

* Pseudo code

① REINFORCE with Baseline (Episodic)

Input: a differentiable policy parameterization $\pi(a|s; \theta)$
 Input: a differentiable state-value parameterization $\hat{V}(s; w)$
 Parameters: step size $\alpha > 0$, $\beta > 0$
 Initialize policy parameter $\theta \in \mathbb{R}^k$ and state-value weights $w \in \mathbb{R}^d$
 Repeat forever:
 Generate an episode $S_0, A_0, R_0, \dots, S_{T-1}, A_{T-1}, R_{T-1}$, following $\pi(\cdot|s; \theta)$
 For each step of the episode $t=0, \dots, T-1$:
 $G_t \leftarrow$ return from step t
 $\hat{V}_t \leftarrow \hat{V}(S_t; w)$
 $w \leftarrow w + \beta \gamma \delta_t \hat{V}(S_t; w)$
 $\theta \leftarrow \theta + \alpha \gamma \sum_{s \sim \pi} \delta_t \nabla_{\theta} \pi(a_t|s; \theta)$
 $\delta_t = \frac{G_t - \hat{V}_t}{\hat{V}(S_t; w)}$

② REINFORCE, MCPG (Episodic)

Input: a differentiable policy parameterization $\pi(a|s; \theta)$, $\forall a \in A, s \in S, \theta \in \mathbb{R}^k$
 Initialize policy weights θ
 Repeat forever:
 Generate an episode $S_0, A_0, R_0, \dots, S_{T-1}, A_{T-1}, R_{T-1}$, following $\pi(\cdot|s; \theta)$
 For each step of the episode $t=0, \dots, T-1$:
 $G_t \leftarrow$ return from step t
 $\theta \leftarrow \theta + \alpha \gamma \sum_{s \sim \pi} \delta_t \nabla_{\theta} \pi(a_t|s; \theta)$

③ DQN (Deep Q-learning) with Experience Replay

Initialize replay memory D to capacity N
 Initialize action-value function Q with random weights
 for episode $e = 1, M$ do
 Initialize sequence $s_1 = s_0$ and preprocessed sequences $\phi_1 = \phi(s_1)$
 for $t = 1, T$ do
 With probability ϵ select a random action a_t
 otherwise select $a_t = \arg \max_a Q(s_t, a; \theta)$
 Execute action a_t in simulator and observe reward r_t and next state
 Set $s_{t+1} = s_t, a_t, r_t$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in D
 Sample random minibatch of transitions $(\phi_i, a_i, r_i, \phi_{i+1})$ from D
 Set $\delta_t = \begin{cases} r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta) - Q(s_t, a_t; \theta) & \text{for terminal } \phi_{t+1} \\ \gamma Q(s_{t+1}, a_{t+1}; \theta) - Q(s_t, a_t; \theta) & \text{for non-terminal } \phi_{t+1} \end{cases}$
 Perform a gradient descent step on $(\delta_t - Q(s_t, a_t; \theta))^2$
 end for
 end for

④ One-Step Actor-Critic (Episodic)

Input: a differentiable policy parameterization $\pi(a|s; \theta)$, $\forall a \in A, s \in S, \theta \in \mathbb{R}^k$
 Input: a differentiable state-value parameterization $\hat{V}(s; w)$, $\forall s \in S, w \in \mathbb{R}^d$
 Parameters: step size $\alpha > 0$, $\beta > 0$
 Initialize policy weights θ and state-value weights w
 Repeat forever:
 Initialize S (first state of episode)
 $t \leftarrow 1$
 while S is not terminal:
 $A \sim \pi(\cdot|S; \theta)$
 Take action A , observe S', R
 $\hat{V}_t \leftarrow \hat{V}(S; w)$
 $\hat{V}_{t+1} \leftarrow \hat{V}(S'; w)$ (if S' is terminal, then $\hat{V}_{t+1} = 0$)
 $w \leftarrow w + \beta \gamma \delta_t \hat{V}(S; w)$
 $\theta \leftarrow \theta + \alpha \gamma \delta_t \nabla_{\theta} \pi(a_t|s; \theta)$
 $t \leftarrow t + 1$
 $S \leftarrow S'$

⑤ A2C (Advantage Actor-Critic)

Assume parameter vectors θ and w
 Initialize step counter $t \leftarrow 1$
 Initialize episode counter $E \leftarrow 1$
 repeat:
 Reset gradients: $d\theta \leftarrow 0$ and $dw \leftarrow 0$
 $t_{start} = t$
 Get state s_t
 repeat:
 Perform a_t according to policy $\pi(a_t|s_t; \theta)$
 Receive reward r_t and new state s_{t+1}
 $t \leftarrow t + 1$
 until terminal s_t or $t - t_{start} == t_{max}$
 $R = \begin{cases} r_t & \text{for terminal } s_t \\ \hat{V}(s_t; w) & \text{for non-terminal } s_t \end{cases}$ // bootstrap from last state
 for $\tau \in \{t-1, \dots, t_{start}\}$ do
 $R_t \leftarrow R + \gamma R_{t+1}$
 Accumulate gradients w.r.t θ : $d\theta \leftarrow d\theta + \gamma \delta_t \nabla_{\theta} \pi(a_t|s_t; \theta) (R - \hat{V}(s_t; w))$
 Accumulate gradients w.r.t w : $dw \leftarrow dw + \gamma \delta_t \nabla_w \hat{V}(s_t; w) (R - \hat{V}(s_t; w))$
 end for
 perform update θ using $d\theta$ and w using dw
 $t \leftarrow t + 1$
 until $E > E_{max}$

⑥ A2C (Asynchronous Advantage Actor-Critic)

Assume global shared parameter vectors θ and w and global shared counter $T \leftarrow 0$
 Assume thread-specific parameter vectors θ' and w'
 Initialize thread step counter $t \leftarrow 1$
 repeat:
 Reset gradients: $d\theta \leftarrow 0$ and $dw \leftarrow 0$
 Synchronize thread-specific parameters $\theta' \leftarrow \theta$ and $w' \leftarrow w$
 $t_{start} = t$
 Get state s_t
 repeat:
 Perform a_t according to policy $\pi(a_t|s_t; \theta')$
 Receive reward r_t and new state s_{t+1}
 $t \leftarrow t + 1$
 until terminal s_t or $t - t_{start} == t_{max}$
 $R = \begin{cases} r_t & \text{for terminal } s_t \\ \hat{V}(s_t; w') & \text{for non-terminal } s_t \end{cases}$ // bootstrap last state
 for $\tau \in \{t-1, \dots, t_{start}\}$ do
 $R_t \leftarrow R + \gamma R_{t+1}$
 Accumulate gradients w.r.t θ' : $d\theta' \leftarrow d\theta' + \gamma \delta_t \nabla_{\theta'} \pi(a_t|s_t; \theta') (R - \hat{V}(s_t; w'))$
 Accumulate gradients w.r.t w' : $dw' \leftarrow dw' + \gamma \delta_t \nabla_{w'} \hat{V}(s_t; w') (R - \hat{V}(s_t; w'))$
 end for
 Perform asynchronous update $\theta \leftarrow \theta + \alpha d\theta$ and $w \leftarrow w + \alpha dw$ using $d\theta$
 until $T > T_{max}$

Policy Gradient

value based

off-policy Policy Gradient

① SAC (Soft Actor-Critic)

Input: The learning rates, λ, λ_a and λ_v for functions π, Q_a and V respectively;
 Input: The weighting factor τ for exponential moving average
 Initialize parameters θ, w, ψ and $\bar{\psi}$
 for each iteration do
 for each environment setup do
 $a_t \sim \pi(a_t|s_t; \theta)$
 $S_{t+1} \sim \mathcal{P}(S_{t+1}|S_t, a_t)$
 $D \leftarrow D \cup \{(s_t, a_t, r(s_t, a_t), S_{t+1})\}$
 for each gradient update step do
 $\psi \leftarrow \psi - \lambda_v \nabla_{\psi} J_{\psi}(D)$
 $w \leftarrow w - \lambda_a \nabla_w J_{Q_a}(D)$
 $\theta \leftarrow \theta - \lambda \nabla_{\theta} J_{\pi}(D)$
 $\bar{\psi} \leftarrow \tau(1-\tau)\bar{\psi} + \tau\psi$

② TRPO (Trust Region Policy Optimization)

Initialize π_0 to anything
 Repeat forever:
 Sample s_t and set $\pi \leftarrow \pi_0$
 Repeat N times:
 Sample $a_t \sim \pi(a_t|s_t)$
 Execute a_t , observe S_{t+1}, r_t
 $\hat{V}_t \leftarrow \tau + \gamma \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t] - \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t]$
 $A(s_t) \leftarrow r_t + \gamma \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t] - \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t]$
 Update Q : $w \leftarrow w + \alpha \nabla_w \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t]$
 $\pi \leftarrow \pi + \alpha \nabla_{\pi} \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t]$
 Update π : $\theta \leftarrow \arg \max_{\theta} \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t]$
 Subject to $\sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t] \leq \delta$

③ PPO (Proximal Policy Optimization)

Input: initial policy parameters θ_0 , initial function parameters ϕ_0
 for $k=0, 1, 2, \dots$ do
 Collect set of trajectories $D_k = \{ \tau_i \}$ by running policy $\pi_k = \pi(\cdot|s; \theta_k)$ in the environment
 Compute returns to go G_t
 Compute advantage estimates, A_t (using any method of advantage estimation) based on current value function V_{ϕ_k}
 Update the policy by maximizing PPO objective:
 $\theta_{k+1} = \arg \max_{\theta} \frac{1}{|D|} \sum_{\tau \in D} \sum_{t \in \tau} \min \left(\frac{\pi(a_t|s_t; \theta) A_t}{\pi(a_t|s_t; \theta_k)}, \frac{\pi(a_t|s_t; \theta)}{\pi(a_t|s_t; \theta_k)} \right) \sum_{k=t}^{\infty} \gamma^k \mathbb{E}_{\pi} [r_{t+k} | S_t, a_t]$
 Fit value function by MSE: $\phi_{k+1} = \arg \min_{\phi} \frac{1}{|D|} \sum_{\tau \in D} \sum_{t \in \tau} (V_{\phi}(s_t) - G_t)^2$
 end for

* MCTS (Monte Carlo Tree Search)

- AlphaGo의 게임 탐색 알고리즘
- 많은 tree made 대응으로 빠른 게임 시뮬레이션은 물론 가장 가치있는 플레이 방법도 탐색
- Min-Max 알고리즘의 성능을 개선하는 것은 평균값 탐색하기 위한 탐색의 최적

* AlphaGo

- MCTS를 deep learning pipeline을 통해 빠른 학습을 가능하게 함
- Policy network (가정), Value network (가치)

* AlphaZero

- AlphaGo는 비강제 학습이었고, AlphaZero는 강제 학습을 함

* A2C vs DQN

- ① A2C Algorithm has 2 networks, a target network and a Q network (CT)
- ② A Double DQN algorithm has 2 networks, a target network and a Q network (CF)
- multiple copies of each network is needed
- ③ A DQN algorithm utilizes replay buffer (CT)
- ④ The only difference between DQN and A2C is that DQN utilizes value function approximation with no additional enhancements (ex. replay buffer) (CF)
- DQN has 2 networks
- ⑤ A Double DQN algorithm utilizes replay buffer (CT)

Kulterm

- #1 Multi-armed bandits are simply more efficient method for Reinforcement Learning that can be used for small domains. T/F
 → MABs have no state so are inherently more limited than full RL. The domains for MABs are simpler but not just because they are small in state or action space.
- #2 Anytime an agent takes an action $b \neq \max_a Q(s,a)$, it could be considered exploration. T/F
 → one example is ϵ -greedy but it could be any other function of the value function.
- #3 Exploration actions only result from the use of an ϵ -greedy policy. T/F
 → we need to get the max of Q value for all states.
- #4 To get the value function from a Q -function, we need to get the average Q value over all actions at a given state. T/F
 → we need to get the max of Q value for all states.
- #5 An example of an optimistic action selection strategy would be to take an action which has never been tried for a given state s . T/F
 → because it assumes unvisited pairs (s,a) are good.
- #6 The "Target" in a TD method is provided as input to the agent, it is the pre-defined goal you seeking to find. T/F
 → The target is the new immediate reward combined with the estimate of return going forward with the current policy.
- #7 The expected SARSA algorithm attempts to use information from all actions by taking the mean over all actions of the value function for the current state. T/F
 → The algorithm takes expected value of all actions, this would only be equal to the mean if all actions had the same probability. The value of each (s,a) pair is weighted by the probability that the action a will be selected by the current policy while s .
- #8 Which of the following best describe the meaning of the values stored in a state-value function or state-action value function?
 (a) an estimate of expected reward (b) an estimate of discounted rewards (c) a function of rewards resulting from following the current policy into the future (d) all of them
- #9 A Monte Carlo RL method is one that
 (a) samples randomly from trajectory of (s,a,r) tuples and uses the average weighted return from those samples.
 (b) uses the complete return from an entire trajectory which is sampled randomly from the current policy
 (c) uses any kind of random sampling policy in its policy.
- #10 Which one of the following parameters was introduced to help the Credit Assignment Problem?
 (a) λ in TD (λ) (b) r in DP/MC/TD methods (c) ϵ in greedy policy definition (d) α in TD methods
- #11 Which one of the following parameters was introduced to encourage more exploration during RL?
 (a) λ in TD (λ) (b) r in DP/MC/TD methods (c) ϵ in greedy policy definition (d) α in TD methods
- #12 Which one of the following parameters was introduced to control the influence new information has on the value function update?
 (a) λ in TD (λ) (b) r in DP/MC/TD methods (c) ϵ in greedy policy definition (d) α in TD methods
- #13 Which one of the following parameters was introduced to control the importance of receiving reward sooner in episode, rather than later?
 (a) λ in TD (λ) (b) r in DP/MC/TD methods (c) ϵ in greedy policy definition (d) α in TD methods
- #14 The main difference between solving an MDP with Dynamic programming and using Reinforcement Learning is:
 (a) In RL we don't need to find an exact solution
 (b) In RL, the dynamics are not available
 (c) when we solve an MDP, we don't discount future rewards
 (d) In RL, we use the policy improvement theorem to greedily select actions based on our value estimate.
- #15 Which of the following calculations could give different answers depending on the starting conditions of the value function:
 (a) Exact solution using the value iteration Algorithm
 (b) Exact solution using the policy iteration Algorithm
 (c) Converged solution of SARSA (tabular representation)
 (d) Converged solution of SARSA (using linear value function approximation)
 (e) Converged solution of Q -Learning (tabular representation)

#16 Let the general form of the return for an MDP from timestep t going forward until the terminal state is reached at $t=T$ be defined recursively for all non-terminal steps as:

$G_t = R_{t+1} + \gamma V_t$

Suppose we are given a sequence of rewards actually received from a policy on an MDP as follows: $R_1=2, R_2=-1, R_3=5, R_4=1, R_5=1$

The MDP has horizon of 5 and we use a discount factor of 0.9. Once the terminal state is reached, there are no rewards received and the episode ends.

① Write one-step, recursive calculation for each return below in symbols and show your calculation.

$$G_0 = R_1 + \gamma G_1 = 2 + 0.9 \times 3.581 = 5.2229$$

$$G_1 = R_2 + \gamma G_2 = -1 + 0.9 \times 5.09 = 3.581$$

$$G_2 = R_3 + \gamma G_3 = 5 + 0.9 \times 0.1 = 5.09$$

$$G_3 = R_4 + \gamma G_4 = 1 + 0.9 \times 0 = 1$$

$$G_4 = R_5 + \gamma G_5 = 1 + 0.9 \times 0 = 1 \quad (\because G_5 = \text{terminal state} \therefore 0)$$

② Write out the full, non-recursive calculation of G_0 in symbols and numbers for the whole domain from the start state.

$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \gamma^4 R_5$$

$$= 2 + (0.9 \times -1) + (0.9^2 \times 5) + (0.9^3 \times 1) + (0.9^4 \times 1) = 5.2229$$

#17 SARSA Value Function Update

Consider a system with 3 states and 2 actions. You will perform update calculations to a tabular state-action value function on an MDP using a learning rate $\alpha=0.15$, and a discount factor $\gamma=0.9$ for each step. For each step, the current state, action, reward, and next state and action are given as $(s_t, a_t, r_t, s_{t+1}, a_{t+1})$. In this MDP, rewards are obtained when the agent enters the state after taking the action.

The pre-existing state-action value function, $Q(s,a)$ is defined by the following table:

$Q(s,a)$	a_1	a_2
s_1	5	-1
s_2	-2	2
s_3	-1	3

Show the calculations for the following updates using the provided trajectory of experiences using the SARSA algorithm:

① General Formula for updating $Q(s,a)$ using $s_t, a_t, r_t, s_{t+1}, a_{t+1}$:

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

② $(s_{11}, a_{11}, 2, s_{11}, a_{11})$, $Q(s_{11}, a_{11})$

$$Q(s_{11}, a_{11}) = Q(s_{11}, a_{11}) + \alpha [r_t + \gamma Q(s_{11}, a_{11}) - Q(s_{11}, a_{11})] = 5 + 0.15 \times (2 + 0.9 \times 2 - 5) = 4.82$$

③ $(s_{21}, a_{21}, -1, s_{11}, a_{11})$, $Q(s_{21}, a_{21})$

$$Q(s_{21}, a_{21}) = Q(s_{21}, a_{21}) + \alpha [r_t + \gamma Q(s_{11}, a_{11}) - Q(s_{21}, a_{21})] = -2 + 0.15 \times (1 + 0.9 \times 4.82 - (-2)) = 2.2$$

④ $(s_{11}, a_{11}, 2, s_{21}, a_{21})$, $Q(s_{11}, a_{11})$

$$Q(s_{11}, a_{11}) = Q(s_{11}, a_{11}) + \alpha [r_t + \gamma Q(s_{21}, a_{21}) - Q(s_{11}, a_{11})] = 4.82 + 0.15 \times (2 + 0.9 \times 2.2 - 4.82) = 4.13$$

⑤ $(s_{21}, a_{21}, -1, s_{31}, a_{31})$, $Q(s_{21}, a_{21})$

$$Q(s_{21}, a_{21}) = Q(s_{21}, a_{21}) + \alpha [r_t + \gamma Q(s_{31}, a_{31}) - Q(s_{21}, a_{21})] = 2.2 + 0.15 \times (-1 + 0.9 \times 1 - 2.2) = -1.69$$

#18 Expected SARSA

Now, in #17, we were instead using Expected SARSA and we know the policy is ϵ -greedy but biased towards a_1 , or more precisely:

- For any state s , if $a_1 = \underset{a}{\operatorname{argmax}} Q(s,a)$ then $\pi(a_1|s) = 0.9$

- For any state s , if $a_2 = \underset{a}{\operatorname{argmax}} Q(s,a)$ then $\pi(a_2|s) = 0.6$

Using the same MDP before, show the formula and just the first value update calculation for Expected SARSA:

① General formula for updating $Q(s,a)$ using s, a, r, s', a' :

$$Q(s,a) = Q(s,a) + \alpha [r_k + \gamma \sum_{a'} \pi(a'|s') Q(s',a') - Q(s,a)]$$

② $Q(s_1, a_1, s_1, s_2, a_2), Q(s_1, a_1)$

$$\begin{aligned} Q(s_1, a_1) &= Q(s_1, a_1) + \alpha [r + \gamma [(1 - \pi(a_2|s_1)) Q(s_2, a_1) + \pi(a_2|s_2) Q(s_2, a_2)] - Q(s_1, a_1)] \\ &= 5 + 0.15 (2 + 0.9 \times (0.4 \times 2 + 0.6 \times 2) - 5) = 4.6 \end{aligned}$$