

EVb Decomposition

Change in the value function due to learning after taking action a^* :

$$\begin{aligned}
 v(ba^*) - v(b) &= \sum_{b'} p(b' | b, a^*) (v(b') - v(b)) \\
 &= \sum_{b'} p(b' | b, a^*) \left(\sum_a \pi(a | b') q(b', a) - \sum_a \pi(a | b) q(b, a) \right) \\
 &= \sum_{b'} p(b' | b, a^*) \sum_a \left((\pi(a | b') - \pi(a | b)) q(b', a) \right. \\
 &\quad \left. + \pi(a | b) (q(b', a) - q(b, a)) \right)
 \end{aligned} \tag{1}$$

Expanding $q(b', a) - q(b, a)$:

$$\begin{aligned}
 q(b', a) - q(b, a) &= \sum_{b''} p(b'' | b', a) [r(b', a) + \gamma v(b'')] \\
 &\quad - \sum_{b'} p(g' | b, a) [r(b, a) + \gamma v(g')] \\
 &= r(b', a) + \gamma \sum_{b''} p(b'' | b', a) v(b'') \\
 &\quad - r(b, a) + \gamma \sum_{g'} p(g' | b, a) v(g') \\
 &= \underbrace{r(b', a) - r(b, a)}_{\text{Difference in the expected immediate return}} + \underbrace{\gamma \left[\sum_{b''} p(b'' | b', a) v(b'') - \sum_{g'} p(g' | b, a) v(g') \right]}_{\text{Difference in the expected future return}}
 \end{aligned} \tag{2}$$

So overall the EVb decomposes as:

$$\begin{aligned}
 v(ba^*) - v(b) &= \mathbb{E}_{b' \sim p(b' | b, a^*)} \left[\sum_a (\pi(a | b') - \pi(a | b)) q(b', a) \right. \\
 &\quad \left. + \mathbb{E}_{a \sim \pi(a | b)} [r(b', a) - r(b, a)] \right. \\
 &\quad \left. + \mathbb{E}_{a \sim \pi(a | b)} [\gamma \sum_{b''} p(b'' | b', a) v(b'') - \gamma \sum_{g'} p(g' | b, a) v(g')] \right]
 \end{aligned} \tag{3}$$

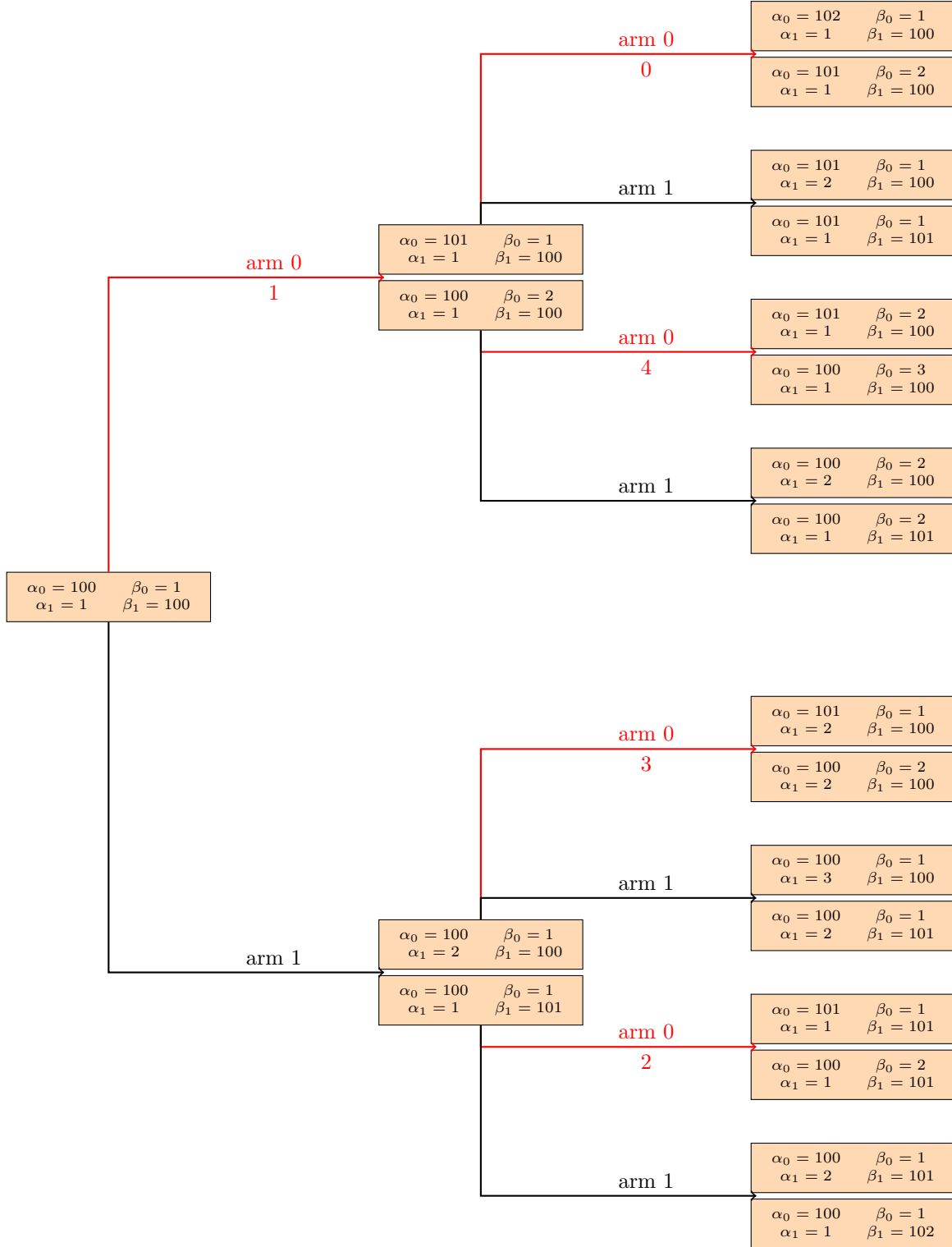
One last thing to consider is how the prioritisation of distal experiences should differ from those that are more immediate. This also has implications for how likely those experiences are to occur according to the current model.

For instance, if the agent considers updating $v(b)$ towards the value that would result from taking a particular action from that belief state – say, $v(ba^*)$ – the EVb associated with that update needs to be weighted by the probability of transitioning into belief state b in the first place (i.e., from the current root of the tree).

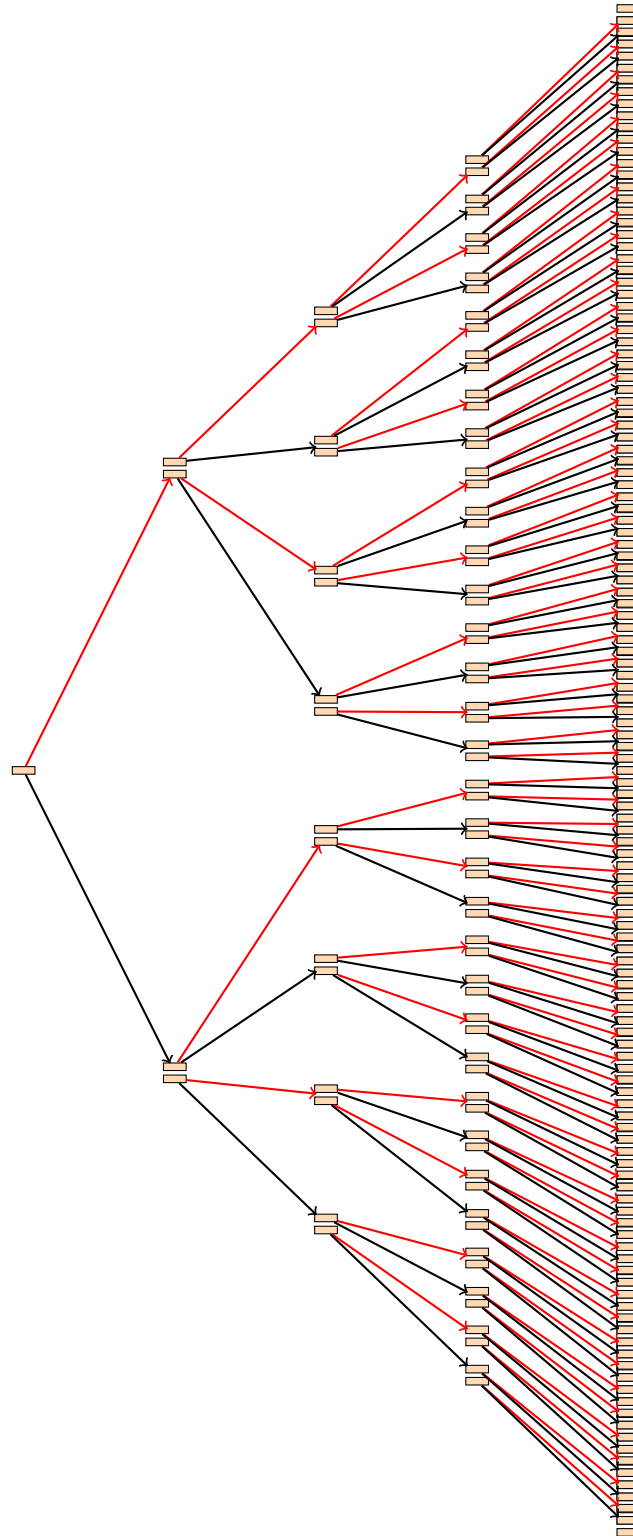
$$\begin{aligned}
 v(ba^*) - v(b) &= p(b_{\text{root}} \rightarrow b) \times \left(\mathbb{E}_{b' \sim p(b' | b, a^*)} \left[\sum_a (\pi(a | b') - \pi(a | b)) q(b', a) \right. \right. \\
 &\quad \left. \left. + \mathbb{E}_{a \sim \pi(a | b)} [r(b', a) - r(b, a)] \right. \right. \\
 &\quad \left. \left. + \mathbb{E}_{a \sim \pi(a | b)} [\gamma \sum_{b''} p(b'' | b', a) v(b'') - \gamma \sum_{g'} p(g' | b, a) v(g')] \right] \right)
 \end{aligned} \tag{4}$$

Simulations

Prioritisation pattern with horizon $h = 2$. Numbers under the arrows show the order in which the replay updates (red arrows) were executed.



Prioritisation pattern with horizon $h = 4$. The prior belief at the root is set to the same values as in the above example with shorter horizon. Note that for this tree, EVB was not scaled by the probability of reaching any belief.



Prioritisation pattern with horizon $h = 4$. Same example as before, however, here, each EVB value $v(ba^*) - v(b)$ is scaled by the probability of reaching belief b from the root of the tree.

