

---

# Active Transfer Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1

## 2 1 Introduction

## 3 2 Preliminaries

4 **Markov Decision Processes** We define a Markov decision process (MDP) as a tuple  $\mathcal{M} =$   
5  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, p_0, \gamma \rangle$ , where  $\mathcal{S}$  is the state-space,  $\mathcal{A}$  is a finite set of actions,  $\mathcal{P}(\cdot | s, a)$  is the distribution  
6 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,  $p_0$   
7 is the initial-state distribution, and  $\gamma \in [0, 1)$  is the discount factor. We assume the reward function  
8 to be uniformly bounded by a constant  $R_{max} > 0$ . A deterministic policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  is a mapping  
9 from states to actions. At the beginning of each episode of interaction, the initial state  $s_0$  is drawn  
10 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  
11 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  
12 the long-term return over a possibly infinite horizon:  $\max_{\pi} J(\pi) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | \mathcal{M}, \pi]$ . To this  
13 end, we define the optimal value function  $Q^*(s, a)$  as the expected return obtained by taking action  
14  $a$  in state  $s$  and following an optimal policy thereafter. Then, an optimal policy  $\pi^*$  is a policy that  
15 is greedy with respect to the optimal value function, i.e.,  $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$  for all states  
16  $s$ . It can be shown (e.g., [4]) that  $Q^*$  is the unique fixed-point of the optimal Bellman operator  $T$   
17 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  
18 adopt the term  $Q$ -function to denote any plausible value function, i.e., any function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$   
19 uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ .

20 We define the Bellman residual of a  $Q$ -function  $Q$  as  $B(Q) \triangleq TQ - Q$ . Notice that a  $Q$ -function  $Q$   
21 is optimal if, and only if,  $B(Q)(s, a) = 0$  for all  $s, a$ .

22 **Multitask Settings** We represent tasks  $\tau$  as MDPs with shared state and action spaces, but with  
23 potentially different values for all other parameters. We assume the existence of a distribution  $\mathcal{D}$  over  
24 tasks, i.e.,  $\tau \sim \mathcal{D}$ , and we suppose that we are able to sample from such distribution.

## 25 3 Approach

26 We start by noticing that the task distribution  $\mathcal{D}$  clearly induces a distribution over optimal  $Q$ -  
27 functions. Then, our goal is to estimate such distribution from the set of source tasks and use it as a  
28 prior for speeding up the learning process in the target task.

29 Consider solving the target task given a dataset  $D$  of  $N$  samples. The posterior distribution over  
30 optimal  $Q$ -functions is:

$$P(Q | D) \propto P(D | Q)P(Q), \quad (1)$$

where  $P(D | Q)$  is the likelihood of observing the dataset  $D$  given that  $Q$  is optimal and  $P(Q)$  is the prior distribution induced by  $\mathcal{D}$ . Computing the likelihood requires knowing the transition and reward models of the current task (i.e., the MDP model), which are not available in practice. However, since we are conditioned on the fact that  $Q$  is optimal, we can derive a simple approximation. We consider the empirical Bellman error  $\|B(Q)\|_{p,D}^p$  under the  $l_p$ -norm defined as:

$$\|B(Q)\|_{p,D}^p = \frac{1}{N} \sum_{i=1}^N \left| r_i + \gamma \max_{a'} Q(s'_i, a') - Q(s_i, a_i) \right|^p \quad (2)$$

Assuming states and actions to be drawn from a fixed joint distribution  $\mu$ , we have  $\|B(Q)\|_{p,\mu}^p = 0$  whenever  $Q$  is optimal. Then, the probability of observing a dataset  $D$  clearly decreases exponentially as  $\|B(Q)\|_{p,D}^p$ , for an optimal  $Q$ , increases. This can be seen, for instance, by applying Hoeffding's inequality. Thus, a natural way to model the likelihood  $P(D | Q)$  is:

$$P(D | Q) \propto e^{-\lambda N \|B(Q)\|_{p,D}^p} \quad (3)$$

where  $\lambda$  is a constant hyperparameter. Intuitively, this indicates that  $D$  is more likely when it induces low empirical Bellman error under an optimal  $Q$ -function. Furthermore, as the number of available samples  $N$  increases, the distribution becomes more peaked at zero since the empirical Bellman error converges to the true Bellman error. In the limit  $N \rightarrow \infty$ :

$$P(D | Q) \propto I_{\{\|B(Q)\|_{p,D}^p=0\}} \quad (4)$$

### 3.1 Regularized Bellman Residual Minimization with Gaussian Priors

Taking the maximum of the log-posterior, we obtain the following optimization problem:

$$\min_{Q \in \mathcal{Q}} \|B(Q)\|_{p,D}^p - \log P(Q) \quad (5)$$

Let us specify a particular hypothesis space  $\mathcal{Q}$ . We consider  $Q$ -functions  $Q_w$  parameterized by the vector  $w$ . Our optimization problem becomes:

$$\min_w \|B(w)\|_{p,D}^p - \log P(w) \quad (6)$$

We model the prior distribution over the optimal parameters  $w$  as a Gaussian. That is, we assume:

$$P(w) = \mathcal{N}(w, \Sigma) \quad (7)$$

Then, our optimization, adopting the  $l_2$ -norm, becomes:

$$\min_w \|B(w)\|_D^2 + \|w - \mu\|_{\Sigma} = \min_w \frac{1}{N} \sum_{i=1}^N |y_i - Q_w(s_i, a_i)|^2 + \|w - \mu\|_{\Sigma} \quad (8)$$

where  $y_i = r_i + \gamma \max_{a'} Q_w(s'_i, a')$  and  $\|x\|_A \triangleq x^T A^{-1} x$  for  $A$  positive definite matrix.

**Linear model** We assume a linear model for the  $Q$ -functions:  $Q_w(s, a) = w^T \phi(s, a)$ . Here  $\phi$  is a  $K$ -dimensional feature vector. Then, the solution to the optimization problem of Eq. (8) can be computed in closed form as follows:

$$w^* = (A^T A + \Sigma^{-1})^{-1} (A^T b + \Sigma^{-1} \mu) \quad (9)$$

where  $A$  is an  $N \times K$  matrix containing the feature vectors at each data point  $(s_i, a_i)$  and  $b$  is an  $N$ -dimensional vector containing their targets  $y_i$ . The resulting algorithm is a regularized version of LSVI.

**Neural network** We model  $Q_w(s, a)$  as a neural network with parameters  $w$ . Then, the gradient of the objective function of Eq. (8) can be easily computed by standard backpropagation. The resulting algorithm is a regularized version of NFQI.

### 3.2 Variational Inference for Efficient Exploration

A major drawback of the maximum-a-posteriori (MAP) approach is that it does not explicitly estimate the posterior distribution  $P(Q | D)$ . Thus, given a dataset  $D$ , the algorithm consistently chooses the same MAP  $Q$ -function, and the only way to allow learning is to introduce a simple form of exploration (e.g.,  $\epsilon$ -greedy or softmax), whose drawbacks are well-known. The advantages of a more time-coherent exploration strategy such as posterior sampling (or Thompson sampling) have already been proved by many existing works (CITE). In particular, direct posterior sampling of  $Q$ -functions (also known as value function randomization), when a distribution is available, have been proven to be very effective and scalable. Thus, we now extend our previous approach to explicitly estimate  $P(Q | D)$ , so as to allow efficient exploration of the target task. The main complication is that, as mentioned earlier, we cannot compute the likelihood  $P(D | Q)$  explicitly, but we only have access to a model that is proportional to it. Then, estimating the posterior distribution is clearly intractable since it requires computing the marginal likelihood  $P(D) = \int P(D | Q) dQ$ . Thus, we resort variational inference (CITE) to approximate the intractable posterior with a simpler distribution from which we are able to sample.

The main idea behind variational inference is to approximate the posterior  $P(Q|D)$  with a simpler distribution  $q_\phi(Q)$ , parameterized by the vector  $\phi$ , from which we can easily get samples. The best approximation is chosen in terms of the Kullback-Leibler (KL) divergence between the two distributions, that is:

$$\min_{\phi} KL(q_\phi(Q) || P(Q | D)) \quad (10)$$

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

$$ELBO(\phi) = \mathbb{E}_{Q \sim q_\phi} [\log P(D|Q)] - KL(q_\phi(Q) || P(Q)) \quad (11)$$

As before, we assume a parametric form  $Q_w$  for representing  $Q$ -functions. Thus, we reduce distributions over functions to distributions over vectors  $w \in \mathbb{R}^K$ . Given a specific form for the approximate posterior  $q_\phi(w)$  and the prior  $P(w)$ , our inference problem reduces to:

$$\min_{\phi} \mathbb{E}_{w \sim q_\phi} [\lambda N \|B(w)\|_D^2] + KL(q_\phi(w) || P(w)) \quad (12)$$

Intuitively, the approximate posterior trades-off between placing probability mass over those weights  $w$  that have low empirical Bellman error (first term), and staying close to the prior distribution that is learned from the set of source tasks (second term). The factor  $N$  multiplying  $\|B(w)\|_D^2$  makes sure that, as the number of samples goes to infinity, only the first term is considered in the optimization. In such case, all probability mass goes over the (approximate) fixed-point of the optimal Bellman operator. On the other hand, for small values of  $N$ , there is less evidence that the empirical Bellman error closely estimate the true one, thus giving more importance to staying close to the prior.

Depending on the choice of the prior and approximate posterior distributions, the optimization problem of Eq. (12) might still be very hard to solve. We now propose some suitable choices that make the problem efficiently solvable.

**Gaussian distributions** We model both the prior and the approximate posterior as Gaussian distributions:  $P(w) = \mathcal{N}(\bar{\mu}, \bar{\Sigma})$  and  $q_\phi(w) = \mathcal{N}(\mu, \Sigma)$ , respectively. For the sake of simplicity, we assume  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_K^2)$  and  $\bar{\Sigma} = \text{diag}(\bar{\sigma}_1^2, \dots, \bar{\sigma}_K^2)$ . Notice, in this case, that the parameter vector  $\phi$  is  $[\mu_1, \dots, \mu_K, \sigma_1^2, \dots, \sigma_K^2]$ . The KL divergence between the two distributions can be computed in closed form as:

$$KL(q_\phi(w) || P(w)) = \frac{1}{2} \left( \log \frac{|\bar{\Sigma}|}{|\Sigma|} + \text{Tr}(\bar{\Sigma}^{-1}\Sigma) + (\mu - \bar{\mu})^T \bar{\Sigma}^{-1}(\mu - \bar{\mu}) - K \right) \quad (13)$$

To solve the optimization problem of Eq. (12) we adopt stochastic backpropagation (CITE). Thus, we need to compute the gradient of the objective function  $\mathcal{L}(\phi)$  of Eq. (12) with respect to each component of  $\phi$ . This can be easily done for the second term:

$$\nabla_{\mu} KL(q_\phi(w) || P(w)) = \bar{\Sigma}^{-1}(\mu - \bar{\mu}) \quad (14)$$

$$\nabla_{\sigma_k^2} KL(q_\phi(w) || P(w)) = \frac{1}{2} \left( \frac{1}{\bar{\sigma}_k^2} - \frac{1}{\sigma_k^2} \right) \quad (15)$$

For the first term, we have to compute the gradient of an expectation with respect to the parameters of the Gaussian distribution under which the expectation is taken. For the mean, we can resort to Thm. X of (CITE Bonnet 1964):

$$\nabla_{\mu_k} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\lambda N \|B(\mathbf{w})\|_D^2] = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\lambda N \nabla_{w_k} \|B(\mathbf{w})\|_D^2] \quad (16)$$

For the covariance, we can resort to Thm. Y of (CITE Price 1958):

$$\nabla_{\sigma_k^2} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\lambda N \|B(\mathbf{w})\|_D^2] = \frac{1}{2} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\lambda N \nabla_{w_k, w_k}^2 \|B(\mathbf{w})\|_D^2] \quad (17)$$

We refer the reader to (CITE Deepmind) for a simple proofs of these two theorems.

It only remains to compute the gradient and Hessian of the empirical Bellman error  $\|B(\mathbf{w})\|_D^2$  with respect to  $\mathbf{w}$ . Recall that:

$$\|B(\mathbf{w})\|_D^2 = \frac{1}{N} \sum_{i=1}^N \left( r_i + \gamma \max_{a'} Q_{\mathbf{w}}(s'_i, a') - Q_{\mathbf{w}}(s_i, a_i) \right)^2 \quad (18)$$

A major complication is that computing the gradient with respect to  $\mathbf{w}$  requires differentiating through the maximum operator. In iterative approaches such as Fitted Q-Iteration (FQI) or Deep Q-Networks (DQNs), this is typically solved by computing the Bellman target using the previous weights and adapting the current weights so that the resulting Q-functions matches the target. In such case, no differentiation through the maximum operator is required. However, in our case there is no previous weight: we simply want to know what is the Bellman error of any given weight  $\mathbf{w}$  and we are thus required to differentiate through the maximum operator when updating the approximate posterior, which is not possible in practice. To solve this issue, we replace the max with the *mellowmax* operator defined by:

$$\text{mm}_a Q_{\mathbf{w}}(s, a) = \frac{1}{\kappa} \log \frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \quad (19)$$

where  $\kappa$  is a hyperparameter and  $|\mathcal{A}|$  is the number of actions. The mellowmax operator has several appealing properties that make it suitable for our settings: (i) it converges to the maximum as  $\kappa \rightarrow \infty$ , (ii) it has a unique fixed point, and (iii) it is *differentiable*. Replacing max with the mellowmax operator, we obtain a new approximation to the Bellman residual:

$$\tilde{B}(\mathbf{w})(s, a) \triangleq \mathcal{R}(s, a) + \gamma \mathbb{E}_{\mathcal{P}} [\text{mm}_{a'} Q_{\mathbf{w}}(s', a')] - Q_{\mathbf{w}}(s, a) \quad (20)$$

and to the empirical Bellman error:

$$\|\tilde{B}(\mathbf{w})\|_D^2 = \frac{1}{N} \sum_{i=1}^N \left( r_i + \gamma \text{mm}_{a'} Q_{\mathbf{w}}(s'_i, a') - Q_{\mathbf{w}}(s_i, a_i) \right)^2 \quad (21)$$

Of fundamental importance, this last quantity is now differentiable with respect to  $\mathbf{w}$ .

Differentiating through the mellowmax operator with respect to  $\mathbf{w}$  is now very simple:

$$\nabla_{w_k} \text{mm}_a Q_{\mathbf{w}}(s, a) = \nabla_{w_k} \frac{1}{\kappa} \log \frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \quad (22)$$

$$= \frac{1}{\kappa} \frac{1}{\frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)}} \frac{1}{|\mathcal{A}|} \sum_a \nabla_{w_k} e^{\kappa Q_{\mathbf{w}}(s, a)} \quad (23)$$

$$= \frac{1}{\kappa} \frac{1}{\frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)}} \frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \kappa \nabla_{w_k} Q_{\mathbf{w}}(s, a) \quad (24)$$

$$= \frac{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \nabla_{w_k} Q_{\mathbf{w}}(s, a)}{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)}} \quad (25)$$

126 The diagonal of the Hessian can also be computed straightforwardly:

$$\nabla_{\mathbf{w}_k, \mathbf{w}_k}^2 \min_a Q_{\mathbf{w}}(s, a) = \nabla_{\mathbf{w}_k} \frac{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \nabla_{\mathbf{w}_k} Q_{\mathbf{w}}(s, a)}{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)}} \quad (26)$$

$$= \frac{\nabla_{\mathbf{w}_k} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \nabla_{\mathbf{w}_k} Q_{\mathbf{w}}(s, a)}{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)}} + \nabla_{\mathbf{w}_k} \left( \frac{1}{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)}} \right) \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \nabla_{\mathbf{w}_k} Q_{\mathbf{w}}(s, a) \quad (27)$$

$$= (\kappa + 1) \frac{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \nabla_{\mathbf{w}_k, \mathbf{w}_k}^2 Q_{\mathbf{w}}(s, a)}{\sum_a e^{\kappa Q_{\mathbf{w}}(s, a)}} - \kappa \frac{\left( \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \nabla_{\mathbf{w}_k} Q_{\mathbf{w}}(s, a) \right)^2}{\left( \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \right)^2} \quad (28)$$

## 127 4 Related Works

128 List of any paper that relates to our approach together with a brief description:

- 129 • [3]: the authors propose a method for efficient exploration via randomized value functions.  
130 The optimal  $Q$ -function is computed by bayesian LSVI and, after each update, parameters  
131 are sampled from the posterior. Then, the agent follows a greedy policy with respect to the  
132 sampled  $Q$ -function. Regret bounds are provided.
- 133 • [2]: the authors extend the idea of randomized value functions to drive exploration in deep  
134 RL. A posterior distribution over  $Q$ -functions is approximated via bootstrapping. In each  
135 episode, the agent acts greedily with respect to a  $Q$ -function sampled from the approximated  
136 posterior.
- 137 • [1]: the authors build on top of bootstrapped DQNs to provide UCB-like exploration bonuses.  
138 In a previous (?) version of the paper, exploration bonuses based on information gain are  
139 also proposed.

## 140 5 Experiments

## 141 6 Conclusion

## 142 References

- 143 [1] Richard Y Chen, John Schulman, Pieter Abbeel, and Szymon Sidor. Ucb and infogain exploration via  
144 q-ensembles. *arXiv preprint arXiv:1706.01502*, 2017.
- 145 [2] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped  
146 dqn. In *Advances in neural information processing systems*, pages 4026–4034, 2016.
- 147 [3] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value  
148 functions. *arXiv preprint arXiv:1402.0635*, 2014.
- 149 [4] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley  
150 & Sons, Inc., New York, NY, USA, 1994.

