

# Simple RL Algorithm

Hyperparameter: Initial policy parameters  $\theta$

linear,  $\theta = 0$ .  
ANN, use Weight init Schemes (He)

For each episode:

For each time  $t$ :

Agent observes  $S_t$

Agent selects  $A_t$  using  $\pi_\theta$

Environment responds by changing to state  $S_{t+1}$  and giving reward  $R_t$ .

$$V_i, \theta_i \leftarrow \theta_i + \alpha \left( \sum_{t=0}^{\infty} \gamma^t R_t \right) \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i}$$

For all times  $t$ .

How much more likely should we make  $A_t$ ?

How to change  $\theta_i$  to make  $A_t$  more likely in  $S_t$ .

$$\pi_\theta(s, a) = P(A_t = a | S_t = s)$$

~~IF  $\sum_{t=0}^{\infty} \gamma^t R_t$  is big  
for all times  $t$   
└ Make  $A_t$  more likely in  $S_t$   
Else  $\sum_{t=0}^{\infty} \gamma^t R_t$  is small  
for all times  $t$   
└ Make  $A_t$  less likely in  $S_t$~~

all policy parameters  $\theta_i$

$$\textcircled{1} V_i, \theta_i \leftarrow \theta_i + \alpha \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i}$$

$$\textcircled{2} V_i, \theta_i \leftarrow \theta_i - \alpha \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i}$$

# Simple RL Algorithm

Hyperparameter: Initial policy parameters  $\theta$

$\theta$  → linear,  $\theta = 0$ .  
 → ANN, use Weight init schemes (He)

For each episode:

For each time  $t$ :

Agent observes  $S_t$   
 Agent selects  $A_t$  using  $\pi_\theta$   
 Environment responds by changing to state  $S_{t+1}$  and giving reward  $R_{t+1}$ .

play the games

$$V_i, \theta_i \leftarrow \theta_i + \alpha \left( \sum_{t=0}^{\infty} \gamma^t R_{t+1} \right) \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i}$$

For all times  $t$ .

How much more likely should we make  $A_t$ ?  
 How to change  $\theta_i$  to make  $A_t$  more likely in  $S_t$ .

For each time  $t$ :

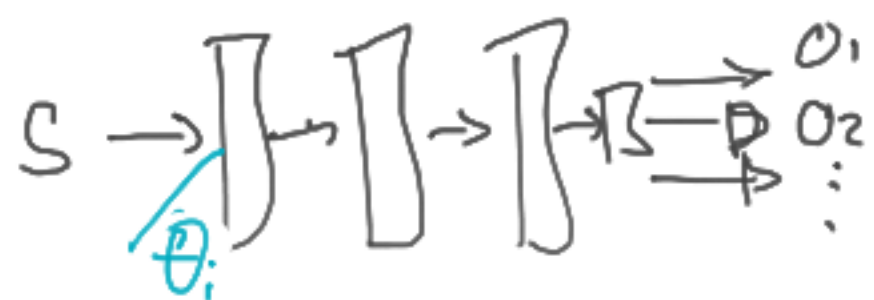
$$V_i, \theta_i \leftarrow \theta_i + \alpha \left( \delta_t \right) \frac{\partial \pi_\theta(S_t, A_t)}{\partial \theta_i}$$

$$\delta_t = R_{t+1} + \gamma V^\pi(S_{t+1}) - V^\pi(S_t)$$

learn from after the game.

only depends/uses most recent episode (game of Tic-Tac-Toe).

$t=5$



$$\pi(S, a) = \frac{e^{\theta a}}{\sum_{a'} e^{\theta a'}}$$

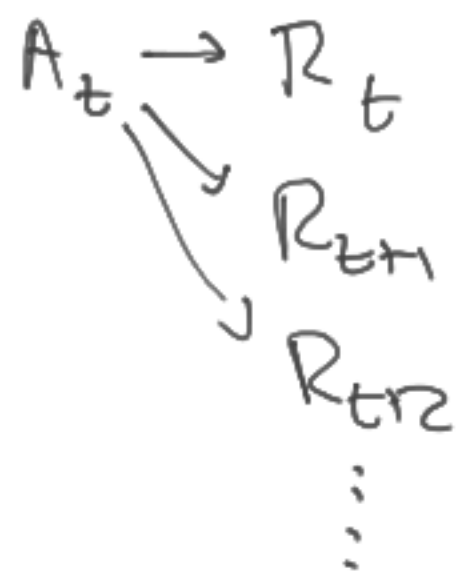
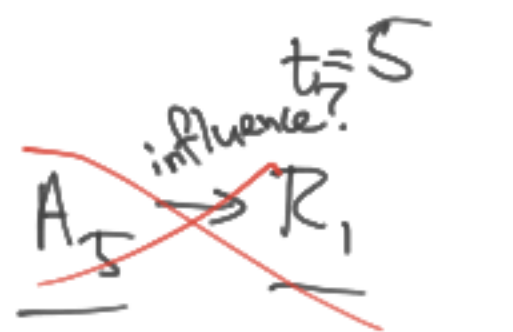
REINFORCE  
 unbiased

Policy update:

$$v_i, \theta_i \leftarrow \theta_i + \alpha \left( \sum_{t'=0}^{\infty} \gamma^{t'} R_{t'} \right) \frac{\partial \pi_{\theta}(S_t, A_t)}{\partial \theta_i}$$

How "good" was the outcome?

How to change  $\theta_i$  to make  $A_t$  more likely in state  $S_t$ .



$$\left( \sum_{k=0}^{\infty} \gamma^k R_{t+k} \right)$$

$$\theta_i \leftarrow \theta_i + \alpha \left( \sum_{k=0}^{\infty} \gamma^k R_{t+k} \right) \frac{\partial \pi_{\theta}(S_t, A_t)}{\partial \theta_i}$$

Note: Discounting starts at time  $t$ .

$$\gamma^0 R_t + \gamma^1 R_{t+1} + \gamma^2 R_{t+2}$$

How do we update before the end of the episode?

→ Not a good "outcome" that matters, but whether the outcome was ~~better~~ or worse than expected.

outcome → outcome relative to what the agent was expecting.

What is the agent expecting?

↳ Not what is  $S_{t+1}$  or  $S_{t+2}$

↳ Want to know what  $\sum_{k=0}^{\infty} \gamma^k R_{t+k}$  is going to be.

The value of state  $S_t$  under policy  $\pi$  is the amount of reward the agent expects to get if it uses policy  $\pi$  starting from state  $S_t$ .

$$v^{\pi}(s) = E \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s; \pi \right]$$

Does not depend on  $t$ .

$v^{\pi}$  → value function

- Assume  $v^\pi$  is given (for now)

Better than expected:

$$(S_t, A_t, R_t, S_{t+1})$$

$$\sum_{k=0}^{\infty} \gamma^k R_{t+k} > v^\pi(S_t)$$

$$R_t + (\gamma R_{t+1} + \gamma^2 R_{t+2}) > v^\pi(S_t)$$

$$R_t + \gamma (R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots) > v^\pi(S_t)$$

Worse than expected:

$$\sum_{k=0}^{\infty} \gamma^k R_{t+k} < v^\pi(S_t)$$

$S_t$  = state at time  $t$   
Rand. Variable.

$S$  = any state.

Not a R.V.

$$v^\pi(S_t)$$

$$v^\pi(S_t)$$

$$v^\pi(S)$$

$$\gamma = 0.5$$

$$\frac{x_1 \gamma}{x_1 + \gamma}$$

$$v^\pi(s) = \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s; \pi \right]$$

$$0 + \gamma (\gamma^0(1) + \gamma^1(2)) > 1$$

$$R_t + \gamma v^\pi(S_{t+1}) > v^\pi(S_t)$$

$$= \gamma^0(0) + \gamma^1(1) + \gamma^2(2)$$

$$= 1(0) + .5(1) + .25(2) = 1$$

$$R_t + \gamma v^\pi(S_{t+1}) < v^\pi(S_t)$$

$$R_t = 0$$

$$R_{t+1} = 1$$

$$R_{t+2} = 2$$

$$R_{t+3} = 2$$

Expecting 1

cookie in

1 min

2 cookies in

2 minutes

(3 total)

$$R_t + \gamma v^\pi(S_{t+1})$$

vs

$$v^\pi(S_t)$$

But not just  $R_t$ .

↑ after one time step.

$\delta_t$  = "was it better or worse than expected"  
= positive if better than expected  
= negative if worse than expected.

$$\delta_t = R_t + \gamma v^\pi(S_{t+1}) - v^\pi(S_t)$$

- 1)  $R_t$  <sup>bigger</sup> smaller than expected
- 2)  $S_{t+1}$  <sup>better</sup> worse than expected.

$\delta_t$  positive → outcome ( $R_t + \gamma v^\pi(S_{t+1})$ ) was better than expected.  
 $\delta_t$  negative → outcome was worse than expected.

$\delta_t$  = "temporal difference error"  
TD error.

$t$  = get MM  
 $t+1$  = April 1.

