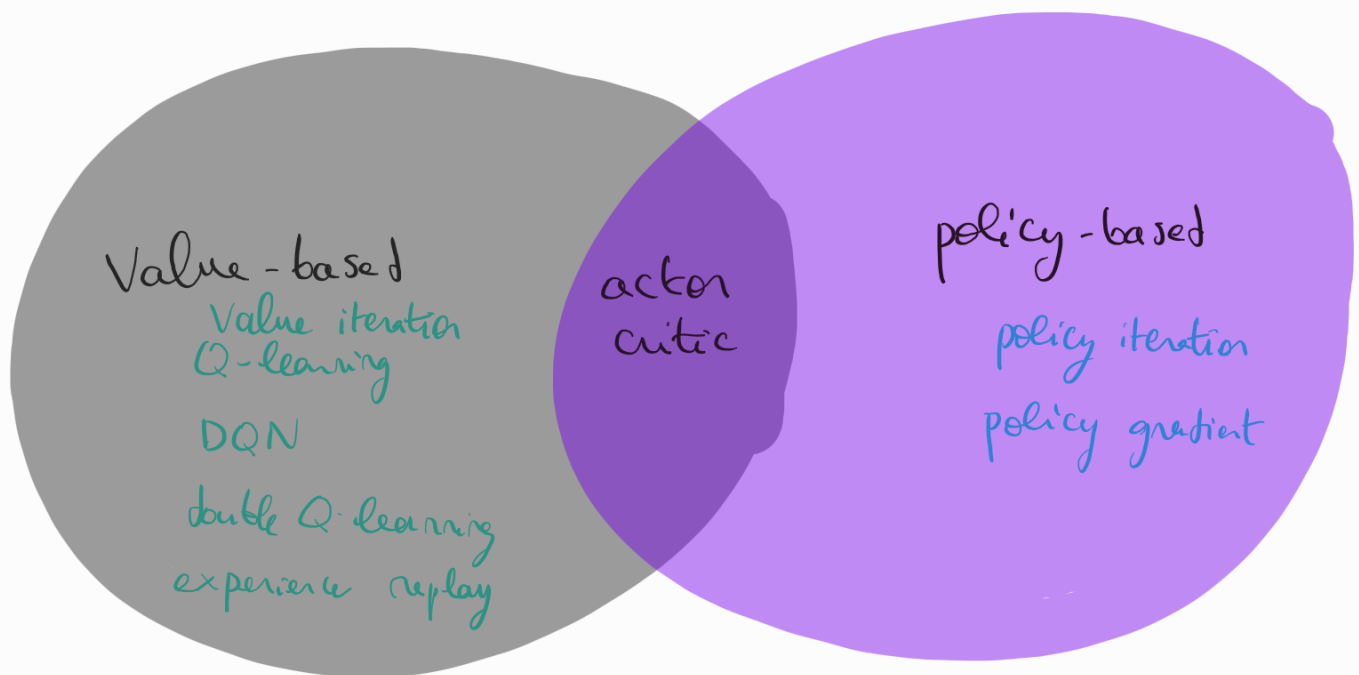


ACTOR CRITIC METHODS



Value-based :

- Monte Carlo
- Temporal difference

Policy-based :

policy gradient

TD :

$v_w : S \rightarrow \mathbb{R}$ represented by neural network

w : parameters

current state s

choose a using ϵ -greedy from $v_w(s)$

sample (s, a, r, s')

$$\delta \leftarrow r + \gamma v_w(s') - v_w(s)$$

$$w \leftarrow w + \alpha \delta \nabla_w v_w(s)$$

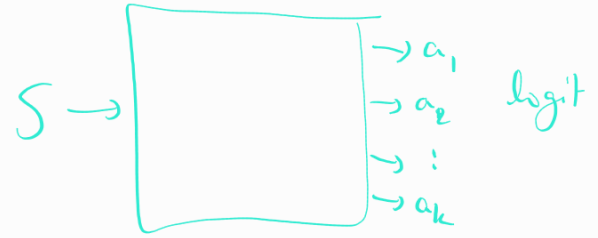
PG

$\pi_{\theta} : S \rightarrow \mathcal{O}(A)$ represented by neural network

θ : parameters

current state s

choose a using $\pi_{\theta}(s)$



sample (s, a, r, s')

$$\delta \leftarrow r + \gamma \max_{a'} \pi_{\theta}(s', a') - \pi_{\theta}(s, a)$$

$$\theta \leftarrow \theta + \beta \delta \nabla_{\theta} \log \pi_{\theta}(a|s)$$

Terminology:

Critic: value function

Actor: policy

One step Actor Critic

$v_w : S \rightarrow \mathbb{R}$ represented by neural network
 w : parameters

$\pi_\theta : S \rightarrow \mathcal{O}(A)$ represented by neural network
 θ : parameters

current state s

choose a using $\pi_\theta(s)$

sample (s, a, r, s')

$$\delta \leftarrow r + \gamma v_w(s') - v_w(s)$$

$$w \leftarrow w + \alpha \delta \nabla_w v_w(s) \quad \text{TD}$$

$$\theta \leftarrow \theta + \beta \delta \nabla_\theta \log \pi_\theta(a|s) \quad \text{PG}$$

INCREASING ROBUSTNESS WITH TRUST REGIONS

Regularisation

Kullback-Leibler divergence:

$$KL(\pi_{\phi_1} \parallel \pi_{\theta}) = \mathbb{E}_S \left[\sum_a \pi_{\phi_1}(a|s) \log \frac{\pi_{\theta}(a|s)}{\pi_{\phi_1}(a|s)} \right]$$

new objective:

$$\mathbb{E}_{\pi} [R(\pi|\theta)] - \eta KL(\pi_{\phi_1} \parallel \pi_{\theta})$$

TRPO / PPO