

Empirical convergence analysis

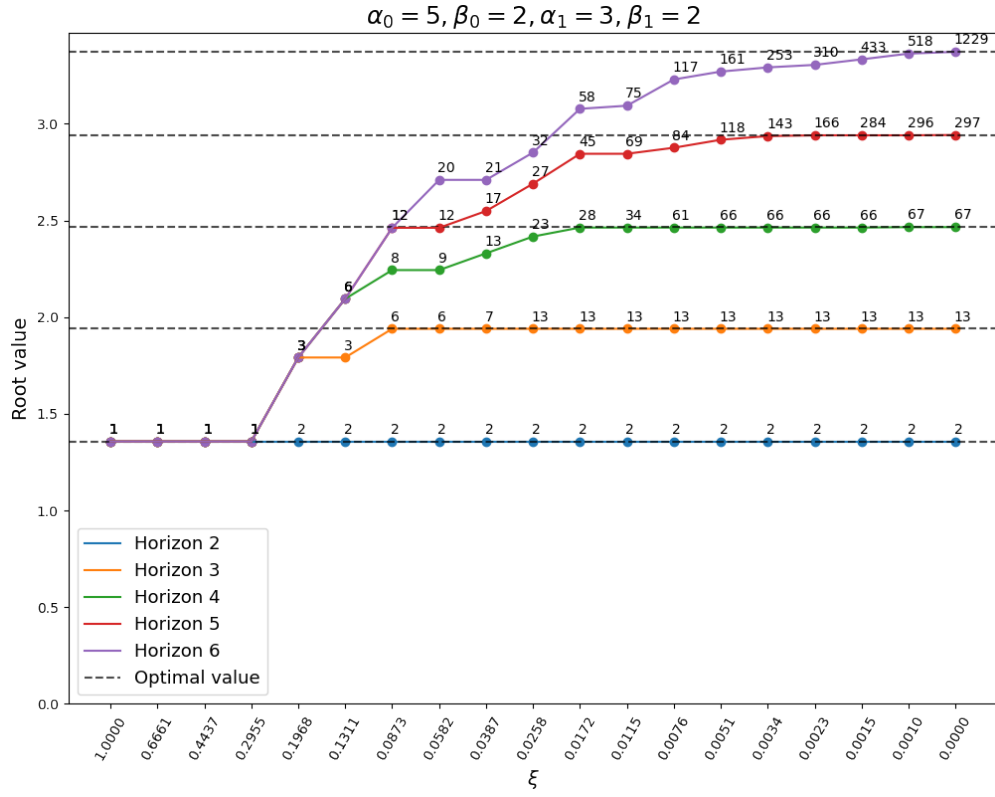
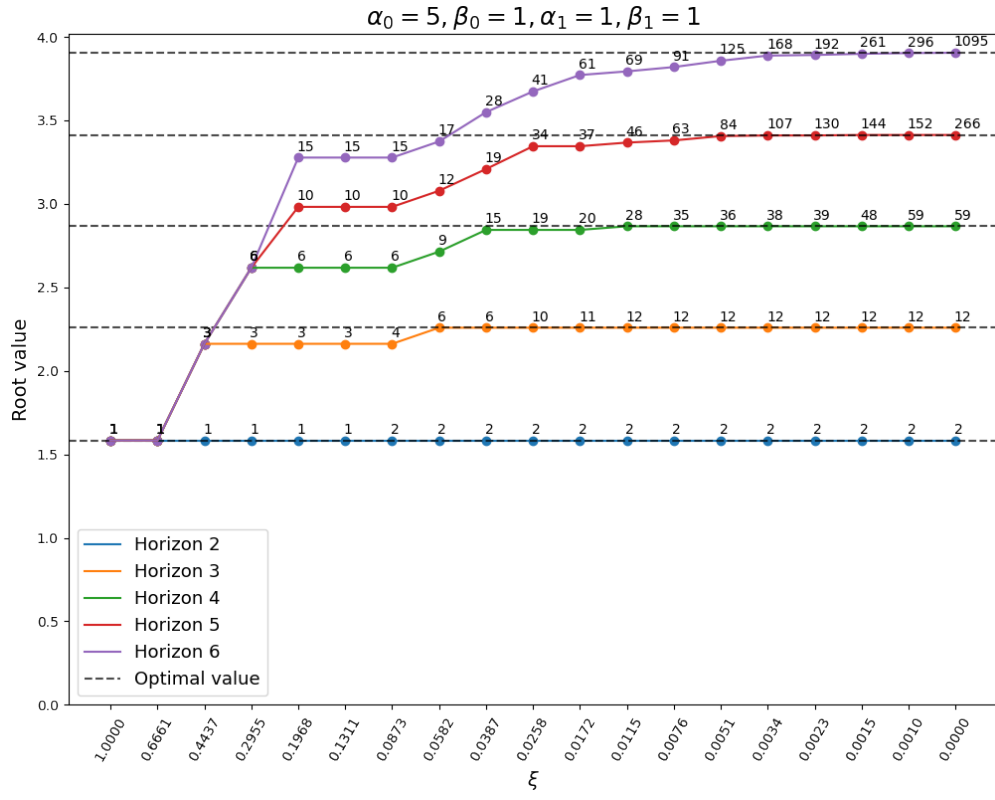
I am still puzzled by the choice of policy *inside* the tree. The options are, for instance, to set all Q -values for the intermediate nodes within the tree to 0 – in which case the policy would be uniform; or, alternatively, those Q -values could be set to the immediate expected reward according to those intermediate belief states. The former choice will force the replays to be initiated at the leaves of the tree (i.e., at the final horizon reached), since only the leaf values would have non-zero values to propagate. As for the other option, the replay patterns would be very different because of the Need term which discounts states that are further away.

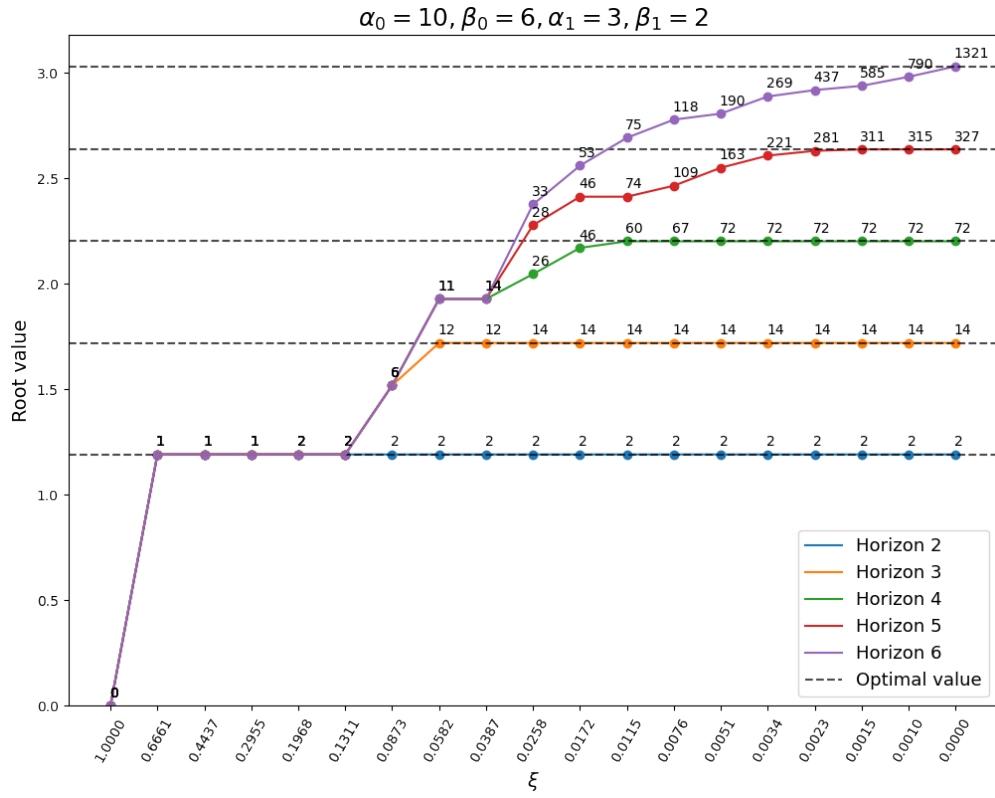
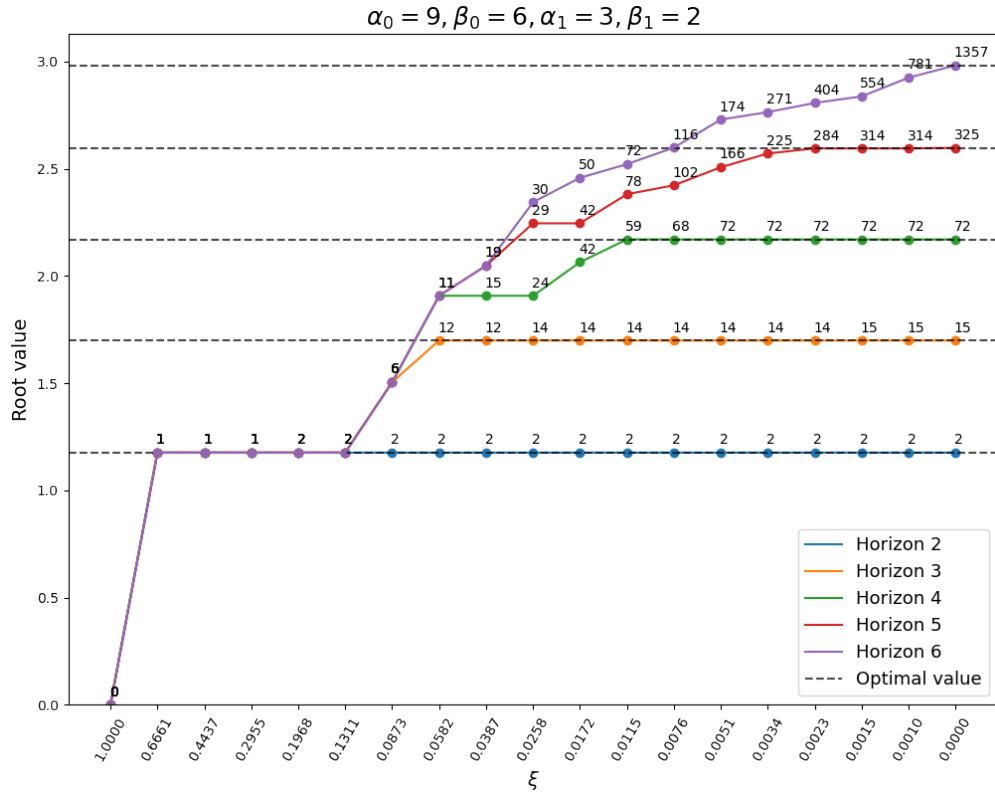
Moreover, I think that the uniform policy would make it less likely that distal beliefs for the sub-optimal action (at the root) get replayed. Indeed, the choice of the tree policy seems to affect the rate of convergence. I generated convergence plots for the two policies for comparison.

Replay threshold

The plots below show how the value of the root state changes with the EVB threshold ξ at multiple horizons and for various prior beliefs. Note that the root Q -values (i.e, the current MF Q -values) were set to 0 everywhere unless specified otherwise.

Numbers on top of each data point indicate the number of replays. Titles specify the prior beliefs at the root.





Number of replays