

Appendix D Discount factor inference

D.1 Model specification

The main idea we will try to develop here is that the discount factor γ can be considered as a latent random variable, so that inference is made about its value in order to get the most precise estimate of the posterior predictive distribution of the long term return of an action $p(V(s, a)|x_{\leq j})$.

For practical reasons, it is easier to model separately the factorized distribution of the current, immediate state-action reward $p(r(s, a))$ and the discounted value of the future state $p(\gamma v(s', a))$. In this notation, $\gamma v(s', a)$ is a random variable representing the discounted long-term value of taking the action a in state s minus the value of the reward received while performing this action $r(s, a)$:

$$V(s, a) = r(s, a) + \gamma v(s', a)$$

Following this logic, we define the following factorized distribution:

$$p(r(s, a), v(s', a)) = p(r(s, a)|\mu^r(s, a), \sigma^r(s, a)) p(v(s', a)|\mu^v(s, a), \sigma^v(s', a), \gamma)$$

where $r(s, a)$ is normally distributed and $v(s', a)$ is distributed according to the following modified Gaussian distribution:

$$p(v(s', a)|\mu^v(s, a), \sigma^v(s, a), \gamma) = \gamma \frac{\exp\left\{-\frac{(\mu^v(s, a) - \gamma v(s', a))^2}{2\sigma^v(s, a)^2}\right\}}{\sqrt{2\pi} \sigma^v(s, a)}.$$

It becomes clear that this distribution integrates to 1 over $v(s', a)$, and that it corresponds to a Normal distribution over $\gamma v(s', a)$. Therefore, the probability distribution of $V(s, a) = r(s, a) + \gamma v(s', a)$ can be written as:

$$p(V(s, a)) = \mathcal{N}(\mu^r(s, a) + \mu^v(s, a), \sigma^r(s, a)^2 + \sigma^v(s, a)^2).$$

This joint probability distribution encodes, for each state-action pair the probability distribution of the associated long-term, discounted return.

We naturally define the prior probability of the discounting factor $\gamma \in [0; 1]$ at $j = 1$ as a Beta distribution with parameters $B(\alpha_0^\gamma, \beta_0^\gamma)$, and a \mathcal{NG}^{-1} prior over the Gaussian parameters $\{\mu^r(s, a), \sigma^r(s, a), \mu^v(s, a), \sigma^v(s, a)\}$.

D.2 Learning from observed transitions and rewards

We now tackle the question of how to estimate the value of the state that has been reached $v(s', a)$. Indeed, this value, contrary to the reward, is not directly observed but must be inferred from the state of arrival s' . We consider the following state-value inference model:

$$\begin{aligned} v(s', a) &= \mathbb{E}_{\pi(a'|s')} [V(s', a')] \\ &= \mathbb{E}_{\pi(a'|s')} [r(s', a') + v(s', a)]. \end{aligned} \tag{9}$$

Eq 9 stated that the agent assumes that the expected value of the next state is equal to the average of the various state-action values available, weighted by the probability of selection the corresponding action (the observed policy $\pi(a'|s')$). We assume that the agent does not have a direct access to the complete description of the policy (meaning that, whereas she can make a decision efficiently, she cannot retrieve the summary

statistics of this policy). This assumption is in accordance with the decision process described in 2.6. This policy can, however, be learned similarly to $r(s, a)$ using a multinomial distribution with a Dirichlet (or Beta in the case of Two-Alternative Forced Choice tasks) prior and approximate posterior: at each trial, the agent observes the action she has performed, and update her belief of the current policy.

Let us now consider how the expected log-joint probability of $v(s'|s, a)$ would appear in the ELBO if this value was observed under the mean-field assumption:

$$\begin{aligned} \mathbb{E}_{q_{j+1}(\mu^v(s,a), \sigma^v(s,a)^2, \gamma)} \left[\log p(v(s'|s, a) | \mu^v(s, a), \sigma^{v^2}, \gamma) \right] &= -\frac{1}{2} \mathbb{E}_{q_{j+1}(\sigma^v(s,a)^2)} \left[\log \sigma^v(s, a)^2 \right] \\ &- \frac{1}{2} \mathbb{E}_{q_{j+1}(\mu^v(s,a), \sigma^v(s,a)^2)} \left[\frac{(\gamma v(s'|s, a) - \mu^v(s, a))^2}{\sigma^v(s, a)^2} \right] + \mathbb{E}_{q_{j+1}(\gamma)} [\log \gamma] + \text{cst.} \end{aligned} \quad (10)$$

Once the agent reaches the next state \hat{s}' , inference of the value $\mathbb{E} [v(s'|s, a) | s' = \hat{s}']$ can be achieved by using the posterior predictive expectation of this formula:

$$\begin{aligned} \mathbb{E} [v(s'|s, a) | s' = \hat{s}'] &= \mathbb{E}_{p(v|x_{\leq j})} [v(\hat{s}'|s, a)] \\ &\approx \mathbb{E}_{q_j(\mathbf{z})} \left[\mathbb{E}_{p(V, \pi | \mathbf{z})} \left[\mathbb{E}_{\pi} [V(\hat{s}', a')] \right] \right] \\ &= \mathbb{E}_{q_j(\mathbf{z})} \left[\mathbb{E}_{p(r, v, \pi | \mathbf{z})} \left[\mathbb{E}_{\pi} [r(\hat{s}', a') + \gamma v(s''|s', a')] \right] \right] \\ &= \mathbb{E}_{q_j(\mathbf{z})} \left[\mathbb{E}_{p(r, v, \pi | \mathbf{z})} \left[\sum_{a' \in A} \pi(a' | \hat{s}') (r(\hat{s}', a') + \gamma v(s''|\hat{s}', a')) \right] \right] \\ &= \sum_{a' \in a_1, a_2} \frac{\delta_j^\pi(\hat{s}', a')}{\sum_{a \in A} \delta_j^\pi(\hat{s}', a)} (\mu_j^r(\hat{s}', a') + \mu_j^v(\hat{s}', a')) \end{aligned} \quad (11)$$

where A is the set of actions available, and where we indexed with j each variational parameter to emphasize the fact that this expectation is made with respect to the previous posterior belief.

Eq 10 being linear and quadratic wrt $v(s'|s, a)$, we also need to solve the expected value $\mathbb{E} [v(s'|s, a)^2 | s' = \hat{s}']$:

$$\begin{aligned} \mathbb{E}_{p(v|x_{\leq j})} [v(s'|s, a)^2] &\approx \\ \mathbb{E}_{q_j(\mathbf{z})} \left[\mathbb{E}_{p(r, v, \pi | \mathbf{z})} \left[\left(\sum_{a' \in A} \pi(a' | s') (r(s', a') + \gamma v(s''|s', a')) \right)^2 \right] \right]. \end{aligned} \quad (12)$$

The expression in Eq 12 can be computed iteratively (Algorithm 1):

Algorithm 1: $\mathbb{E}_{p(v(s'|s,a)|x_{\leq j})} [v(s'|s,a)^2 | s' = \widehat{s}']$ iterative computation

```

input: Factorized approximate posterior distribution  $q_{j-1}$ , state reached at trial  $j$ 
 $\widehat{s}'$ 
1 Results: Posterior predictive estimate  $\nu = \mathbb{E}_{p(v(s'|s,a)|x_{\leq j})} [v(s'|s,a)^2 | s' = \widehat{s}']$ 
2  $\nu = 0$ 
3 for  $a \in A$  do
4   for  $a' \in A$  do
5     if  $a == a'$  then
6        $\nu += \mathbb{E}_{q_{j-1}} [\pi(a|s')^2] \times$  // Eq 17
7          $(\mathbb{E}_{q_{j-1}} [\mathbb{E}_p [r(s', a)^2]] + \mathbb{E}_{q_{j-1}} [\mathbb{E}_p [(\gamma v(s''|s', a))^2]])$ 
8         // Eq 14, Eq 16
9          $+ 2\mathbb{E}_{q_{j-1}} [\mathbb{E}_p [r(s', a)]] \mathbb{E}_{q_{j-1}} [\mathbb{E}_p [\gamma v(s''|s', a)]]$  // Eq 13, Eq 15
10    else
11       $\nu += \mathbb{E}_{q_{j-1}} [\pi(a|s')\pi(a'|s')] \times$  // Eq 18
12         $(\mathbb{E}_p [r(s', a)] + \mathbb{E}_{q_{j-1}} [\mathbb{E}_p [\gamma v(s''|s', a)]]) \times$  // Eq 13, Eq 15
13         $(\mathbb{E}_p [r(s', a')] + \mathbb{E}_{q_{j-1}} [\mathbb{E}_p [\gamma v(s''|s', a')]])$  // Eq 13, Eq 15
14    end
15  end
16 end

```

using the following equivalences

$$\mathbb{E}_{q_{j-1}} [\mathbb{E}_{p(r|\mu^r, \sigma^r)} [r]] = \mu_{j-1}^r \quad (13)$$

$$\mathbb{E}_{q_{j-1}} [\mathbb{E}_{p(r|\mu^r, \sigma^r)} [r^2]] = \mu_{j-1}^r{}^2 + \frac{\beta_{j-1}^r}{(\alpha_{j-1}^r - 1)\kappa_{j-1}^r} + \frac{\beta_{j-1}^r}{(\alpha_{j-1}^r - 1)} \quad (14)$$

$$\mathbb{E}_{q_{j-1}} [\mathbb{E}_{p(v|\mu^v, \sigma^v, \gamma)} [\gamma v]] = \mu_{j-1}^v \quad (15)$$

$$\mathbb{E}_{q_{j-1}} [\mathbb{E}_{p(v|\mu^v, \sigma^v, \gamma)} [(\gamma v)^2]] = \mu_{j-1}^v{}^2 + \sigma_{j-1}^v{}^2 \quad (16)$$

$$\mathbb{E}_{q_{j-1}} [\pi(a|s)^2] = \left(\frac{\delta_{j-1}^{\pi(a|s)}}{\delta_0} \right)^2 + \frac{\delta_{j-1}^{\pi(a|s)}(\delta_0 - \delta_{j-1}^{\pi(a|s)})}{\delta_0^2(\delta_0 + 1)} \quad (17)$$

$$\mathbb{E}_{q_{j-1}} [\pi(a|s)\pi(a'|s)] = \frac{\delta_{j-1}^{\pi(a|s)}\delta_{j-1}^{\pi(a'|s)}}{\delta_0^2} - \frac{\delta_{j-1}^{\pi(a|s)}\delta_{j-1}^{\pi(a'|s)}}{\delta_0^2(\delta_0 + 1)} \quad (18)$$

$$\text{with } \delta_0 = \sum_{a' \in A} \delta^{\pi(a'|s)}$$

where we have dropped several state-action indices for clarity.

We can now take the expectation of Eq 10 under Eq 11 and Eq 12.

This formula has the advantage that the learning of the current distribution of $V(s, a)$ takes into account the uncertainty about the policy, about the reward distribution at the next trial and about the long-term discounted return of the actions in the next state.