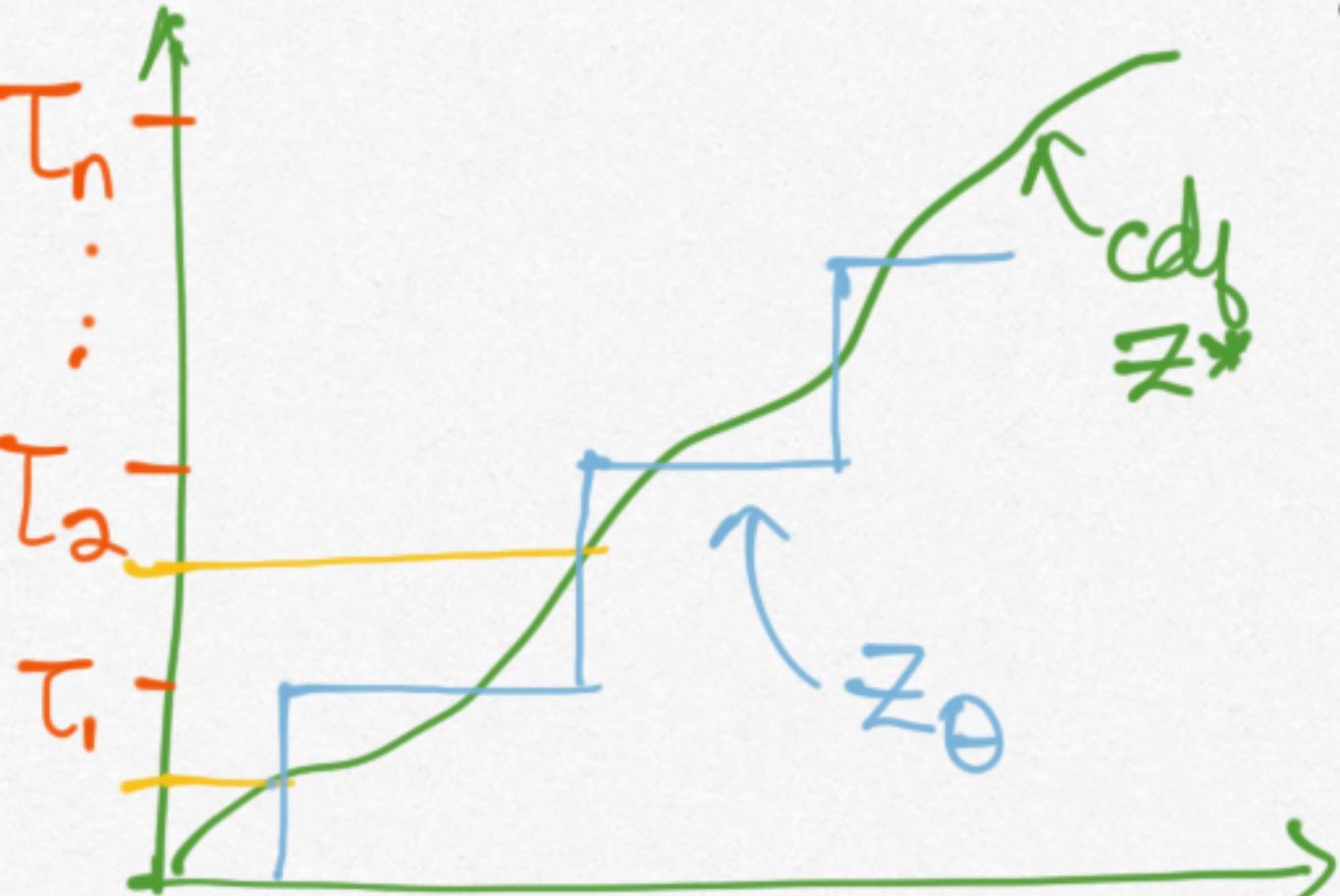


QR-DQN

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%

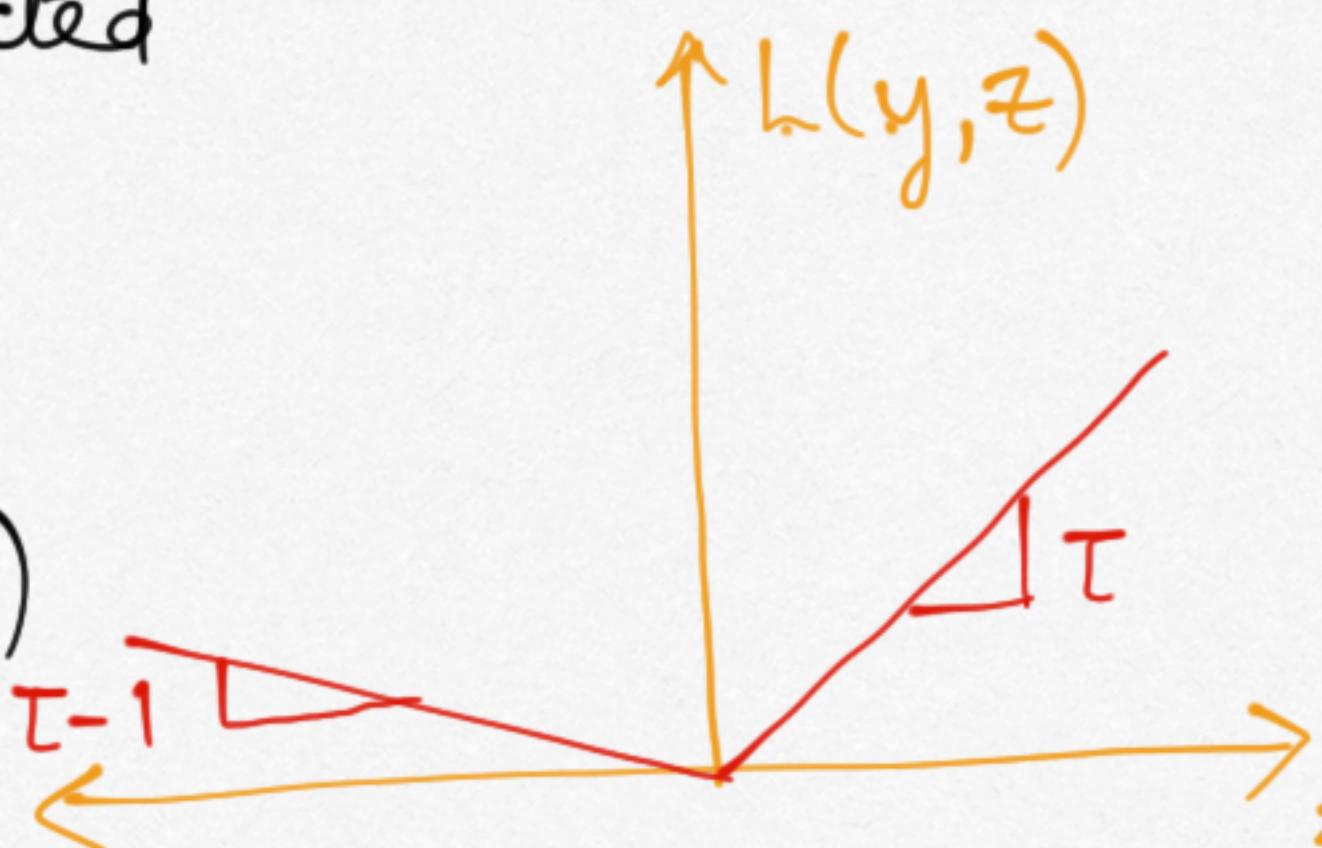
QR Loss (pinball loss)

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 .

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

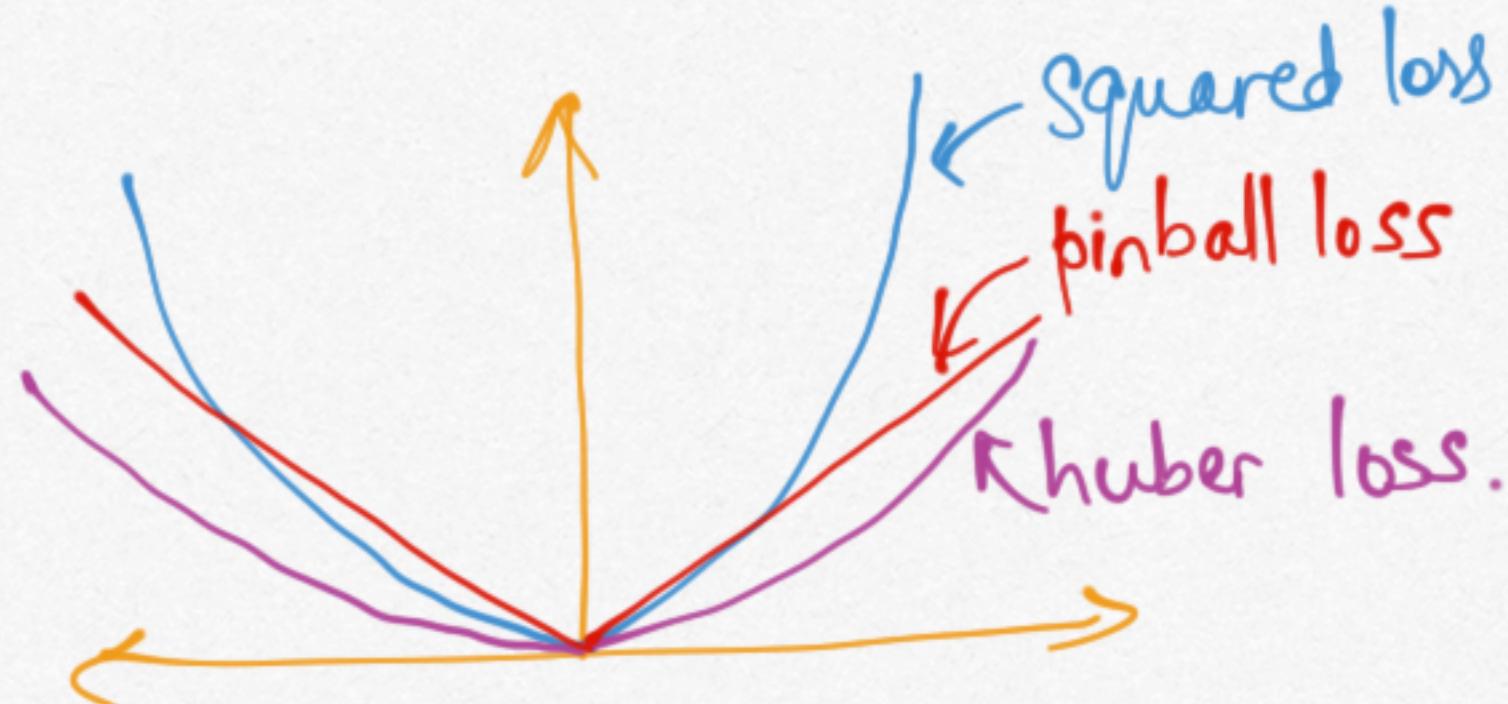
$$z = 11, L_\tau = (11-10) \cdot (1-0.05) = 0.95$$

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