## ⟨ QR-DQN: How I coded. ⟩

※ The first dim of tensors should be mb_size, but we are ignoring it.

target_q_dist = $\begin{pmatrix} target_1 & \cdots & target_N \\ \vdots & \ddots & \vdots \\ target_1 & & target_N \end{pmatrix}$    (N = # of quantiles)
= n_quant.

eval_q_dist = $\begin{pmatrix} eval_1 & \cdots & eval_1 \\ \vdots & \ddots & \vdots \\ eval_N & & eval_N \end{pmatrix}$

If we broadcast it,
target_q_dist
$\in \mathbb{R}^{mb \times 1 \times N\_quant}$ (target)
$\to \mathbb{R}^{mb \times N\_quant (eval) \times N\_quant (target)}$

eval_q_dist expands as well!

u_values = $\begin{pmatrix} target_1 - eval_1 & \cdots & target_N - eval_1 \\ \vdots & \ddots & \\ target_1 - eval_N & & target_N - eval_N \end{pmatrix}$
$\in \mathbb{R}^{mb \times N\_quant (eval) \times N\_quant (target)}$

tau_values = $\begin{pmatrix} \tau_1 & \cdots & \tau_1 \\ \vdots & \ddots & \vdots \\ \tau_N & & \tau_N \end{pmatrix} \in \mathbb{R}^{mb \times N\_quant (eval) \times N\_quant (target)}$

u_values.le(0) = $\begin{pmatrix} \mathbb{1}(target_1 - eval_1 < 0) & \cdots & \mathbb{1}(target_N - eval_1 < 0) \\ \vdots & \ddots & \\ \mathbb{1}(target_1 - eval_N < 0) & & \mathbb{1}(target_N - eval_N < 0) \end{pmatrix}$
$\in \mathbb{R}^{mb \times N\_quant (eval) \times N\_quant (target)}$

weight = $\begin{pmatrix} |\tau_1 - \mathbb{1}(target_1 - eval_1 < 0)| & \cdots & |\tau_1 - \mathbb{1}(target_N - eval_1 < 0)| \\ \vdots & \ddots & \\ |\tau_N - \mathbb{1}(target_1 - eval_N < 0)| & & |\tau_N - \mathbb{1}(target_N - eval_N < 0)| \end{pmatrix}$
$\in \mathbb{R}^{mb \times N\_quant (eval) \times N\_quant (target)}$

Although N_quant (eval) = N_quant (target), we need to differentiate which dimension corresponds to which (eval/target).

$$\text{rho\_values} = \begin{pmatrix} \mathcal{L}_k(\text{eval}_1 - \text{target}_1) & \cdots & \mathcal{L}_k(\text{eval}_1 - \text{target}_N) \\ \vdots & \ddots & \\ \mathcal{L}_k(\text{eval}_N - \text{target}_1) & & \mathcal{L}_k(\text{eval}_N - \text{target}_N) \end{pmatrix}$$

$$\in \mathbb{R}^{mb \times N\_quant(\text{eval}) \times N\_quant(\text{target})}$$

weight * rho_values

$$= \begin{pmatrix} |\tau_1 - \mathbb{1}(\text{target}_1 - \text{eval}_1 < 0)| \cdot \mathcal{L}_k(\text{eval}_1 - \text{target}_1) & \cdots & |\tau_1 - \mathbb{1}(\text{target}_N - \text{eval}_1 < 0)| \\ & & \cdot \mathcal{L}_k(\text{eval}_1 - \text{target}_N) \\ \vdots & \ddots & \\ |\tau_N - \mathbb{1}(\text{target}_1 - \text{eval}_N < 0)| \cdot \mathcal{L}_k(\text{eval}_N - \text{target}_1) & & |\tau_N - \mathbb{1}(\text{target}_N - \text{eval}_N < 0)| \\ & & \cdot \mathcal{L}_k(\text{eval}_N - \text{target}_N) \end{pmatrix}$$

$$\in \mathbb{R}^{mb \times N\_quant(\text{eval}) \times N\_quant(\text{target})}$$

· Now we can compute our objective function

$$\sum_{i=1}^{N} E_j \left[ \rho_{\hat{\tau}_i}^k (T\theta_j - \theta_i)(x, a) \right]$$

eval_index    target_index                 (weight * rho_values)                 necessary!

$$= \sum_{i=1}^{N} E_j \left[ |\hat{\tau}_i - \mathbb{1}\{(T\theta_j - \theta_i)(x, a) < 0\}| \cdot \mathcal{L}_k \{(T\theta_j - \theta_i)(x, a)\} / R \right].$$

→ depends on sample. $(m = 1, \cdots, mb\_size)$

sum over eval $(\dim = 1)$          mean over target $(\dim = 2)$

python's zero-indexing          python's zero-indexing

$$\mathcal{L}_k(a) = \begin{bmatrix} \frac{1}{2} a^2 & |a| \le k \\ k(|a| - \frac{1}{2}k) & \text{o.w.} \end{bmatrix}$$

$$\frac{\mathcal{L}_k(x)}{k} \longrightarrow \text{approximates } |a| \text{ better!}$$



$\frac{1}{2}k^2$   slope = $k \ll 1$.

$-k$   $k$

$k$   $\frac{1}{2}k$   → difference $= \frac{k}{2} \approx 0$

slope 1.

$-k$   $k$

$\curvearrowright$ "smoothly" (differentiable)

$\hookrightarrow$ Since the goal of Huber loss is to approximate $|x|$, which is the objective in Quantile-regression, it makes more sense to use $\mathcal{L}_k(x)/\kappa$ instead of $\mathcal{L}_k(x)$.

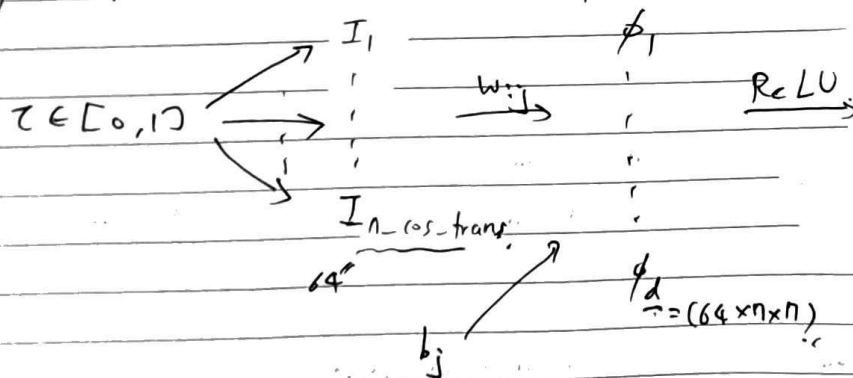⟨IQN: understanding the mechanism⟩ — ignoring the first dimension (mb_size).

· Let $\Psi: \mathbb{R}^{4\times84\times84} \longrightarrow \mathbb{R}^{(64\times n\times n)}$ : 3 repetitive pairs of (CNN, ReLU),
  $\quad$ (3d-tensor) $\quad$ (1d-tensor) followed by Flatten $\to$ (1d-tensor)

$\quad f: \mathbb{R}^{(64\times n\times n)} \longrightarrow \mathbb{R}^{|A|}$ : 2 fully-connected layers each being
  $\qquad\qquad\qquad\qquad\qquad f_1: \mathbb{R}^{(64\times n\times n)} \to \mathbb{R}^{512}, \quad f_2: \mathbb{R}^{512} \to \mathbb{R}^{|A|}$.
  $\qquad\qquad\qquad\qquad\qquad (f = f_2 \circ f_1)$

$\quad \phi: [0,1] \longrightarrow \mathbb{R}^{(64\times n\times n)}$ : to be explained later on

· For a fixed value of $\tau_0 \in [0,1]$,

$$Z_{\tau_0}(x,a) = f\left(\Psi(x) \odot \phi(\tau_0)\right) \quad (\in \mathbb{R}^{|A|})$$
$$= \left(\underbrace{f_1(\Psi(x)\odot\phi(\tau_0))}_{\in \mathbb{R}}, \cdots, \underbrace{f_{|A|}(\Psi(x)\odot\phi(\tau_0))}_{\in\mathbb{R}}\right)^T \quad (\in \mathbb{R}^{|A|})$$

· $\phi: [0,1] \to \mathbb{R}^{(64\times n\times n)}$ can be explained as follows,

$\tau\in[0,1]$
$\quad I_1 \quad\quad\quad\quad \phi_1$
$\quad\quad\quad\quad w_{ij} \quad\quad\quad ReLU$
$\quad I_{n\_cos\_trans}$
$\quad 64$
$\quad b_j$
$\quad\quad\quad \phi_d$
$\quad\quad\quad = (64\times n\times n)$

where $I_i = \cos(\pi i \cdot \tau)$

$(w_{ij})_{\substack{i=1,\cdots, n\_cos\_trans \\ j=1,\cdots, d}}$

$(b_j)_{j=1,\cdots,d}$

· Then our final NN can be summarized as

$$F: \mathbb{R}^{4 \times 84 \times 84} \times [0,1] \longrightarrow \mathbb{R}^{|A|}$$, which is defined as follows

there will be mb_size samples of images

there will be N_quant (either eval or target) samples.
for online network, for target network

$$F(x: \tau) = Z_\tau(x) = [Z_\tau(x, a_1) \cdots Z_\tau(x, a_{|A|}))]^T (\in \mathbb{R}^{|A|})$$
↳ $\tau$-th quantile of $Z(x, a_1)$

$$= \left( f_1(4(x) \odot \phi(\tau)), \cdots, f_{|A|}(4(x) \odot \phi(\tau)) \right)^T (\in \mathbb{R}^{|A|})$$

· Note that all the fcts $f = (f_1 \cdots f_{|A|})^T$, $4$, $\phi$ are NN's that we should train.

Since there are "mb_size" many samples of $x_i$,
and "N_quant" (eval or target) many samples of $\tau_i$,
there are (mb_size × N_quant) many samples that we need
↳ or (mb_size × N_quant (eval) × N_quant (target))
to be mindful of.


· cos_trans

|  | $i=0$ | . . . | $i = n\_cos\_trans - 1$ |
|---|---|---|---|
| $\tau_1$ | $\cos(\pi \cdot 0 \cdot \tau_1)$ | . . . $\cos(\pi \cdot (n-1) \cdot \tau_1)$ | |
| ⋮ | ⋮ | ⋱ | |
| $\tau_{N\_quant}$ | $\cos(\pi \cdot 0 \cdot \tau_{N\_quant})$ | $\cos(\pi \cdot (n-1) \cdot \tau_{N\_quant})$ | |

$\in \mathbb{R}^{N\_quant \times n\_cos\_trans}$

↳ treated as # of samples m    network within $\phi: [0,1] \rightarrow \mathbb{R}^{n\_cos\_trans} \rightarrow \mathbb{R}$

deterministic fct, network part
that we don't learn.