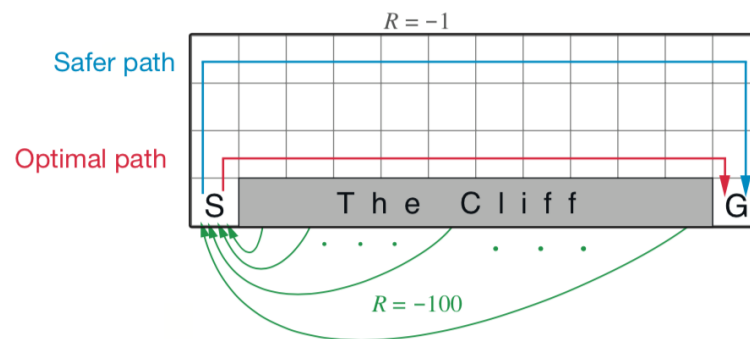




Обучение с подкреплением на распределениях

Осминин Константин
10 сентября 2019

- ✓ Более точная оценка вознаграждений
- ✓ Допускает мультимодальность вознаграждений
- ✓ Допускает управление риском
- ✓ Более стабильное обучение





Учим распределения, а не среднее

$$Q(x, a) = \mathbb{E}_{\pi} \sum_i \gamma^i R(x_i, a_i)$$

$$\mathbb{E}_{\pi} \sum_i \gamma^i R(x_i, a_i) \sim Z(x, a)$$

$$Q(x, a) = \mathbb{E} R(x, a) + \gamma \mathbb{E} Q(X', A').$$

$$Z(x, a) \stackrel{D}{=} R(x, a) + \gamma Z(X', A').$$

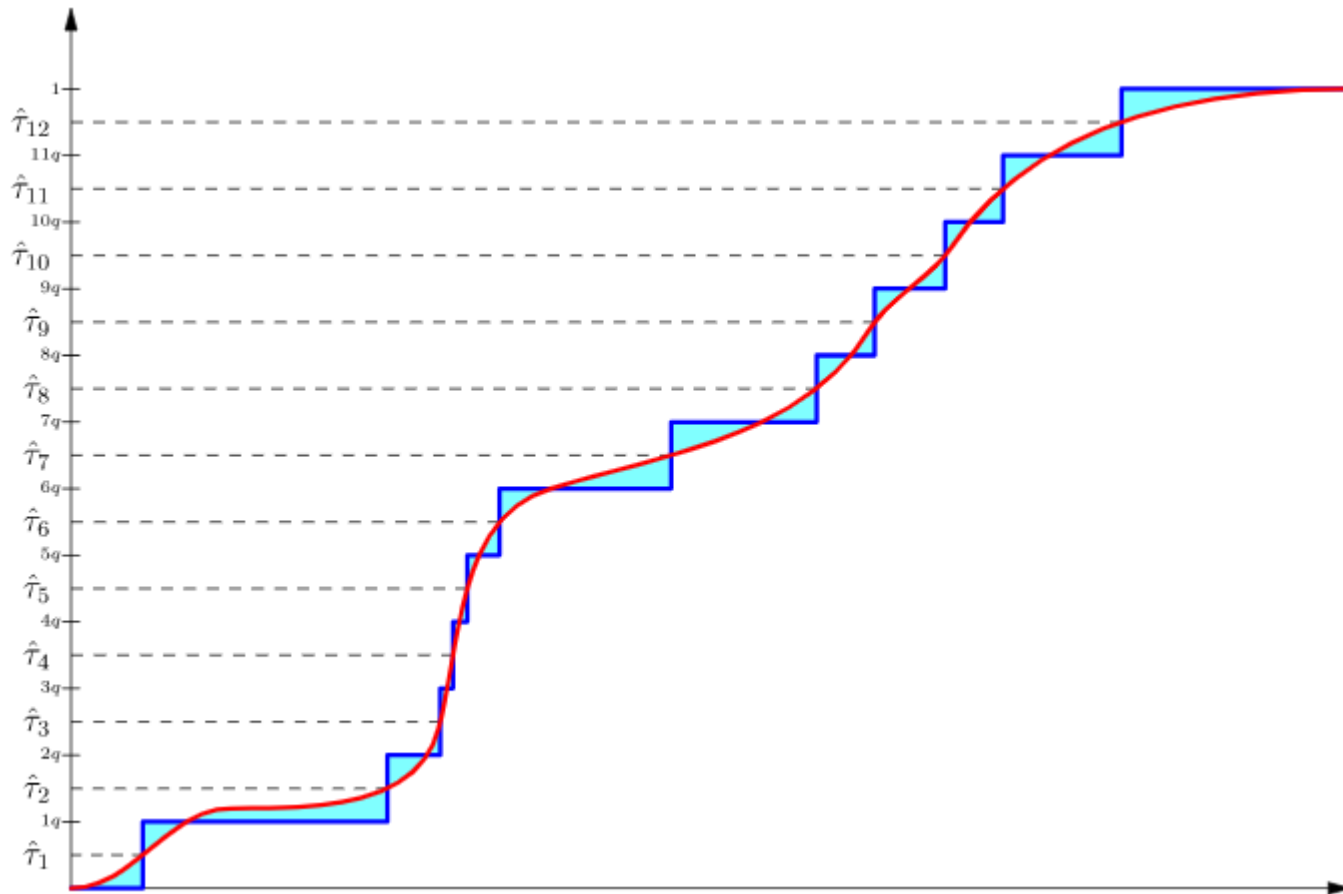


по распределению

$$Q(x, a) = \mathbb{E} Z(x, a)$$



Как учить распределение



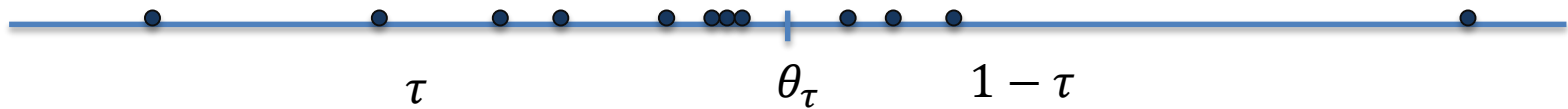
Учим N квантилей $\{\theta_i = F^{-1}(i/N)\}$



Квантильная регрессия

θ_τ есть τ -квантиль выборки $X \Rightarrow$

$$\theta_\tau = \operatorname{argmax} L_\theta = \operatorname{argmax} \mathbb{E}_X |X - \theta| \cdot |\tau - \delta_{X < \theta}|$$



$$\nabla_\theta L_\theta = \begin{cases} 1 - \tau, X < \theta \\ \tau, X \geq \theta \end{cases}$$



Как учить квантили

- ✓ Используем квантильную регрессию, только для гладкости применяем Huber Loss

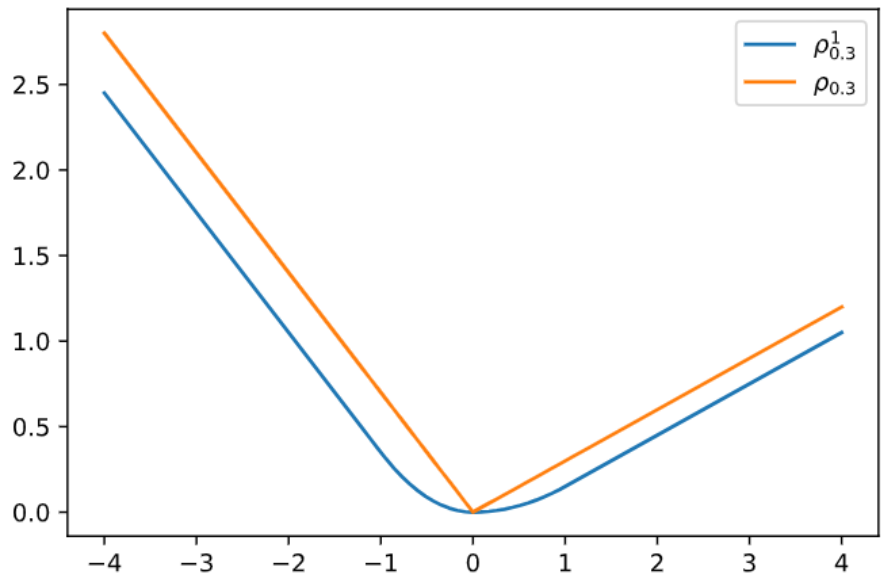
$$\mathcal{L}_\kappa(u) = \begin{cases} \frac{1}{2}u^2 & \text{if } |u| \leq \kappa \\ \kappa(|u| - \frac{1}{2}\kappa) & \text{otherwise} \end{cases}$$

- ✓ Для квантили τ

$$L_\theta = \mathbb{E}_X \mathcal{L}_\kappa(X - \theta) \cdot |\tau - \delta_{X < \theta}|$$

- ✓ Для всех квантилей

$$L_\theta = \sum_{i=1}^N \mathbb{E}_X \mathcal{L}_\kappa(X - \theta) \cdot |\tau_i - \delta_{X < \theta}|$$





Алгоритм QR-DQN

- ✓ $s = s_0$
- ✓ Пока не сошлись:
 - ✓ $a = \operatorname{argmax}_{a'} \frac{1}{N} \sum_{j=1}^N \theta_j(s, a')$
 - ✓ Из состояния s делаем шаг a и наблюдаем s', r .
 - ✓ Обучение:
 - ✓ $a' = \operatorname{argmax}_{a'} \frac{1}{N} \sum_{j=1}^N \theta_j(s, a')$
 - ✓ $\theta_j^* = r + \gamma \theta_j(s', a')$
 - ✓ $L = \frac{1}{N^2} \sum_{i,j=1}^N \mathcal{L}_k(\theta_j^* - \theta(x, a)) \cdot |i/N - \delta_{x < \theta(x, a)}|$
 - ✓ Шаг градиентного спуска по параметрам $\theta(x, a)$
- ✓ $s \leftarrow s'$

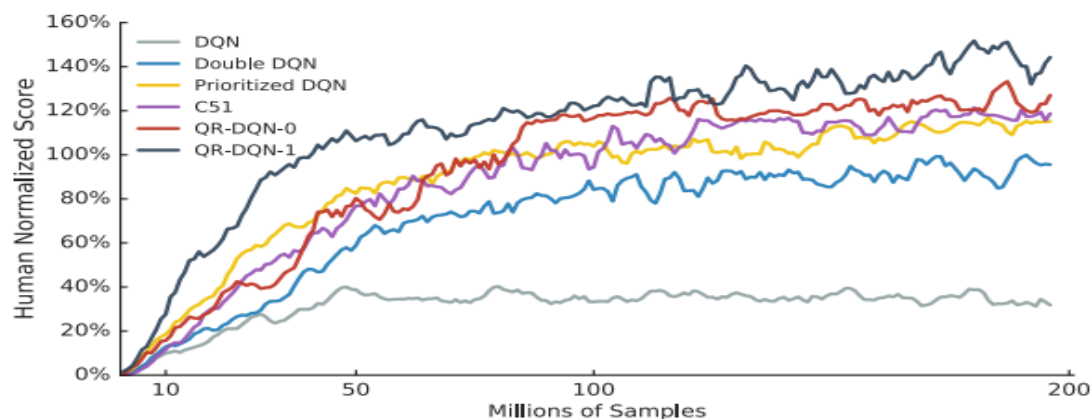


Figure 4: Online evaluation results, in human-normalized scores, over 57 Atari 2600 games for 200m train samples.

	Mean	Median	>human	>DQN
DQN	228%	79%	24	0
DDQN	307%	118%	33	43
DUEL.	373%	151%	37	50
PRIOR.	434%	124%	39	48
PR. DUEL.	592%	172%	39	44
C51	701%	178%	40	50
QR-DQN-0	881%	199%	38	52
QR-DQN-1	915%	211%	41	54

Table 1: Mean and median of *best* scores across 57 Atari 2600 games, measured as percentages of human baseline (Nair et al. 2015).



- ✓ Marc G Bellemare, Will Dabney, and Rémi Munos, A distributional perspective on reinforcement learning, arXiv preprint arXiv:1707.06887 (2017).
- ✓ Will Dabney, Mark Rowland, Marc G Bellemare, and Rémi Munos, Distributional reinforcement learning with quantile regression, arXiv preprint arXiv:1710.10044 (2017).



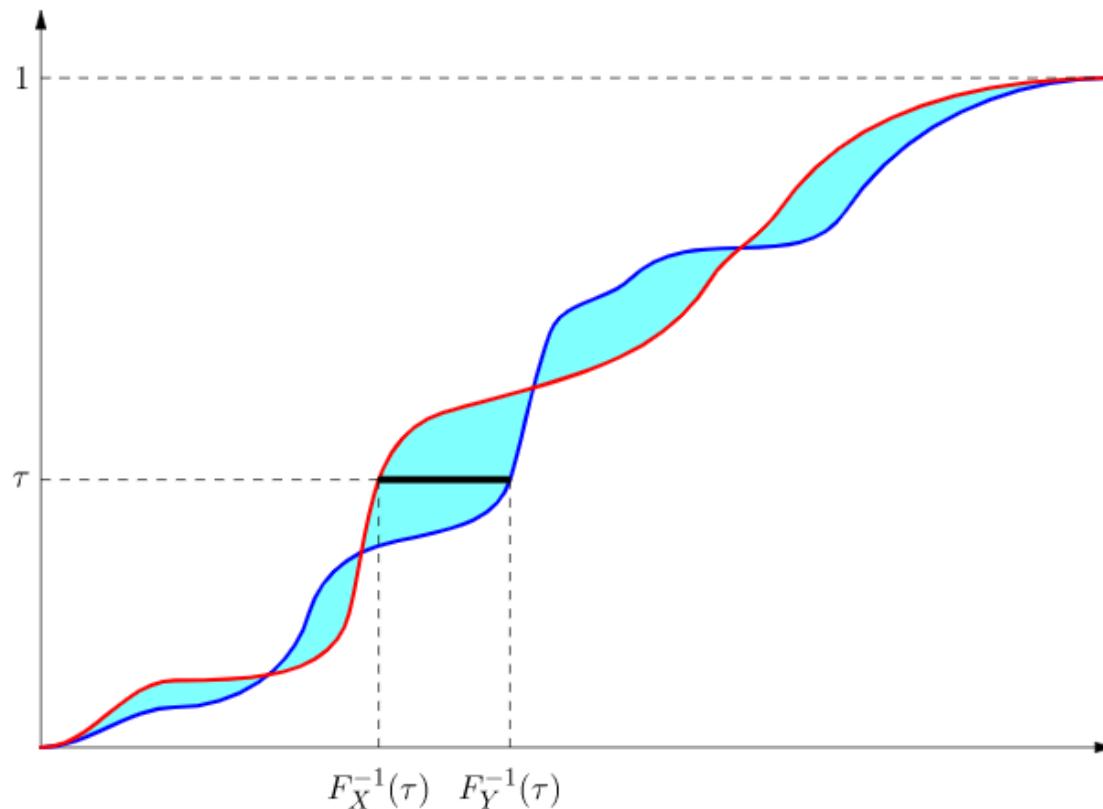
https://github.com/bayesgroup/deepbayes-2018/blob/master/day3_qr-qnetwork/qr_dqn-local.ipynb



Метрика Вассерштайна

Let X and Y be two *scalar* random variables and F_X and F_Y their CDFs. Then, their p -Wasserstein distance is

$$\mathcal{W}_p(X, Y) = \left(\int_0^1 |F_X^{-1}(u) - F_Y^{-1}(u)|^p du \right)^{1/p}$$





Спасибо