

Distributional RL

① Typical Objective fn: Maximize expected discounted return.

In Value-function based,

$$\text{Bellman Expectation Eq, } Q(s, a) = E \left[R(s, a) + \gamma \max_{s'} Q(s', a') \right]$$

$$\begin{aligned} \text{Bellman Optimality Eq, } Q(s, a) &= E \left[R(s, a) + \gamma \max_{a'} Q(s, a') \right] \\ &= E[R(s, a)] + \gamma E \left[\max_{a'} Q(s, a') \right] \end{aligned}$$

In Distributed RL,

$$\begin{aligned} Q(s, a) &= E[Z(s, a)] \\ &\quad \xrightarrow{\text{return, random variable from some dist.}} \\ &= E[R(s, a) + \gamma Z(s', a')] \end{aligned}$$

We want to learn this!

② Algorithm (C5I)

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

Supports' disjoint, cannot compute KL loss. (goes to ∞)

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

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

Minimize loss using gradient descent.

② Algorithm (C5I)

a) Objective functions for learning

In Value-based, tabular: TD error

In Approx setting: MSE

In value-dist, propose W.distance BUT

In practice, KL divergence used instead of W distance.

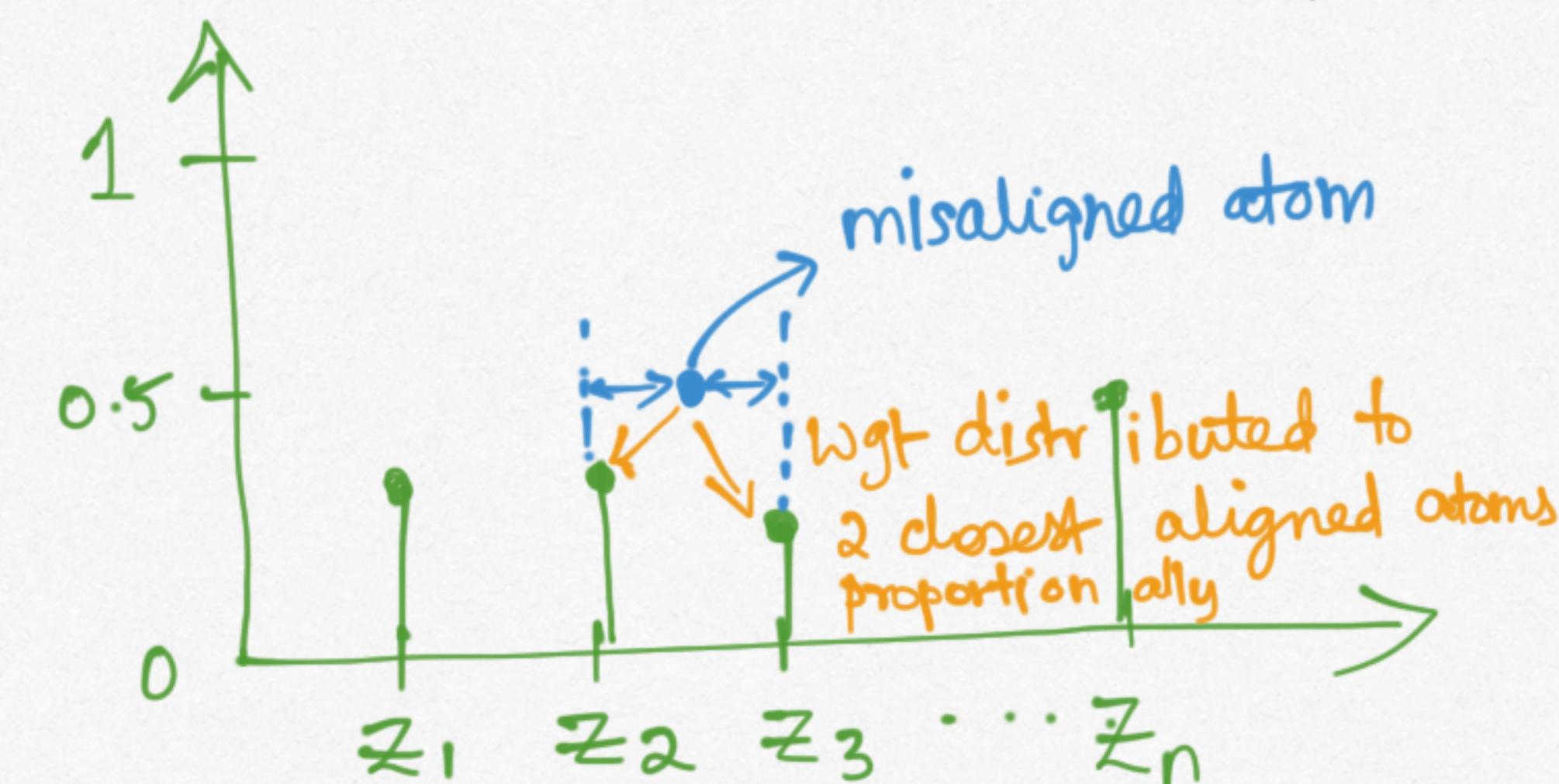
b) Defining $Z(s, a)$ ($N=51$ in C5I)

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.

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

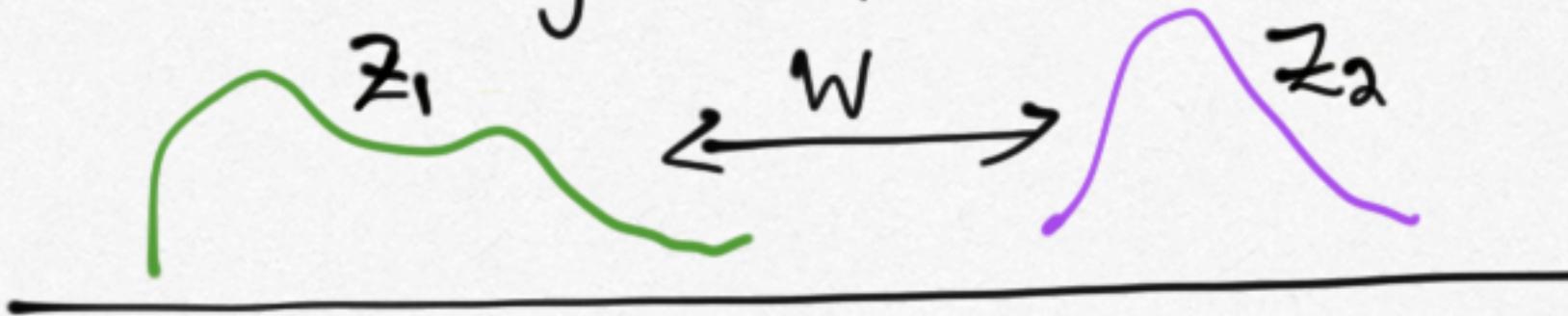
Π - set of transport plans to move "earth" from one dist to other

π - transport plan w/ min. cost/distance

(x, y) sampled from π

$\|x - y\|$ - dist b/w x & y . Can also be

L_p norm to give p -Wasserstein distance



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

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

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 γ -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 γ -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 distribution, we can use Wasserstein distance, dual form is γ -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 distance cannot be used w/ samples. Hence KL divergence used instead.]

But, dual of W distance is in terms of Expectaf, which means we can use samples to estimate, & same measure is used in WGANs.]

Implicit Quantile Networks (IQN)

- Prev work (QR-DQN) learnt a distributⁿ over returns but still used mean returns for policy.
- IQN learns a risk sensitive policy.

Risk-Sensitive RL - I

- Instead of $\max E[\text{Returns}]$, $\max E[\text{Utility}]$.
- E.g.s of Utility: Exp/Log(Returns), Measure that includes variance etc.
- In general, utility f^n satisfies 4 axioms
 - ① Monotonicity: lower cost \Rightarrow lower risk
 - ② Sub-additivity: Diversification \Rightarrow lower risk
 - ...
- Can be concave(risk averse)/convex(seeking) or both.

Risk-Sensitive RL - II

- One more axiom of utility f^n is of independence:
if $X \succ Y$ (preferred over) Z , then $pX + (1-p)Z \succ pY + (1-p)Z$
(Any mixture of $X \& Z$ preferred over same mixture over $Y \& Z$).
We can think of utility as mixture of distributions, which we want to learn, to learn a risk-sensitive policy.

(As before, we want to use quantiles to learn distributⁿ)

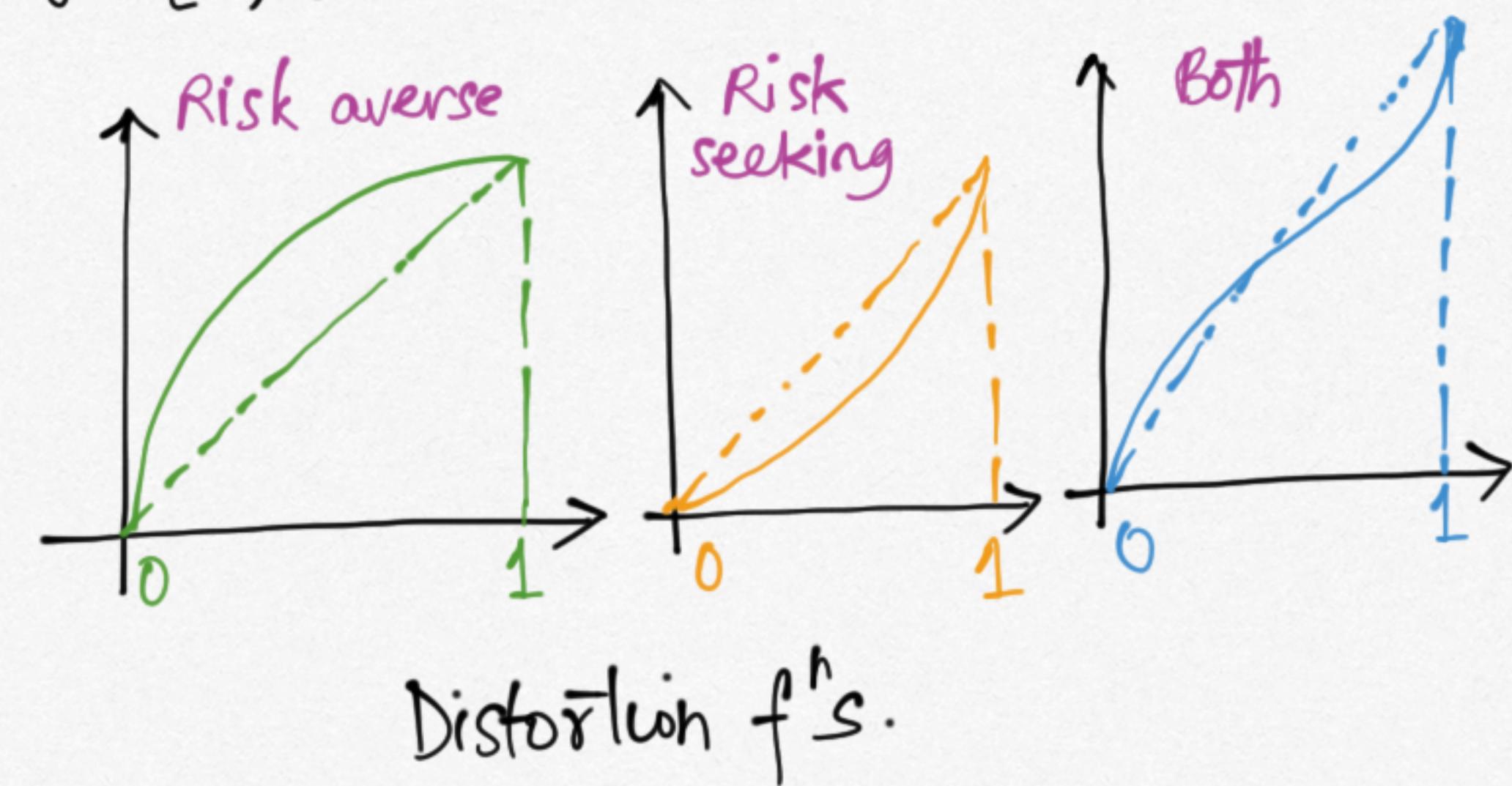
Learning mixture of distributions (MoD)

- MoD can be represented in terms of cdf.
 - cdf of MoD is convex combination of its individual distributions.
- $$F(x) = \sum_i p_i F_i(x)$$
- MoD is also represented as convex comb. of quantiles (in insurance/reinsurance...) where weighting is some risk fⁿ. (e.g. VaR)

Distribution Risk Measures (DRM)

- Var, Conditional Var are DRMs.
- In general, DRM is a continuous non-decreasing fⁿ.

$$g: [0,1] \rightarrow [0,1] \text{ s.t } g(0)=0, g(1)=1.$$



Dual theory of Choice under Risk

- An alternate objective for risk sensitive policy:
Instead of maximizing expected utility,
maximize expectation of distortion risk measure
 \Rightarrow i.e. Distorted Expectation.

Learning a Distorted Expectation

Say $Z(s, a)$ is our usual returns distribution.
for some quantile τ , $E[Z_\tau(s, a)]$ is expected return
 $\tau \in U(0, 1)$ for that quantile.

We can reweight this expectation w/ β (DRM),
 $Q_\beta(s, a) := E[z_{\beta(\tau)}(s, a)]$

[Can think of this as Inverse Transform Theorem:
To generate samples from some distribution, sample
from $Unif(0, 1)$, & transform them through inverse
CDF (quantile) of the prob. distribution.]

$Q_\beta(s, a)$ can be used to follow policy (action-selection)
 $\pi_\beta(s) = \operatorname{argmax}_{a \in A} Q_\beta(s, a)$

IQN Loss function

- Similar to QR-DQN we want to minimize loss b/w pairs of quantiles. But we don't have a fixed set of quantiles as in QR-DQN
- Instead, parametric $f^n \phi(t)$ used to represent quantiles (embedding layer in NN).
- So we can generate samples, N for τ , N' for τ' & calc. loss.
$$L = \frac{1}{N'} \sum_{i=1}^N \sum_{j=1}^{N'} p_{\tau_i}^k (\delta^{\tau_i, \tau'_j})$$

 $\delta^{\tau, \tau'} - TD$ error

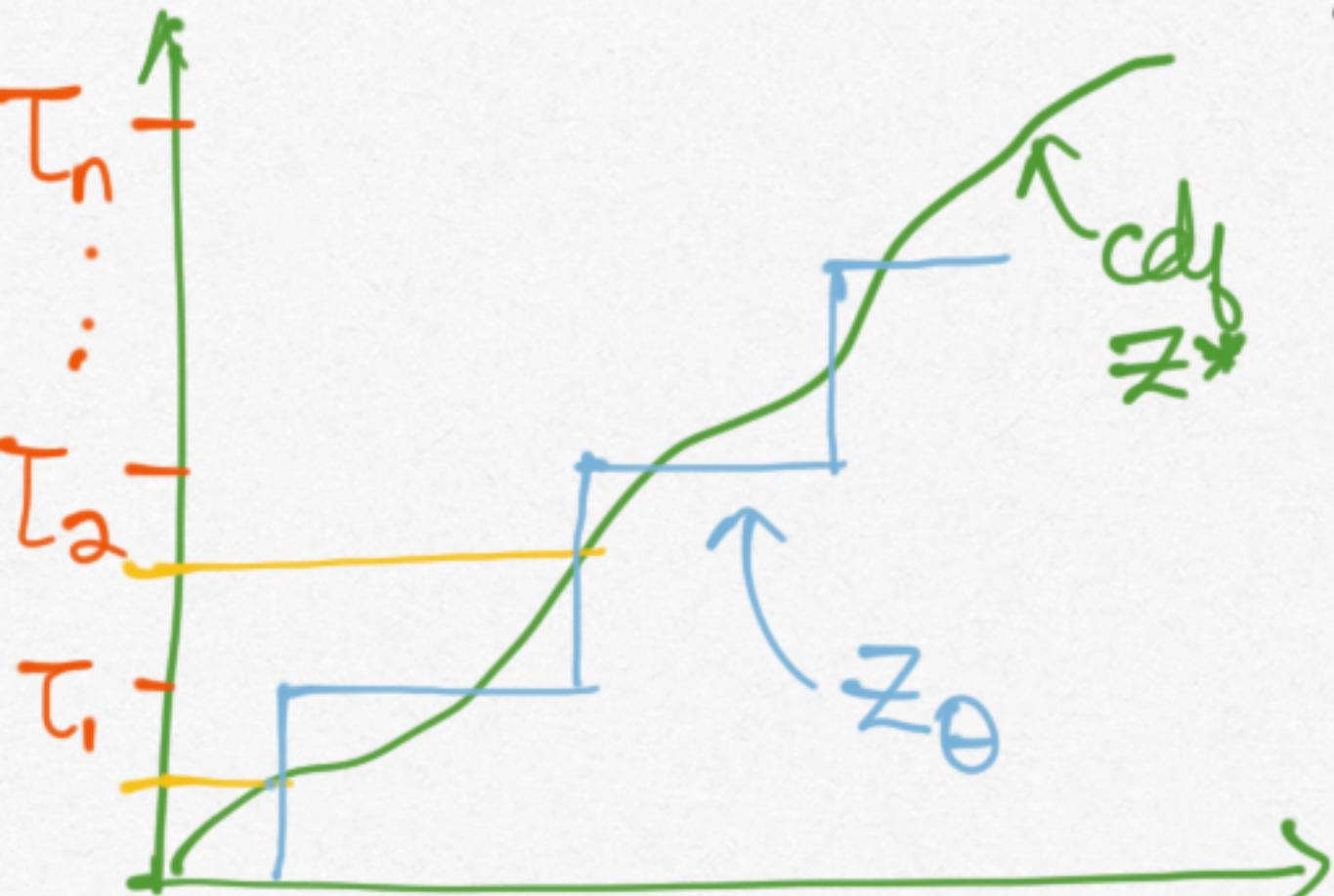
- Similarly to select actions, use K samples of $\tilde{\tau}$ & take argmax of mean $Q_{\beta(\tilde{\tau})}(s, a)$
- How to generate these samples?
Sample from $U(0, 1)$ & reweight using the parametric f .

QR-DQN

1. We want to learn distribution of return Z
2. We can start with estimate \hat{Z}_θ (θ : model params) & iteratively improve using bellman operator style updates.
3. It was proved that p-Wass. distance is a γ -contraction mapping & the distributional Bellman operator $T^\pi z : \mathcal{R} + \gamma \mathcal{Z}(s', a')$ can converge to a fixed point (using p-Wass in its dual form)

4. $d_p(z_1, z_2) = \sup_{s, a} W_p(z_1(s, a), z_2(s, a))$

5. Consider some true distribution we want to learn z^* . If we want to use Wass. dist to learn, Wass. dist is the max difference in CDF's of the 2 distributions.



What is the best approx distribution?

If we divide cdf into equal prob intervals, say N , then $T_i = i/N$. It turns out that best approx Z_θ is at quantile midpoints $T' = \frac{T_{i-1} + T_i}{2}$

So we can learn quantiles using QR.

Quantile Regression

Linear regression estimates "conditional mean" (mean/expectation over set of conditions)

If Least absolute deviatⁿ is used (MAE), we estimate median of target.

QR - we estimate conditional quantiles. (Median is special case, $q = 0.5$)

Given R.V. X : drive time Everett to Seattle

$q \in (0, 1)$, $P(X \leq x) = q$

say $P(X \leq 42 \text{ min}) = 0.95$ implies

Prob. of getting to Seattle in ≤ 42 min is 95%

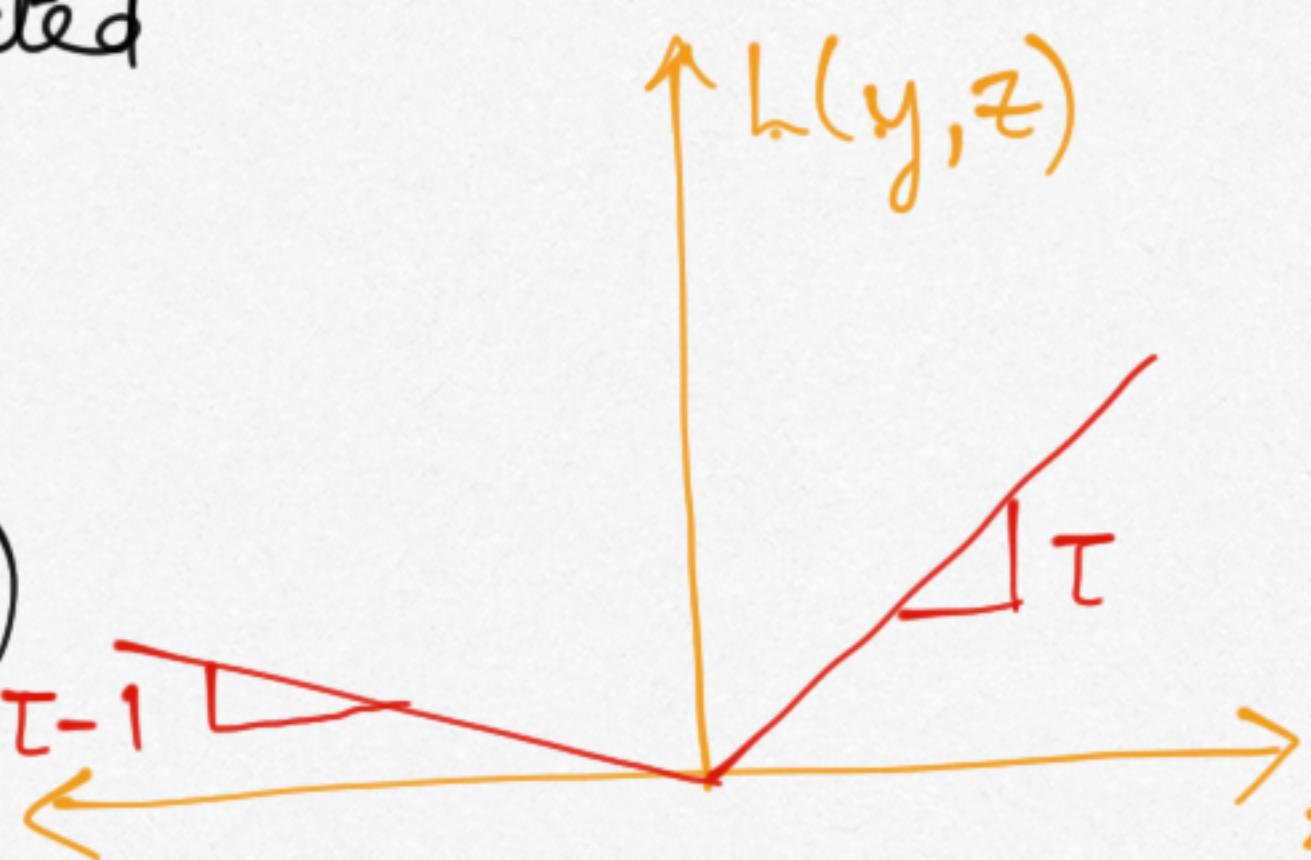
QR Loss (pinball loss)

τ -target quantile, y -actual, z -predicted

$$L_\tau(y-z) = \begin{cases} |y-z| & y \geq z \\ |z-y|(1-\tau) & y < z \end{cases}$$

or written as $\rho_\tau(u) = u(\tau - \mathbb{I}_{u<0})$

u -residual, \mathbb{I} -indicator function.



E.g. $\tau = 0.05$, 5%, quantile value = y of test scores.

Only 5% of students scored below y .

and 95% of students scored above y .

$$\text{Let } y = 10. \text{ If } z = 9, L_\tau = (10-9) \cdot (0.05) = 0.05$$

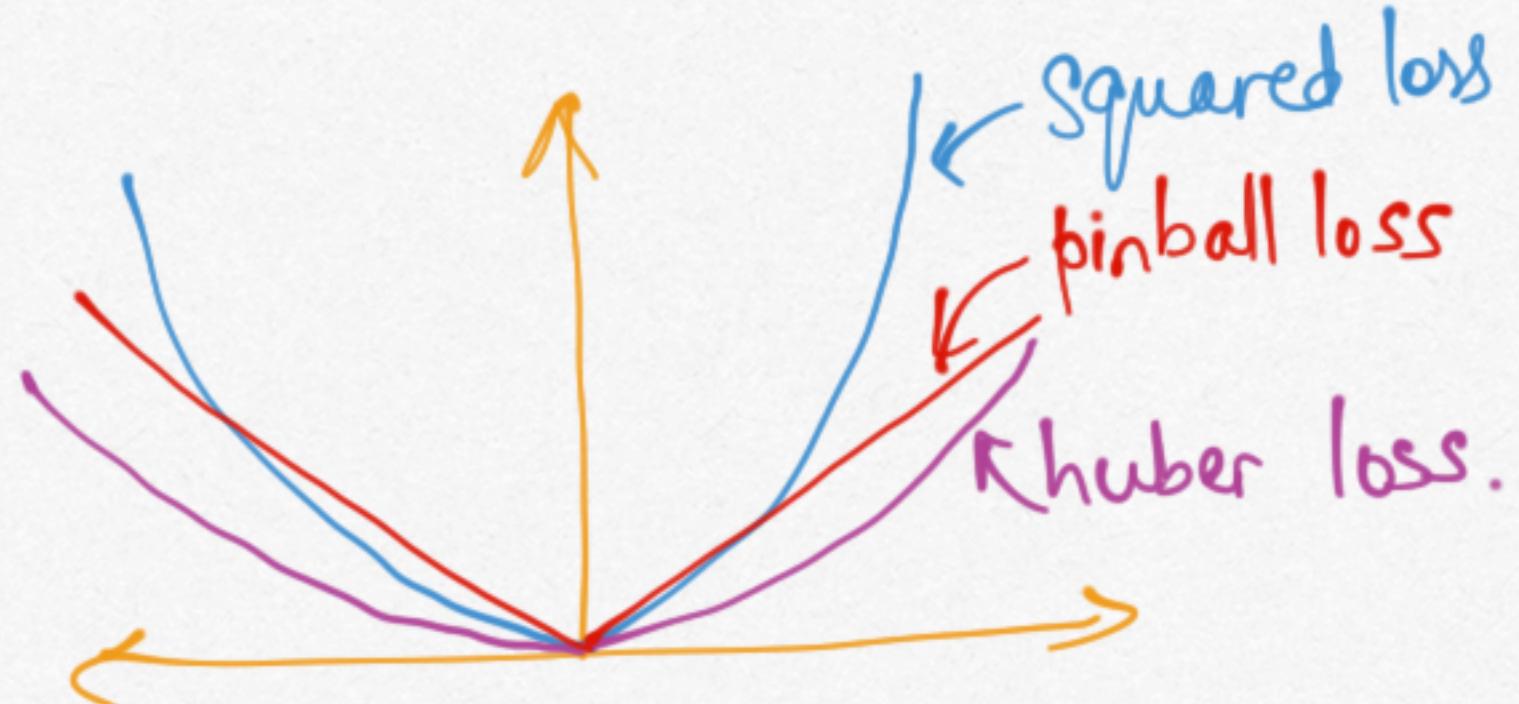
$$z = 11, L_\tau = (11-10) \cdot (1-0.05) = 0.95$$

L_τ is higher when we predict on "wrong" side of quantile.

Huber loss

QR loss is not zero everywhere.

When $u \rightarrow 0$, derivative is constant. So authors combine w/ huber loss.



$$l_k(u) = \begin{cases} \frac{1}{2}u^2 & |u| \leq k \\ k(|u| - \frac{1}{2}k) & \text{otherwise} \end{cases}$$

u -residual, $k=0$ or 1

$$\text{Quantile huber loss} = |\tau - \mathbb{I}_{u<0}| l_k(u)$$

In practice we compute the loss f^n between all pairs (θ_i, θ_j) of quantiles of estimated & target distribution.

$$\text{i.e. } \rho_{\hat{\tau}_i}^k(T\theta_j - \theta_i(s,a))$$

$T\theta_j$: target (Bellman update)

$\theta_i(s,a)$: current estimate