

## ① Distributional RL

Typical Objective fn: Maximize expected discounted return.

In Value-function based,

$$\text{Bellman Expectation Eq}^n \quad Q(s, a) = E[\overbrace{R(s, a) + \gamma \sum_{s'} Q(s', a')}^{\text{return}}]$$

$$\text{Bellman Optimality Eq}^n \quad Q(s, a) = E[R(s, a) + \gamma \max_{a'} Q(s, a')] \\ = E[R(s, a)] + \gamma E[\max_{a'} Q(s, a')]$$

In Distributed RL,

$$Q(s, a) = E[Z(s, a)] \\ = E[R(s, a) + \gamma \underbrace{Z(s', a')}_{\substack{\text{return, random variable} \\ \text{from some dist.}}}]$$

↑  
We want to learn this!

## ③ Projection step

Estimated  $Z$ :  $Z_i(s, a)$  Target:  $R_i + \gamma Z_i(s', a')$

Supports' disjoint, cannot compute KL loss. (goes to  $\infty$ )

⇒ Projection: Transfer prob. mass from misaligned target atoms to closest neighboring estimate 'aligned' atoms.

Then compute KL loss.  $KL(p||q)$  b/w target  $p$  and estimate  $q$ .

Minimize loss using gradient descent.

## ② Algorithm (C51)

### ① Objective functions for learning

In Value based, tabular: TD error

In Approx setting: MSE

In value-dist, propose w. dist. BUT

In practice, KL divergence used instead of W distance.

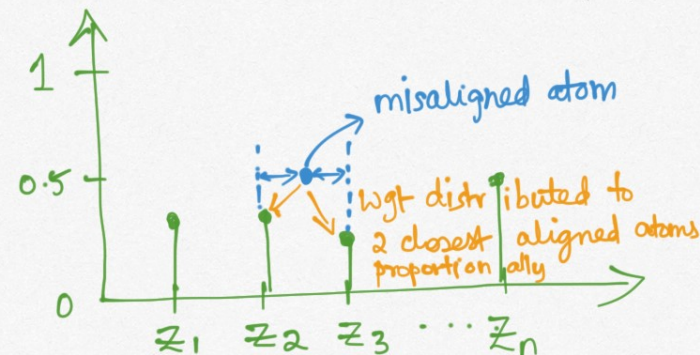
### ② Defining $Z(s, a)$ ( $N=51$ in C51)

$Z$  defined as discrete dist of  $N$  fixed atoms (support)

$$Z_\theta(s, a) = Z_i \text{ w.p. } p_i = \frac{e^{\theta_i(s, a)}}{\sum_j e^{\theta_j(s, a)}} \quad (\text{softmax})$$

$Z_i$ 's bounded in  $[V_{\min}, V_{\max}]$

$$Z_i = V_{\min} + i \Delta Z \quad \Delta Z = \frac{V_{\max} - V_{\min}}{N-1}$$





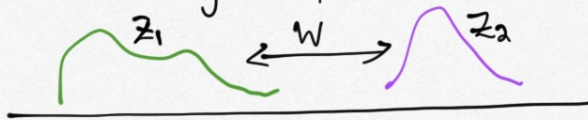
### ③ Wasserstein Distance

Dist b/w 2 prob distributions.

$$W(z_1, z_2) = \inf_{\pi \in \Pi(z_1, z_2)} \mathbb{E}_{(x, y) \sim \pi} \|x - y\|$$

$\Pi$  - set of transport plans to move "earth" from one dist to other

$\pi$  - transport plan w/ min. cost/dist  
 $(x, y)$  sampled from  $\pi$   
 $\|x - y\|$  - dist b/w  $x$  &  $y$ . Can also be  $L_p$  norm to give  $p$ -Wasserstein dist



Primal is intractable, joint dist. We can use dual.  
 Dual also has convergence guarantees in policy evaluation case.

$$W(z_1, z_2) = \sup_{f \in \mathcal{F}_{Lip}} \mathbb{E}_{x \sim z_1} [f(x)] - \mathbb{E}_{y \sim z_2} [f(y)]$$

Or written as: (HOW?)

$$\bar{W}(z_1, z_2) = \sup_{s, a} W(z_1(s, a), z_2(s, a))$$

This dual form can be proved to be  $\gamma$ -contraction in  $W$  in policy eval case.  
 In control, cannot guarantee.

### ⑤ Why should this work?

In value function based,

a) Bellman operator  $T$  is  $\gamma$ -contraction mapping

$$\|TF - TG\|_{\infty} \leq \gamma \|F - G\|_{\infty}$$

(has convergence properties  $F^* = TF^*$ ,  $F_t = TF_{t-1}$ )  
 in tabular setting. Converges to fixed point.

In value distribut<sup>n</sup>, we can use Wasserstein dist; dual form is  $\gamma$ -contraction mapping (PROOF). Only in policy evaluation case.

Banach's fixed pt theorem:  $d(T(F), T(G)) \leq \gamma d(F, G)$   
 ( $T$  is contraction mapping in  $d$ )  $d$  - some distance measure  
 So, also convergence to fixed point.

[But authors claim  $W$  dist cannot be used w/ samples. Hence KL divergence used instead.

But, dual of  $W$  dist is in terms of Expect<sup>n</sup> which means we can use samples to estimate, & same measure is used in WGANs.]