
Algorithm 1: Distributional Soft Actor-Critic (DSAC)

Input: Initial policy parameters ϕ , quantile value network parameters θ_1, θ_2 , target networks $\bar{\phi}, \bar{\theta}_1, \bar{\theta}_2$, polyak coefficient ι , number of quantiles N , Huber loss parameter κ , temperature α , discount γ

Data: Batch of transitions (s, a, r, s')

1 **for** each batch **do**

2 Generate quantile fractions $\tau_i, i = 0, \dots, N$ and $\hat{\tau}_j, j = 0, \dots, N$;

 // Update Quantile Value Networks $Z(s, a; \theta_k), k = 1, 2$

3 Get next actions for target: $a' \sim \pi(\cdot | s'; \bar{\phi})$;

4 **for** $i = 0$ **to** $N - 1$ **do**

5 **for** $j = 0$ **to** $N - 1$ **do**

6 Compute target quantile:

$$y_i = \min_{k=1,2} Z_{\hat{\tau}_i}(s', a'; \bar{\theta}_k)$$

 Compute TD error:

$$\delta_{ij}^k = r + \gamma [y_i - \alpha \log \pi(a' | s'; \bar{\phi})] - Z_{\hat{\tau}_j}(s, a; \theta_k), \quad k = 1, 2$$

7 **end**

8 **end**

9 Compute quantile regression loss:

$$\mathcal{J}_Z(\theta_k) = \frac{1}{N} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} (\tau_{i+1} - \tau_i) \rho_{\hat{\tau}_j}^{\kappa}(\delta_{ij}^k), \quad k = 1, 2$$

10 Update quantile network parameters:

$$\theta_k \leftarrow \theta_k - \lambda_Q \nabla \mathcal{J}_Z(\theta_k), \quad k = 1, 2;$$

11 Update target quantile networks: $\bar{\theta}_k \leftarrow \iota \theta_k + (1 - \iota) \bar{\theta}_k, k = 1, 2$;

 // Update Policy Network $\pi(a | s; \phi)$

12 Sample reparameterized actions: $\tilde{a} \sim \pi(\cdot | s; \phi)$;

13 Compute expected Q-value:

$$Q(s, \tilde{a}) = \sum_{i=0}^{N-1} (\tau_{i+1} - \tau_i) \min_{k=1,2} Z_{\hat{\tau}_i}(s, \tilde{a}; \theta_k)$$

14 Compute policy loss: $\mathcal{J}_\pi(\phi) = \alpha \log \pi(\tilde{a} | s; \phi) - Q(s, \tilde{a})$;

15 Update policy parameters: $\phi \leftarrow \phi + \lambda_\pi \nabla \mathcal{J}_\pi(\phi)$;

16 Update target policy: $\bar{\phi} \leftarrow \iota \phi + (1 - \iota) \bar{\phi}$;

17 **end**
