

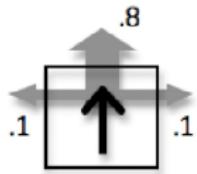
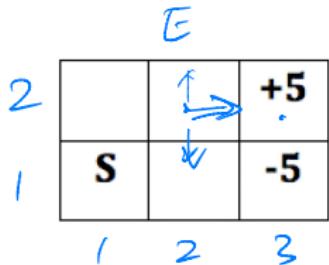
PART 5

Reinforcement Learning

MDP

- * State space: S
- * Action space: A
- * Transition prob: $\underline{p(s'|s,a)}$
Action Model
- * Reward function: $r(s,a,s')$

$$\underline{r(s,a)} = \sum_{s'} p(s'|s,a) \underline{r(s,a,s')}$$



s	a	s'	p(s' s,a)	r(s,a,s')
(2,2)	<i>E</i>	{ (3,2) (2,2) (2,1) }	0.8 0.1 0.1	+5 0 0

$$\gamma((2,2), E) = 0.8 \times 5 + 0.1 \times 0 + 0.1 \times 0 = 0.4$$

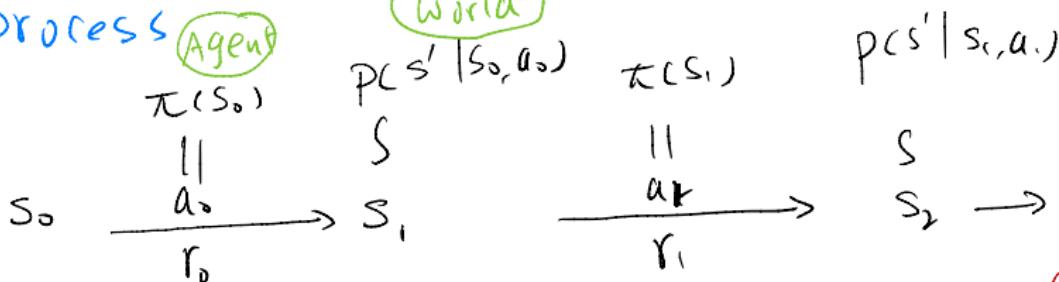
One-step expected reward

Policy: $\pi: S \rightarrow A$

* process

Agent

World



Discount factor

$$0 < \gamma < 1$$

- * Math reason
- * Intuitive

Discounted total rewards

$$R^\pi(s_0) = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$$

(state)
Value

$$V^\pi(s) = E[R^\pi(s)]$$

random Yes
No

Function: Total rewards for following π from s

Theorem: Exist π^* s.t. $\forall \pi$

$$V^{\pi^*}(s) \geq V^\pi(s) \quad \forall s$$

optimal policy: π^*

Optimal Value: $V^*(s) = V^{\pi^*}(s)$
function

Task 1: Planning

$$p(s'|s, a), r(s, a) \Rightarrow \pi^*$$

Task 2: Reinforcement Learning

$$\{ (s, a, r, s') \} \xrightarrow{\text{Experience tuple}} \pi^*$$

Task 1 : $p(s'|s,a), r(s,a) \Rightarrow \pi^*$

(1) Compute V^*

(2) Obtain π^* from V^*

(2) Optimal action-state value function

In s , take a , then follow π^*

$$\underline{Q^*(s,a)} = \frac{\text{current reward} + \gamma \sum_{s'} p(s'|s,a) V^*(s')}{\text{future reward}}$$

$$\pi^*(s) = \arg \max_a \underline{Q^*(s,a)}$$

Task 1, Step 1: Compute V^*

Value Iteration

* pick V_0

Bellman Operator

* For $k=0 \rightarrow$ termination

$$V_{k+1}(s) = \max_a \left\{ \underbrace{r(s,a) + \gamma \sum_{s'} p(s'|s,a) V_k(s')}_{\text{current Future}} \right\}$$

improve

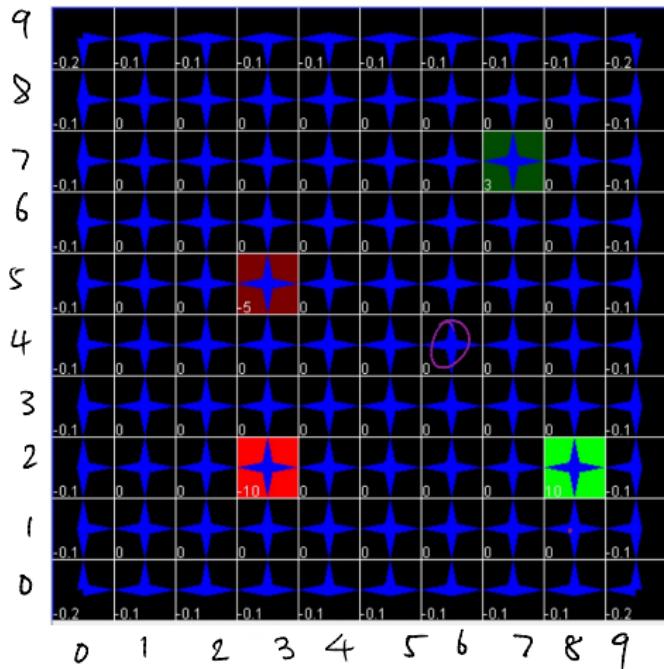
Theorem: $\lim_{k \rightarrow \infty} V_k(s) = V^*$

$$V_{k+1}(s) = \max_a \left\{ r(s, a) + \gamma \underbrace{\sum_{s'} p(s'|s, a) V_k(s')}_{\text{Bellman Optimality Operator}} \right\}$$

$\Downarrow k \rightarrow \infty$

$$V^*(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^*(s') \right\}$$

Bellman's Optimality Equation

V_1 

$S = (8, 1)$

$V_1(S) = 0$

	Current	Future	$\max a$
a	up	0	large +
	down	0	small -
	left	0	small
	right	0	small

 $V_2(S) \uparrow$

$S = (3, 1)$

$V_1(S) = 0$

	Current	Future	$\max a$
a	up	large -	X
	down	small -	
	left	small -	
	right	small -	

 $V_2(S) \downarrow$

$S = (6, 4)$

 \uparrow

$V_1(S) = 0$

$V_2(S) = 0$

 $\Rightarrow ?$

Value iteration for $\underline{Q^*(s,a)}$

$\underline{V^*(s)}$

* PICK Q_0

$Q(s,a)$ Table

	a_1	a_2	-
s_1	0.5	2	-
s_2	3.2	0.1	-

* for $k=0 \rightarrow$ termination

$$Q_{k+1}(s,a) = r(s,a) + \gamma \sum_s p(s'|s,a) \max_{a'} Q_k(s', \frac{a'}{?})$$

— — ↗

Current Future improve

Theorem : $\lim_{k \rightarrow \infty} Q_k(s,a) = Q^*(s,a)$

Greedy policy

$\pi_k(s) = \arg \max_a Q_k(s,a)$ based on Q_k

Task 2: $\{ (s, a, r, s') \} \Rightarrow Q^*(s, a)$
 Experience Tuples

* Initialize $Q(s, a)$ Q-table

* Collect (s, a, r, s')
 Feed back

$$Q(s, a) \leftarrow Q(s, a) + \alpha \Delta$$

↑ learning rate

$$\Delta = \underbrace{r(s, a) + \gamma \max_{a'} Q(s', a')}_{\text{TD target}} - \frac{Q(s, a)}{\text{current estimate}}$$

Q-learning: Converges to Q^* if each (s, a) visited infinitely often

$$Q_{k+1}(s, a) = r(s, a) + \gamma \sum_{s'} p_{cs'}(s, a) \max_{a'} Q_k(s', a')$$

$$\approx r(s, a) + \gamma \frac{1}{m} \sum_{j=1}^m \max_{a'} Q_k(s'_j, a') \quad s'_1, \dots, s'_m \sim p_{cs'}(s, a)$$

$$\Rightarrow r(s, a) + \gamma \max_{a'} Q_k(s'_1, a') \quad \text{when } m=1$$

* high Variance

$$(s, a, \underline{r}, \underline{s'_1}) \quad * \text{ unbiased}$$

Improve $Q(s, a)$ using (s, a, r, s')

$$r + \gamma \max_{a'} Q(s', a')$$

TD target (TD: Time difference)

Q-learning is off-policy

- Choose a for the state s (ϵ -greedy with $\arg \max_a Q(s, a)$)
- Take action a , observe r and s'
- Update:

$$\begin{aligned} Q(s, a) &\leftarrow \underline{Q(s, a)} + \alpha [r + \gamma \max_{a'} Q(s', a') - \underline{Q(s, a)}] \\ s &\leftarrow s' \end{aligned}$$

Q_{Next}

Action for the next time step ?

* For estimation: $\arg \max_a Q(s, a')$

* Actually execute: $\arg \max_a Q_{\text{Next}}(s, a')$

L13 Value - Based Deep RL

Deep Q-Net Work (DQN)

$$S \rightarrow [\theta] = \begin{cases} Q(s, a_1) \\ Q(s, a_2) \end{cases}$$

* Initialize θ

* Repeat: Get (s, a, r, s')

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} (y_{\theta} - Q(s, a; \theta))^2$$

* Target Net Work

* Replay Buffer

$$y_{\theta} = r + \gamma \max_{a'} Q(s', a'; \theta)$$

TD target

* Prioritized Replay

Moving Target

θ_0

↓ Target

$$\theta_1 \leftarrow \theta_0 - \alpha \nabla_{\theta} (\underline{y}_{\theta_0} - Q(s_1, a_1; \theta))^2$$

$$\theta_2 \leftarrow \theta_1 - \alpha \nabla_{\theta} (\underline{y}_{\theta_1} - Q(s_2, a_2; \theta))^2$$

... -

Solution: keep a slow-changing

target network

D / NO

Review : P11, P13

Update rule for DQN:

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} ([\underbrace{r(s, a)}_{\text{current}} + \gamma \max_{a'} \underbrace{Q(s', a'; \theta^-)}_{\text{Future estimate}}] - Q(s, a; \theta))^2$$

TD target

① Pick next action: $a^* = \arg \max_a Q(s', a; \theta^-)$

② Evaluate a^* : $Q(s', a^*; \theta^-)$

\Rightarrow over-estimate of future reward

Update rule for Double DQN:

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} ([r(s, a) + \gamma Q(s', a^*; \theta^-)] - Q(s, a; \theta))^2$$

where $a^* = \arg \max_{a'} Q(s', a'; \theta)$.

Analogy: Estimate highest score of a class
in final exam

- ① Teacher θ^- :
 - (a) Pick best student a_1^*
 - (b) Estimate score for a_1^*
 - ②
 - (a) Teacher θ pick best student a_2^*
 - (b) Teacher θ^- estimate score for a_2^*
- Score 1 > Score 2



Which one
trust more?

DQN: $s \rightarrow \boxed{\theta} = \frac{Q(s, a_1)}{Q(s, a_2)}$ Discrete Action Space

what about continuous action space?

DDPG

L14 : Policy-Based Deep RL

Value-Based Deep RL:

$$s \rightarrow [\theta] \rightarrow q$$

$$\{ (s, a, r, s') \} \Rightarrow Q(s, a; \theta)$$

$$\Rightarrow \pi(s) = \arg \max_a Q(s, a; \theta)$$

Deterministic

policy-based Deep RL

stochastic

$$\{ (s, a, r, s') \} \Rightarrow \pi_\theta(a|s)$$

$$s \rightarrow [\theta] \rightarrow \pi_\theta(a|s)$$

↓
loss

A	T
N	0.8
S	0.1
E	0.05
W	0.05



$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \int \pi_{\theta}(\tau) r(\tau) d\tau && \text{log gradient trick} \\
 &= \int r(\tau) \underbrace{\nabla_{\theta} \pi_{\theta}(\tau)}_{\nabla_{\theta} \log \pi_{\theta}(\tau)} d\tau \\
 &= \int r(\tau) \pi_{\theta}(\tau) \underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau)}_{\nabla_{\theta} \pi_{\theta}(\tau)} d\tau \\
 &= \int \pi_{\theta}(\tau) \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r(\tau) d\tau && \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)}
 \end{aligned}$$

$$T = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_t, a_t, r_t, \dots)$$

$$\pi_{\theta}(\tau) = p(s_0) \prod_t \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

$$\begin{aligned}
 \nabla_{\theta} \log \pi_{\theta}(\tau) &= \nabla_{\theta} \log p(s_0) + \nabla_{\theta} \sum_t \log \pi_{\theta}(a_t | s_t) \\
 &= 0 + \nabla_{\theta} \underbrace{\sum_t \log p(s_{t+1} | s_t, a_t)}_{=0}
 \end{aligned}$$

$$\begin{aligned}
 \mathbb{V}_\theta J &= \int \pi_\theta(\tau) \sum_t P_\theta \log \pi_\theta(a_t | s_t) r(\tau) d\tau \\
 &\approx \frac{1}{N} \sum_{i=1}^N \sum_t P_\theta \log \pi_\theta(a_t^i | s_t^i) r(\tau^i)
 \end{aligned}$$

$$\tau^1, \dots, \tau^N \sim \pi_\theta(\tau)$$

$$\tau^i = (s_0^i, a_0^i, r_0^i, \dots, s_t^i, a_t^i, r_t^i, \dots)$$

Supervised Learning: $\{x_i, y_i\}_{i=1}^N \Rightarrow p(Y|x; \theta)$

$$\theta \leftarrow \theta + \alpha \frac{1}{N} \sum_{i=1}^N D_\theta \log p(Y_i|x_i; \theta)$$

Maximize prob of data (y_i)

Current: $\{\underbrace{(s_t^i, a_t^i)}_{\bar{x}}, \underbrace{y_i}_{\bar{Y}}\} | i=1, \dots, N, t=1, \dots, T\}$

$$\theta \leftarrow \theta + \alpha \frac{1}{N} \sum_t \sum_i D_\theta \log p(a_t^i | s_t^i; \theta)$$

maximize prob of the a_t^i 's

Imitation Learning: Good if trajectories
are expert demonstration

Policy Gradient (REINFORCE)

$$\theta \leftarrow \theta + \alpha \frac{1}{N} \sum_i \sum_t V_\theta \log p(a_t^i | s_t^i; \theta) r(\tau^i)$$

$r(\tau^i) > 0$: positive experience

Increase prob of actions

$r(\tau^i) < 0$: Negative experience

Decrease prob of action

Actor-Critic Algo

$V^*(s)$: optimal total reward starting from s

$Q^*(s, a)$: total reward for taking a in s
and then act optimally

$$V^*(s) = \max_a Q^*(s, a)$$

π : stochastic $\pi(a|s)$

$V^\pi(s)$: total reward for following π
Starting from s

$Q^\pi(s,a)$: total reward for taking a in s
and then following π

$$V^\pi(s) = \sum_a Q^\pi(s,a) \pi(a|s) : \begin{array}{l} \text{Average} \\ \text{of } Q \\ \text{over actions} \end{array}$$

$$A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$$

Advantage function of a in s

Policy Gradient (REINFORCE)

$$\theta \leftarrow \theta + \alpha \frac{1}{N} \sum_t \sum_i \nabla_{\theta} \log p(a_t^i | s_t^i; \theta) r(\tau^i)$$

Actor-Critic (AC) Algo

$$\theta \leftarrow \theta + \alpha \frac{1}{N} \sum_t \sum_i \nabla_{\theta} \log p(a_t^i | s_t^i; \theta) A^{\pi}(s_t^i, a_t^i)$$

$A^{\pi}(s_t^i, a_t^i) > 0$: a_t^i better than average
increase its prob

$A^{\pi}(s_t^i, a_t^i) < 0$: a_t^i worse than average
decrease its prob

How to estimate:

$$A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t)$$

$$= r_t + \gamma V^\pi(\underline{s_{t+1}}) - \underline{V^\pi(s_t)}$$

$(s_t, a_t, r_t, \underline{s_{t+1}})$

$$Q^\pi(s_t, a_t) \approx \underbrace{r_t}_{\text{current}} + \underbrace{\gamma V^\pi(s_{t+1})}_{\text{future}}$$

TD target

\Rightarrow P28, L14

Actor-Critic Algo

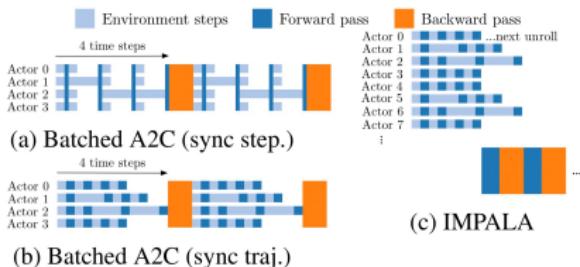
- * On-policy, low sample efficiency

A3C

- * Multiple actors updates the same parameters

IMPALA:

- * One learner learning from experiences from multiple actors



PPO: Makes training more stable

Recall Actor-Critic Algorithm:

$$\begin{aligned}\nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T [\gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) A^{\pi}(s_t^i, a_t^i)] \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T [\gamma^t \frac{\nabla_{\theta} \pi_{\theta}(a_t^i | s_t^i)}{\pi_{\theta}(a_t^i | s_t^i)} A^{\pi}(s_t^i, a_t^i)]\end{aligned}$$

$$\frac{d \log y}{dx} = \frac{1}{y} \frac{dy}{dx}$$

PPO

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T [\gamma^t \nabla_{\theta} \frac{\pi_{\theta}(a_t^i | s_t^i)}{\pi_{old}(a_t^i | s_t^i)} A^{\pi}(s_t^i, a_t^i)]$$

Ratio

$$r_t = \frac{\text{New prob of taking } a_t^i}{\text{old prob of taking } a_t^i}$$

$$1 - \epsilon \leq r_t \leq 1 + \epsilon$$

Soft Actor-Critic

* Better exploration

Hide and Seek

* Multiagent RL