# CHAPTER 1

## DEEP REINFORCEMENT LEARNING

## 1 Neural Networks

...

## 2 Distributional Reinforcement Learning

Remember that we defined $\mathbb{P}_{s,a}^{\pi} := \mathbb{P}^{\pi} \otimes \delta_{S_0}(s) \otimes \delta_{A_0}(a)$ as the probability measure of the Markov reward process $(S, A, R)$ started in $(s, a)$. We define the random variable of the return under policy $\pi$ as

$$Z^{\pi} := \sum_{t=0}^{\infty} \gamma^t R_t, \quad \gamma \in (0, 1).$$

Unlike the methods before, where we were interested in the expected reward $Q^{\pi}(s, a) = \mathbb{E}_{s,a}^{\pi}[Z^{\pi}]$, we are now interested in the distribution of these cumulative rewards. For that define

$$\eta_{s,a}^{\pi}(B) := \mathbb{P}_{s,a}^{\pi}(Z^{\pi} \in B), \quad B \in \mathcal{B}(\mathbb{R}).$$

In analogy to weak solutions for PDEs, i.e. in sense of distribution, the return distribution $\eta_{s,a}^{\pi}$ satisfies the Bellman equation in sense of distribution:

$$\forall \phi \in C_b(\mathbb{R}): \quad \int_{\mathbb{R}} \phi(z) d\eta_{s,a}^{\pi}(z) = \mathbb{E}_{s,a}^{\pi}[\int_{\mathbb{R}} \phi(R + \gamma z) d\eta_{S',A'}^{\pi}(z)].$$

We define $f_{r,\gamma}(z) := r + \gamma z$ and the push forward

$$((\eta_{s,a}^{\pi})_{f_{r,\gamma}})(B) := \eta_{s,a}^{\pi}(f_{r,\gamma}^{-1}(B)), \quad B \in \mathcal{B}(\mathbb{R}).$$

then, the above can be written as

$$\forall \phi \in C_b(\mathbb{R}): \quad \int_{\mathbb{R}} \phi(z) d\eta_{s,a}^{\pi}(z) = \mathbb{E}_{s,a}^{\pi}[\int_{\mathbb{R}} \phi(R + \gamma z) d\eta_{S',A'}^{\pi}(z)] = \mathbb{E}_{s,a}^{\pi}[\int_{\mathbb{R}} \phi(z) d(\eta_{S',A'}^{\pi})_{f_{R,\gamma}}(z)]$$
$$\iff \forall \phi \in C_b(\mathbb{R}): \quad \langle \phi, \eta_{s,a}^{\pi} \rangle = \mathbb{E}_{s,a}^{\pi}[\langle \phi, (\eta_{S',A'}^{\pi})_{f_{R,\gamma}} \rangle].$$

In the weak sense we define the Bellman operator as

$$\mathcal{T}^{\pi} \eta_{s,a}^{\pi} := \sum_{s',a',r} \pi(a'; s') p(s', r; s, a) (\eta_{s',a'}^{\pi})_{f_{r,\gamma}}.$$