

RL with Continuous Action Sets

Chen-Yu Wei

3 main challenges in online RL: Exploration, Generalization, (Temporal) Credit Assignment

+ Generalization over actions

$$\mu_{\theta}(x) \approx \operatorname{argmax}_a R_{\phi}(x, a)$$

Finite actions

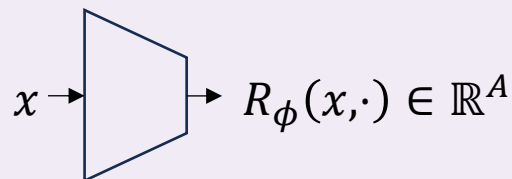
Infinite actions

MAB

Exploration

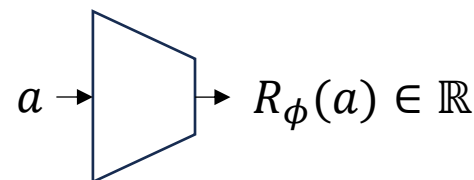
CB

+ Generalization over contexts



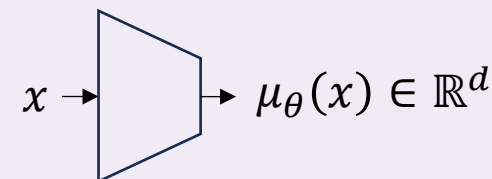
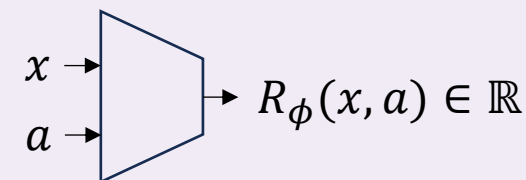
DQN

MAB



$\mu \in \mathbb{R}^d$

CB



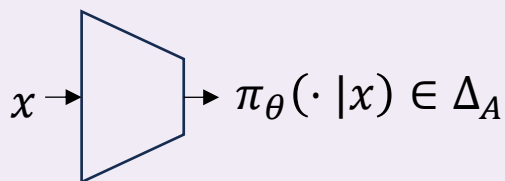
DDPG, TD3, SAC

VB

$R(\cdot) \in \mathbb{R}^A$

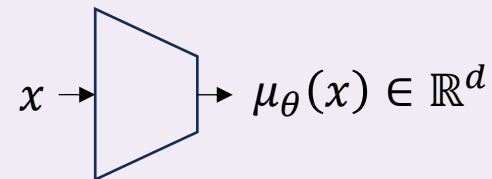
PB

$\pi(\cdot) \in \Delta_A$



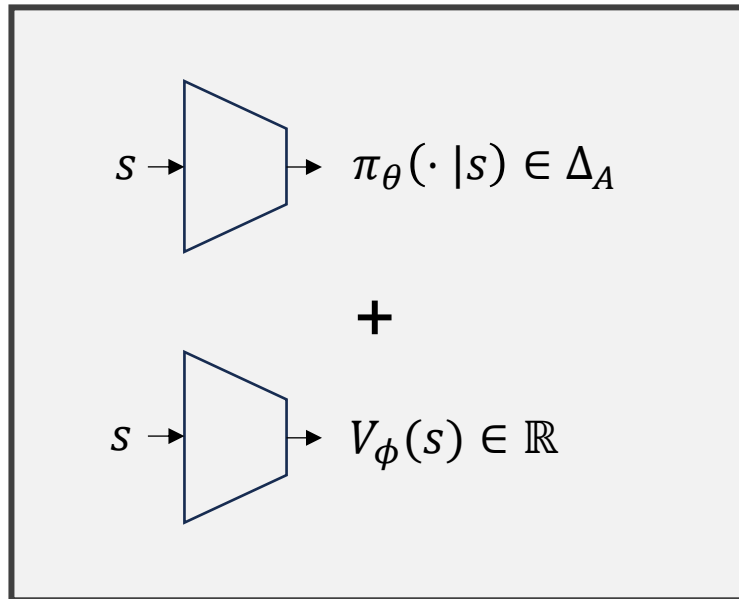
PPO, PG, A2C

$\mu \in \mathbb{R}^d$



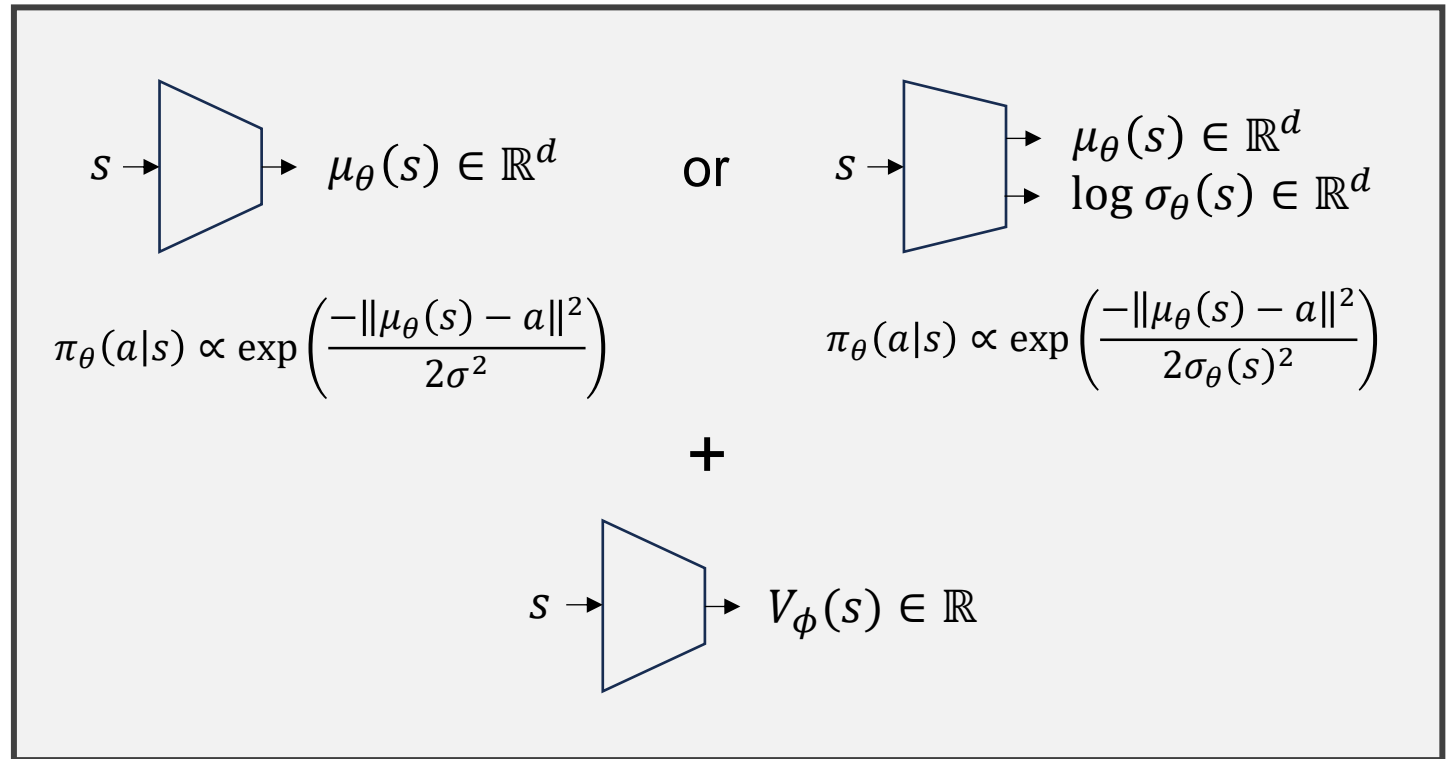
PPO, PG, A2C

PPO / PG / A2C in Discrete / Continuous Action Sets



Discrete actions

$\pi_{\theta}(a|s)$



Continuous actions

Algorithms involving a policy and value network where the value is used in the policy update are called **actor-critic** algorithms.

PPO / PG / A2C in Continuous Action Sets

$$\theta_{k+1} \leftarrow \operatorname{argmax}_{\theta} \left\{ \sum_{i=1}^N \left(\frac{\pi_{\theta}(a_i | s_i)}{\pi_{\theta_k}(a_i | s_i)} A_i - \frac{1}{\eta} \operatorname{KL} \left(\pi_{\theta}(\cdot | s_i), \pi_{\theta_k}(\cdot | s_i) \right) \right) \right\} \quad \text{PPO}$$

$$\theta_{k+1} \leftarrow \theta_k + \eta \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(a_i | s_i) \Big|_{\theta=\theta_k} A_i \quad \text{PG, A2C}$$

where A_i is a weighted average of the following (GAE):

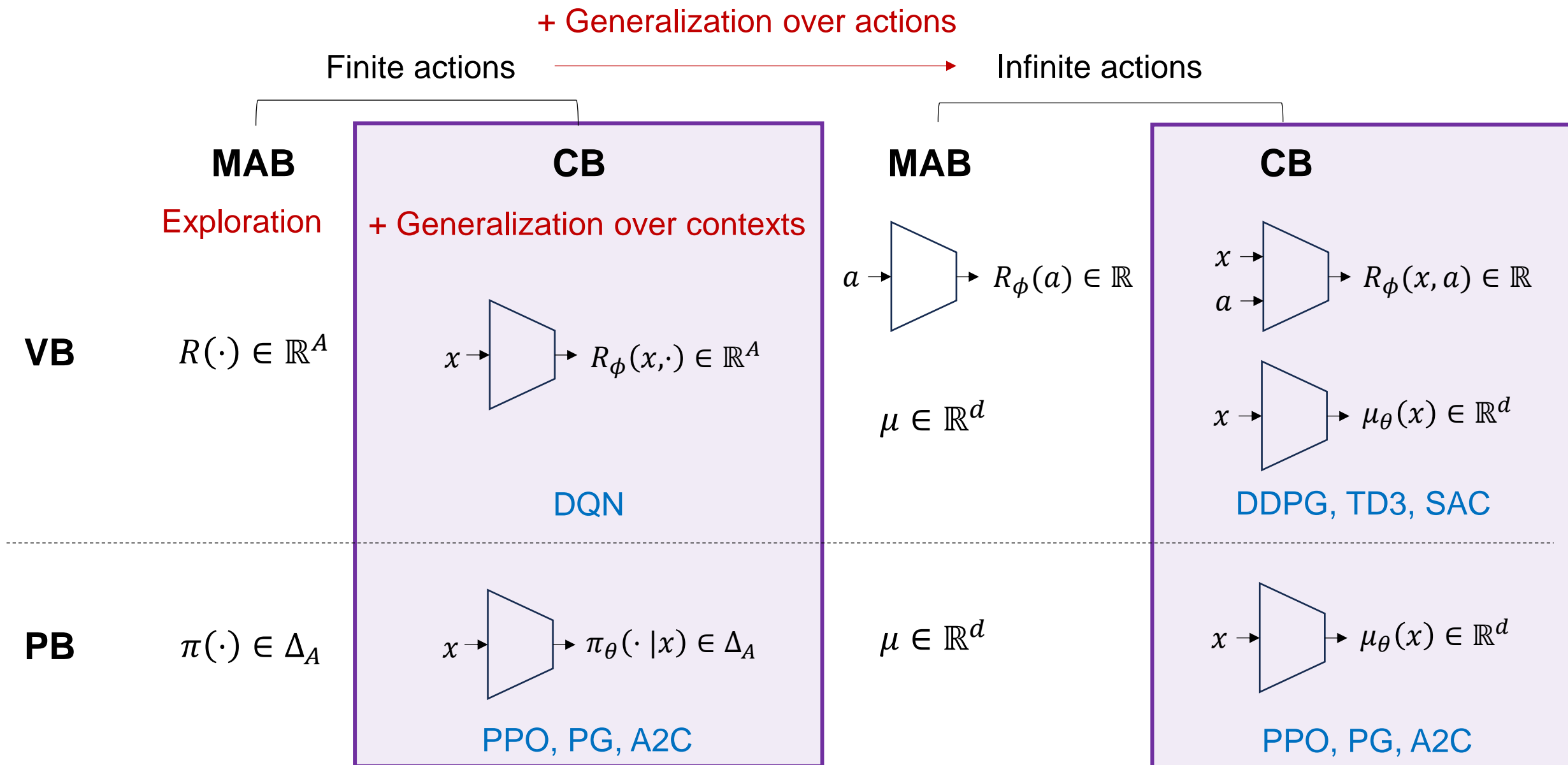
$$r_i + \gamma V_{\phi}(s_{i+1}) - V_{\phi}(s_i)$$

$$r_i + \gamma r_{i+1} + \gamma^2 V_{\phi}(s_{i+2}) - V_{\phi}(s_i)$$

$$r_i + \gamma r_{i+1} + \gamma^2 r_{i+2} + \gamma^3 V_{\phi}(s_{i+3}) - V_{\phi}(s_i)$$

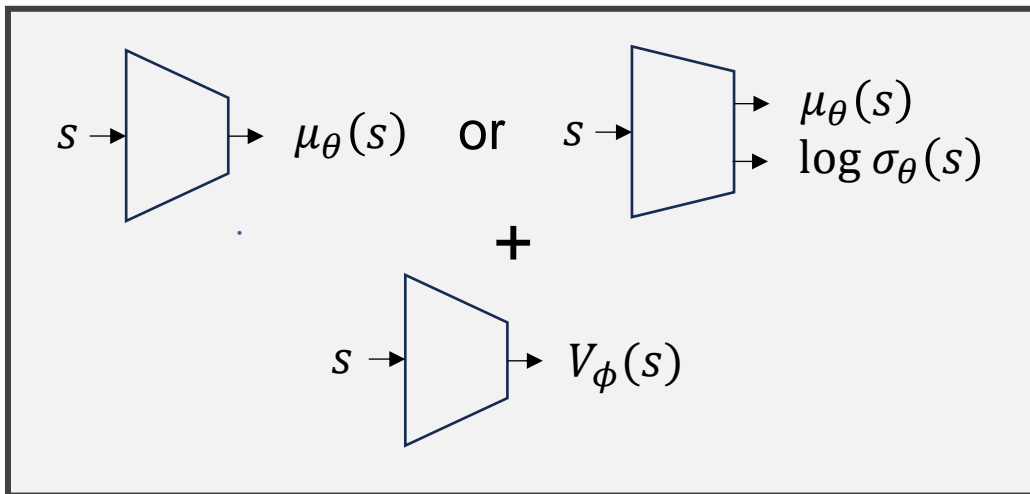
...

3 main challenges in online RL: Exploration, Generalization, (Temporal) Credit Assignment



Q^* Q^z

Two Types of Actor-Critic Algorithms



PPO / PG / A2C

Update θ with

$$\frac{\pi_{\theta}(a_i|s_i)}{\pi_{\theta_k}(a_i|s_i)} (r_i + \gamma V_{\phi}(s_{i+1}) - V_{\phi}(s_i))$$

Handwritten notes: $+ \gamma V_{\phi}(s_{i+1}) + \gamma^2 V_{\phi}(s_{i+2}) \dots$ and Q^z under the term $V_{\phi}(s_{i+1})$.

Idea more aligned with

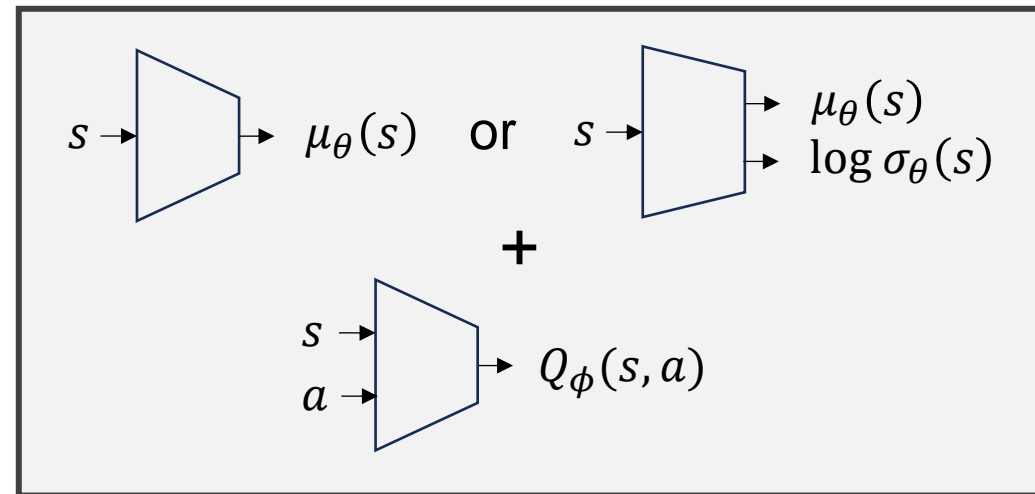
Policy-based bandits (forming unbiased reward estimator)

Policy Iteration (policy improvement based on $Q^{\pi}(s, a)$)

Training type

On-policy

V^{π} Q^z



DDPG / TD3 / SAC

$$\sum_a \pi_{\theta}(a|s_i) Q_{\phi}(s_i, a)$$

Value-based bandits (forming reward estimator from regression)

Policy Iteration or Value Iteration (policy improvement based on $Q^*(s, a)$) – e.g. DQN

On-policy or off-policy (using data collected from previous policies)

DDPG

Deep Deterministic Policy Gradient (DDPG)

For $k = 1, 2, \dots$

Use $\mu_\theta(s) + \mathcal{N}(0, \sigma^2)$ to collect samples and place them in **replay buffer**

Sample a batch $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ from the replay buffer

$$\phi \leftarrow \phi - \lambda \nabla_\phi \sum_{i=1}^n \left(Q_\phi(s_i, a_i) - r_i - \gamma \underbrace{Q_{\bar{\phi}}(s'_i, \mu_{\bar{\theta}}(s'_i))}_{\substack{\text{argmax}_a \\ Q_{\bar{\phi}}(s'_i, a)}} \right)^2$$

$$\theta \leftarrow \theta + \eta \sum_{i=1}^n \nabla_\theta Q_\phi(s_i, \mu_\theta(s_i))$$

$$\bar{\phi} \leftarrow \tau \phi + (1 - \tau) \bar{\phi}, \quad \bar{\theta} \leftarrow \tau \theta + (1 - \tau) \bar{\theta}$$

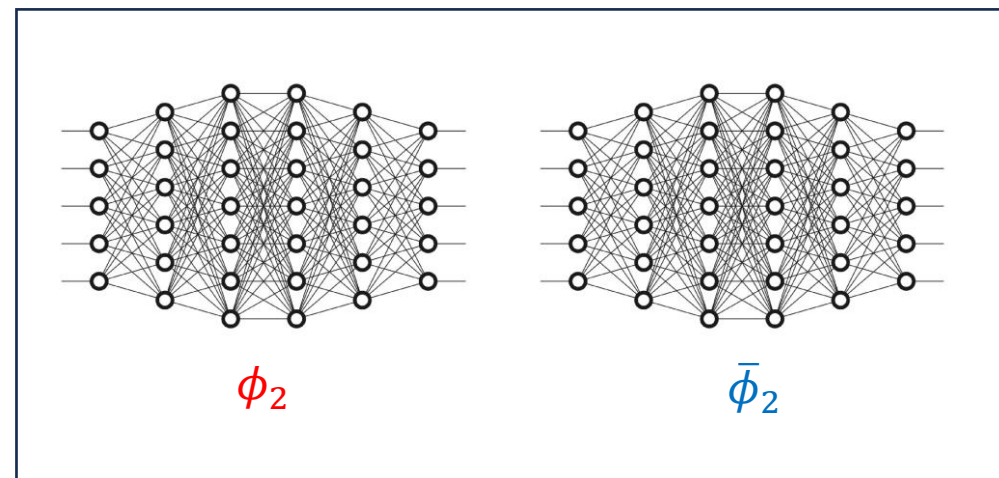
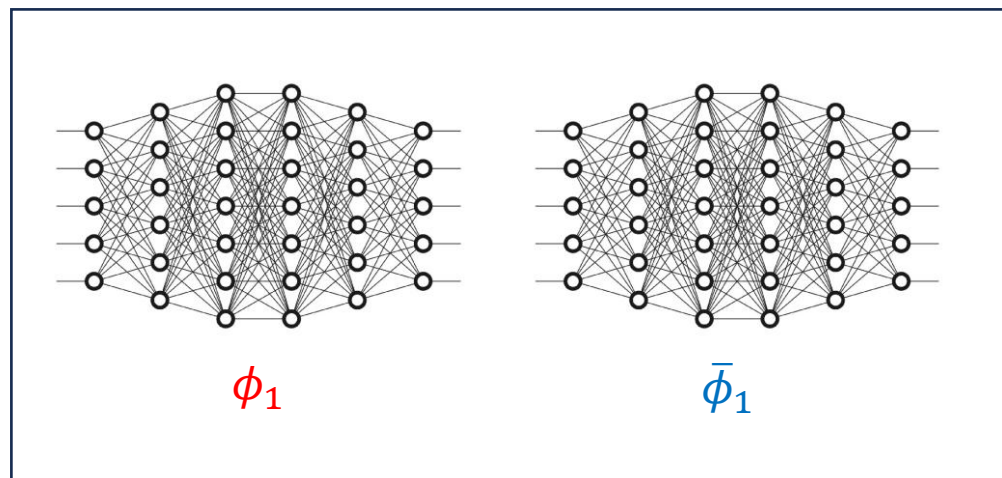
The bandit version of this algorithm: Page 11 [here](#)

Lillicrap et al., Continuous control with deep reinforcement learning. 2015.

TD3

Further Stabilizing DDPG (1/3): Twin Delayed DDPG

- Double Q-learning



Double Q-learning: When training ϕ_1 , instead of using $Q_{\bar{\phi}_1}$ to evaluate the regression target, use $Q_{\bar{\phi}_2}$

$$\text{TD3: } \min \{Q_{\bar{\phi}_1}, Q_{\bar{\phi}_2}\}$$

Double Q-learning: Use independent samples to train ϕ_1 and ϕ_2

TD3: Use the same set of samples
(the independence between ϕ_1 and ϕ_2 only comes from random initialization)

Further Stabilizing DDPG (2/3): Twin Delayed DDPG

- Target policy smoothing

DDPG: use $Q_{\bar{\phi}}(s', \mu_{\bar{\theta}}(s'))$ as the regression target

TD3: sample $a' = \mu_{\bar{\theta}}(s') + \mathcal{N}(0, \sigma^2)$
use $Q_{\bar{\phi}}(s', a')$ as the regression target

Handwritten diagram illustrating the TD3 target policy sampling process. The diagram is enclosed in a purple box and contains the expression $\mathbb{E}_{a \sim \bar{\mu}} Q_{\bar{\phi}}(s', a)$. An arrow points from the sampling notation $a \sim \bar{\mu}$ to the action variable a in the Q-function, indicating that the action is sampled from the target policy.

Further Stabilizing DDPG (3/3): Twin Delayed DDPG

- Delayed policy updates: running multiple steps of value updates before running one step of policy update

Remark: all three changes make it harder for the policy $\mu_\theta(s)$ to exploit the error of the Q function $Q_\phi(s, a)$

Twin Delayed DDPG (TD3)

$$\phi_1, \phi_2, \bar{\phi}_1, \bar{\phi}_2, \theta, \bar{\theta}$$

For $k = 1, 2, \dots$

Use $\mu_\theta(s) + \mathcal{N}(0, \sigma^2)$ to collect samples and place them in replay buffer

Sample a batch $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ from the replay buffer

For each sample i , draw $a'_i \sim \mu_{\bar{\theta}}(s'_i) + \mathcal{N}(0, \sigma^2 I)$

$$\phi_j \leftarrow \phi_j - \lambda \nabla_{\phi_j} \sum_{i=1}^n \left(Q_{\phi_j}(s_i, a_i) - r_i - \gamma \min_{\ell=1,2} Q_{\bar{\phi}_\ell}(s'_i, a'_i) \right)^2 \quad \forall j = 1, 2$$

If $k \bmod M = 0$:

$$\theta \leftarrow \theta + \eta \sum_{i=1}^n \nabla_{\theta} Q_{\phi_1}(s_i, \mu_{\theta}(s_i))$$

$$\bar{\theta} \leftarrow \tau \theta + (1 - \tau) \bar{\theta}$$

$$\bar{\phi}_j \leftarrow \tau \phi_j + (1 - \tau) \bar{\phi}_j \quad \forall j = 1, 2$$

$$a \sim \mu_{\bar{\theta}}(s)$$

$$\min \left\{ Q_{\phi_1}(s, a), Q_{\phi_2}(s, a) \right\}$$

SAC

Soft Actor-Critic (SAC)

- TD3 / DDPG: modeling $\mu_\theta(s)$ + additional noise for exploration
- SAC: modeling $\mu_\theta(s)$ and $\sigma_\theta(s)$ + adding entropy as an exploration bonus

Entropy Bonus (\approx Boltzmann Exploration)

Bandit

$$H(\pi) = -\sum_a \pi(a) \log \pi(a)$$

$$\pi = \operatorname{argmax}_{\pi} \sum_a \pi(a) R(a) + \alpha H(\pi) = \operatorname{argmax}_{\pi} \mathbb{E}_{a \sim \pi} [R(a) - \alpha \log \pi(a)]$$

MDP

$$\pi(a) \propto \exp\left(\frac{1}{\alpha} R(a)\right)$$

$$\begin{aligned} \pi &= \operatorname{argmax}_{\pi} \mathbb{E}^{\pi} \left[\sum_{h=0}^{\infty} \gamma^h \left(\sum_a \pi(a|s_h) R(s_h, a) + \alpha H(\pi(\cdot | s_h)) \right) \right] \\ &= \operatorname{argmax}_{\pi} \mathbb{E}^{\pi} \left[\sum_{h=0}^{\infty} \gamma^h (R(s_h, a_h) - \alpha \log \pi(a_h | s_h)) \right] \end{aligned}$$

Bellman Equation with Entropy Bonus

$$Q^{\pi}(s,a) \triangleq \mathbb{E}^{\pi} \left[\sum_{h=1}^{\infty} \gamma^{h-1} R(s_h, a_h) \mid (s_1, a_1) = (s, a) \right]$$

$$\Rightarrow Q^{\pi}(s,a) = R(s,a) + \mathbb{E}^{\pi} \left[\sum_{h=2}^{\infty} \gamma^{h-2} R(s_h, a_h) \right] = R(s,a) + \gamma \mathbb{E}_{s'} \sum_{a'} \pi(a'|s') Q^{\pi}(s', a')$$

$$\left\{ \begin{array}{l} Q^{\pi}(s,a) = R(s,a) + \underbrace{\alpha H(\pi(\cdot|s))}_{\text{Entropy Bonus}} + \gamma \mathbb{E}_{s'} \sum_{a'} \pi(a'|s') Q^{\pi}(s', a') \\ Q^{\pi}(s,a) = R(s,a) + \gamma \mathbb{E}_{s'} \left(\sum_{a'} \pi(a'|s') Q^{\pi}(s', a') + \underbrace{\alpha H(\pi(\cdot|s'))}_{\text{Entropy Bonus}} \right) \end{array} \right.$$

TD3 vs. SAC

- Value update

TD3: Sample $a' \sim \mu_\theta(s') + \mathcal{N}(0, \sigma^2)$

Use $Q_{\bar{\phi}}(s', a')$ as the regression target

SAC: Sample $a' \sim \pi_\theta(\cdot | s') = \mu_\theta(s') + \mathcal{N}(0, \sigma_\theta^2(s'))$

Use $Q_{\bar{\phi}}(s', a') - \alpha \log \pi_\theta(a' | s')$ as the regression target

Soft Actor-Critic (SAC)

For $k = 1, 2, \dots$

Use $\mu_\theta(s) + \mathcal{N}(0, \sigma_\theta^2(s))$ to collect samples and place them in replay buffer

Sample a batch $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ from the replay buffer

For each sample i , draw $a'_i \sim \mu_\theta(s'_i) + \mathcal{N}(0, \sigma_\theta^2(s'_i))$

$$\phi_j \leftarrow \phi_j - \lambda \nabla_{\phi_j} \sum_{i=1}^n \left(Q_{\phi_j}(s_i, a_i) - r_i - \gamma \left(\min_{\ell=1,2} Q_{\bar{\phi}_\ell}(s'_i, a'_i) + \alpha \log \pi_\theta(a'_i | s'_i) \right) \right)^2 \quad \forall j = 1, 2$$

Perform Policy (θ) Update (to be specified later)

$$\bar{\phi}_j \leftarrow \tau \phi_j + (1 - \tau) \bar{\phi}_j \quad \forall j = 1, 2$$

TD3 vs. SAC

- Policy update

TD3: Do not view $-\alpha \log \pi_{\theta}(a|s)$ as part of the reward
Only train $\mu_{\theta}(s)$

$$\theta \leftarrow \theta + \eta \nabla_{\theta} Q_{\phi}(s, \mu_{\theta}(s))$$

SAC: View $-\alpha \log \pi_{\theta}(a|s)$ as part of the reward
Train both $\mu_{\theta}(s)$ and $\log \sigma_{\theta}(s)$

Sample $a_{\theta}(s) = \mu_{\theta}(s) + \epsilon \sigma_{\theta}(s)$ where $\epsilon \sim \mathcal{N}(0,1)$

$$\theta \leftarrow \theta + \eta \nabla_{\theta} (Q_{\phi}(s, a_{\theta}(s)) - \alpha \log \pi_{\theta}(a_{\theta}(s)|s))$$

Soft Actor-Critic (SAC)

$$\text{Further using } \pi_{\theta}(a|s) = \frac{1}{(2\pi\sigma_{\theta}(s)^2)^{d/2}} \exp\left(-\frac{\|a-\mu_{\theta}(s)\|^2}{\sigma_{\theta}(s)^2}\right)$$

For $k = 1, 2, \dots$

Use $\mu_{\theta}(s) + \mathcal{N}(0, \sigma_{\theta}^2(s))$ to collect samples and place them in replay buffer

Sample a batch $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ from the replay buffer

For each sample i , draw $a'_i \sim \mu_{\theta}(s'_i) + \mathcal{N}(0, \sigma_{\theta}^2(s'_i))$

$$\phi_j \leftarrow \phi_j - \lambda \nabla_{\phi_j} \sum_{i=1}^n \left(Q_{\phi_j}(s_i, a_i) - r_i - \gamma \left(\min_{\ell=1,2} Q_{\bar{\phi}_{\ell}}(s'_i, a'_i) + \alpha \log \pi_{\theta}(a'_i | s'_i) \right) \right)^2 \quad \forall j = 1, 2$$

Let $a_{\theta}(s_i) = \mu_{\theta}(s_i) + \epsilon \sigma_{\theta}(s_i)$ where $\epsilon \sim \mathcal{N}(0, I)$

$$\theta \leftarrow \theta + \eta \sum_{i=1}^n \nabla_{\theta} (Q_{\phi_1}(s, a_{\theta}(s_i)) - \alpha \log \pi_{\theta}(a_{\theta}(s_i) | s_i))$$

$$\bar{\phi}_j \leftarrow \tau \phi_j + (1 - \tau) \bar{\phi}_j \quad \forall j = 1, 2$$