

Beyond Average Return in Markov Decision Processes

Alexandre Marthe¹, Aurélien Garivier¹, Claire Vernade²
¹ ENS de Lyon, ² University of Tuebingen



Setting

We consider the setting of **Undiscounted Finite-Horizon** Tabular Markov Decision Processes $\mathcal{M}(\mathcal{X}, \mathcal{A}, P, R, H)$.

We are interested in the **Return** seen as a random variable :

$$Z^\pi(x) = \sum_{h=0}^H R_h, \quad X_0 = x$$

Most of the literature focuses on optimizing the **Expected Return** :

$$\max_{\pi} \mathbb{E}[Z^\pi(x)]$$

→ can we optimize other *risk-aware* functionals ψ of the Return, such as CVaR, quantiles, etc?

$$\max_{\pi} \psi(Z^\pi(x))$$

Motivation

For some environments with a complex Return distribution, optimizing the mean might be arbitrary. The functional ψ can be a risk measure of the Return, leading to risk-dependent strategies.

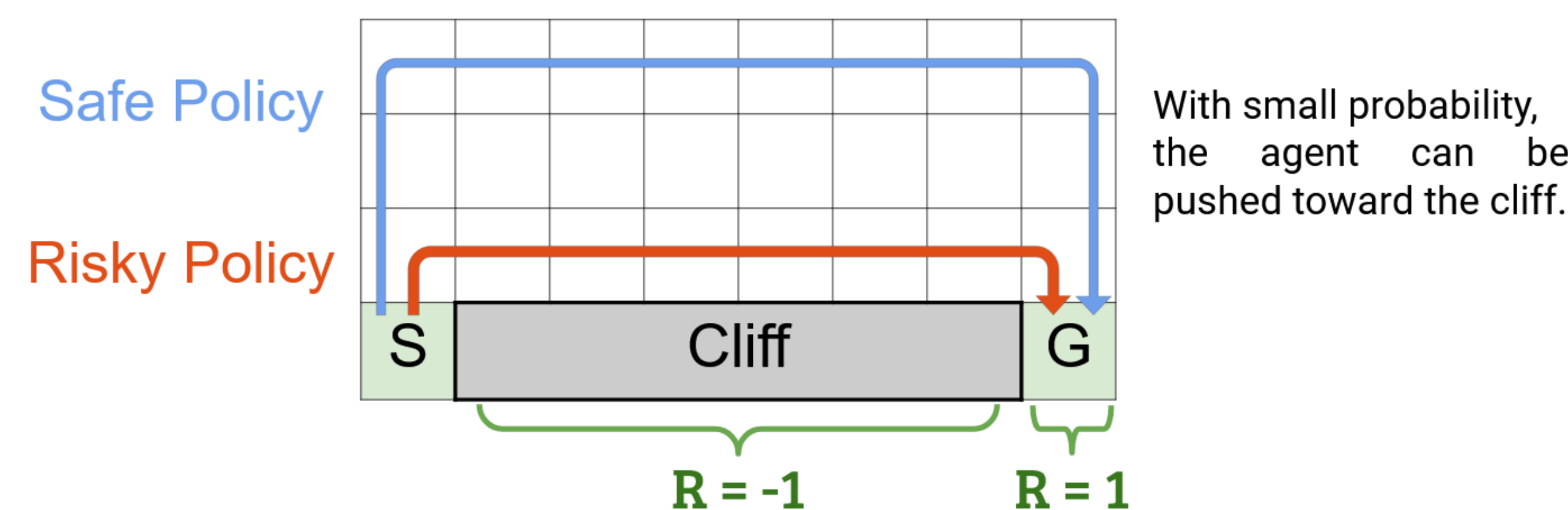


Figure 1: Cliff MDP where the agent should go from S to G.

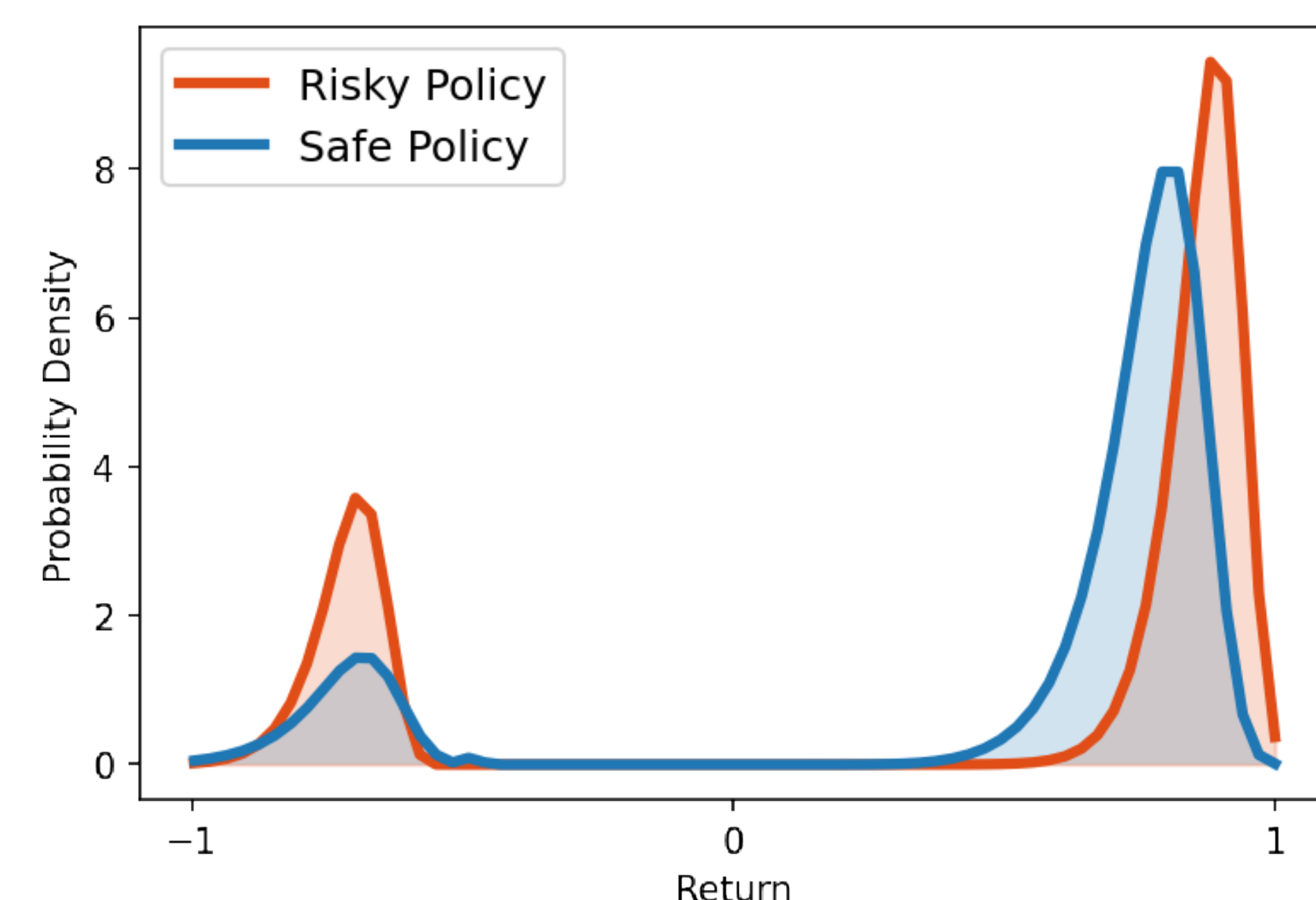


Figure 2: Distribution of the return for the policy in Cliff.

Can such policy be computed through dynamic programming ?

Question

What functionals of the return can be *optimized exactly* through Dynamic Programming ?

Main Result

The only continuous **Bellman Optimizable** functionals are:

- The Expected Return : $\mathbb{E}[Z]$
- The Exponential Utilities : $\mathbb{E}[\exp(\lambda Z)]$

→ The statistics optimizable through Distributional RL are the same as with Classical RL.

→ Exponential utilities allow for risk-dependant strategies
 e.g. if $Z \sim \mathcal{N}(\mu, \sigma)$, $\lambda^{-1} \ln \mathbb{E}[\exp(\lambda Z)] = \mu + \lambda \sigma$

Distributional RL

→ Objective of Distributional RL: Estimate the whole distribution of the return, instead of only its expectation.

→ There exists a Distributional Bellman Equation:

$$Z_h^\pi(x, a) = R_h + Z_{h+1}^\pi(X', A')$$

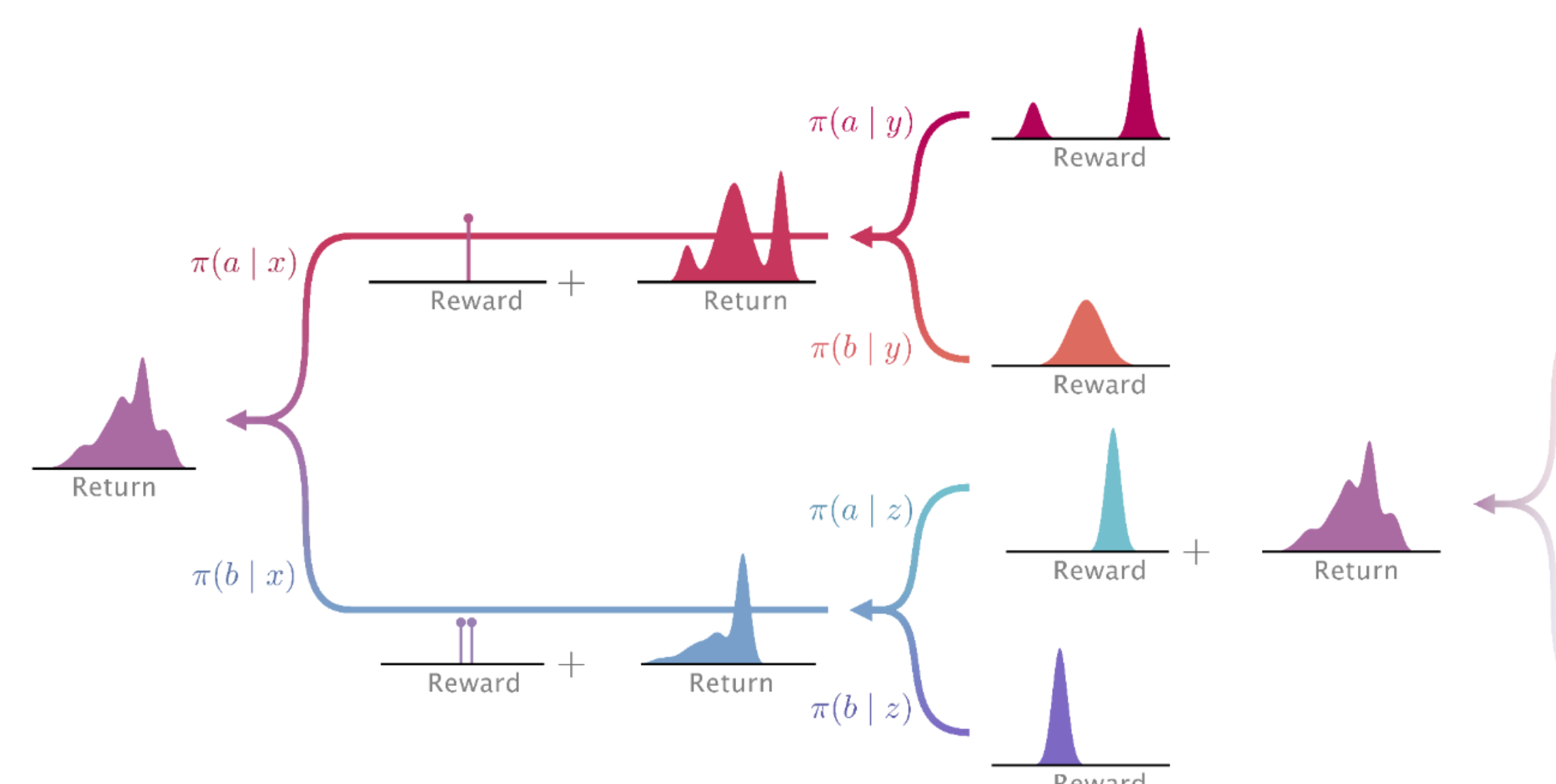


Figure 3: Distributional Dynamic Programming. [Bellemare et. al, 2023]

→ Convenient tool for Dynamic Programming: $\forall h, x$, choose $a^* = \operatorname{argmax}_a \psi(Z_h^\pi(x, a))$. From $h = H$ down to 0, compute recursively a policy π_{DP}

Definition : Bellman Optimizable

A functional ψ is said to be *Bellman Optimizable* if π_{DP} is optimal for any Markov Decision Processes.

Dynamic Programming

Bellman Equation and Dynamic Programming :

$$Q_h^*(x, a) = \mathbb{E}[R_h] + \sum_{x'} p(x, a, x') \max_{a'} Q_{h+1}^*(x', a')$$

Choice independent on the previous reward (1) Choice independent on other potential states (2)

Figure 4: Dynamic Programming for Expected Return

The distributional form, using $Z_h^{\pi_{\text{DP}}}(x, a) \sim \eta_h^{\pi_{\text{DP}}}(x, a)$ and $R_h(x, a) \sim \varrho_h(x, a)$ is:

$$\psi(\eta_h^{\pi_{\text{DP}}}(x, a)) = \psi\left(\varrho_h(x, a) * \sum_{x'} p_h(x, a, x') \eta_{h+1}^{\pi_{\text{DP}}}(x', a_{h+1}^*(x'))\right)$$

Bellman Optimizable Properties

A Bellman Optimizable functional ψ necessarily verifies 2 properties:

- **1. Independence Property:** Mixing in other distributions should not change the choice of action.

$$\psi(\nu_1) \geq \psi(\nu_2) \implies \forall \nu_3, \forall \lambda, \psi(\lambda \nu_1 + (1-\lambda) \nu_3) \geq \psi(\lambda \nu_2 + (1-\lambda) \nu_3)$$

→ allows to apply the Expected Utility Theorem: ψ can be written in the form $E[f(\cdot)]$ for some f .

- **2. Translation Property:** Translating by a constant should not change the choice of action.

$$\psi(\nu_1) \geq \psi(\nu_2) \implies \forall c, \psi(\nu_1(\cdot + c)) \geq \psi(\nu_2(\cdot + c))$$

→ allows to find a differential equation verified by f . The solutions are the Exponential and Linear functions.

Important References

- Bellemare, M. G., Dabney, W., & Rowland, M. (2023). *Distributional reinforcement learning*. MIT Press.
- R. A. Howard and J. E. Matheson. *Risk-Sensitive Markov Decision Processes*. 1972.

