})\|_\nu^2 \right] + \mathbb{E}_q [v(\mathbf{w})] + \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1 - \gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (15)$$

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

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

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

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

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

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

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

$$\left| \|B(\mathbf{w})\|_\nu^2 + v(\mathbf{w}) - \|B(\mathbf{w})\|_D^2 \right| \leq 4 \frac{R_{max}^2}{(1 - \gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (16)$$

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

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

163 Thus, multiplying both sides of (16) by  $\lambda^{-1}N$  and applying the change of measure inequality with  
 164  $f(\mathbf{w}) = \lambda^{-1}N \left| \|B(\mathbf{w})\|_\nu^2 + v(\mathbf{w}) - \|B(\mathbf{w})\|_D^2 \right|$ , 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

165 where the second inequality holds since the right-hand side of (16) does not depend on  $\mathbf{w}$ . Finally,  
 166 we can explicitly write:

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

167 from which the lemma follows straightforwardly.  $\square$

168 **Lemma 3.** *Let  $p$  be a prior distribution over the parameter space  $\mathcal{W}$ , and  $\nu$  be a probability measure*  
 169 *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*  
 170 *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*  
 171 *least  $1 - \delta$ :*

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

172 *Proof.* Let us use Lemma 2 for the specific choice  $q = q_{\hat{\xi}}$ . From Eq. (14), we have:

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

173 where the second inequality holds since  $v(\mathbf{w}) > 0$ , while the equality holds from the definition of  $\hat{\xi}$ .

174 We can now use Eq. (15) to bound  $\mathbb{E}_{q_{\hat{\xi}}} \left[ \left\| B(\mathbf{w}) \right\|_D^2 \right]$ , thus obtaining:

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

175 This concludes the proof.  $\square$