

- Midterm & course curve.

lecture 15  
- Policy  
Gradient



$$J(\theta) = \mathbb{E} \left[ \underbrace{\sum_{t=0}^{\infty} \gamma^t r_t}_{\text{Return}} ; \theta \right]$$

$$\pi(s, a) = P_r(A_t = a | S_t = s)$$

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

# 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  $S_t$ .

~~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  $\theta_i$

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

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

# Simple RL Algorithm

Hyperparameter: Initial policy parameters  $\theta_0$

$\theta_0$  → 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  $R_t$ .

play the games

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

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  $t$ :

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

learn from after the game.

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

REINFORCE  
with bias issue

# Make $A_t$ more likely in $S_t$

$$V_i, \quad \theta_i \leftarrow \theta_i + \alpha \frac{\partial \pi_{\theta}(S, a)}{\partial \theta_i}$$

all weights in policy.

"how much to increase action probability."

How to change  $\theta_i$  to increase  $\pi_{\theta}(S, a) = P(A_t = a | S_t = S)$

Hint:  $\frac{\partial f(x, y)}{\partial y}$  = How to change  $y$  to increase  $f(x, y)$  as quickly as possible.

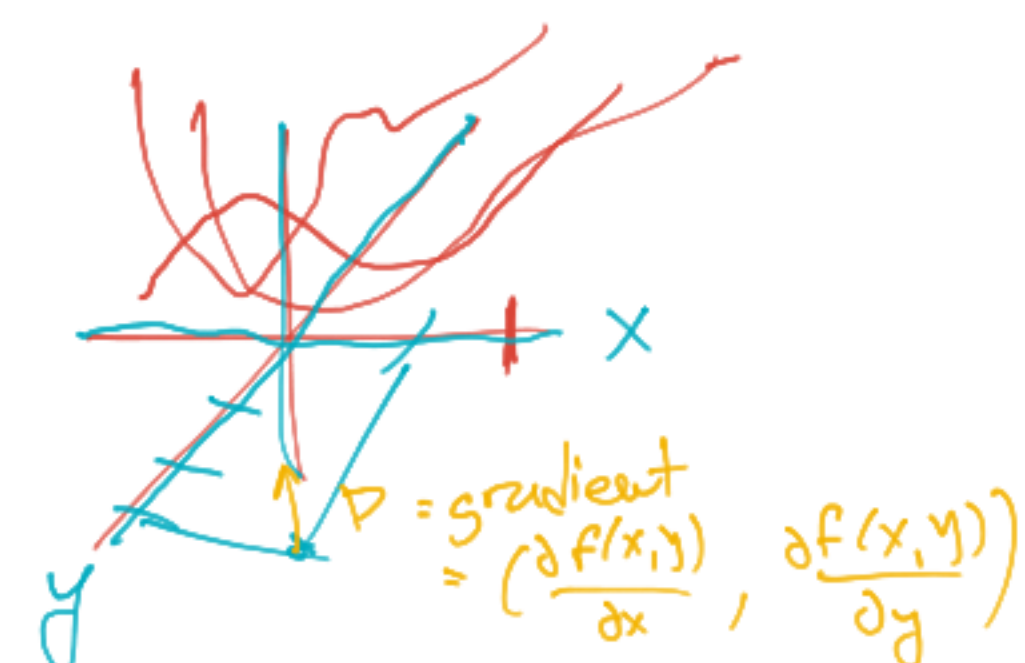
# Make $A_t$ less likely in $S_t$

$$V_i, \quad \theta_i \leftarrow \theta_i - \alpha \frac{\partial \pi_{\theta}(S, a)}{\partial \theta_i}$$

$$f(x, y) = x^2 + y^2 + 7 \sin(y)$$

$$(x, y) = (1, 3)$$

$$\frac{\partial f(x, y)}{\partial y} = 2y + 7 \cos(y) = -9.3$$



$\theta^i$  = weights for  $i$ th node in output layer.

$$\theta^i = (\theta_{i1}, \theta_{i2}, \dots) \in \mathbb{R}$$

$$\theta = (\theta^1, \theta^2, \theta^3, \dots)$$

$$\pi_{\theta}(S, a)$$

$$\pi(S, a, \theta)$$

Is  $\sum_{t=0}^{\infty} \gamma^t R_t$  big or small?

Idea: Scale how much we increase or decrease action probabilities by  $\sum_{t=0}^{\infty} \gamma^t R_t$ .

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

*handles both good/bad outcomes! (positive and negative  $\sum \gamma^t R_t$ )*

Step size to control learning speed. "learning rate"

like a step size, saying how much of a change to make to  $P_i(A_t | S_t)$



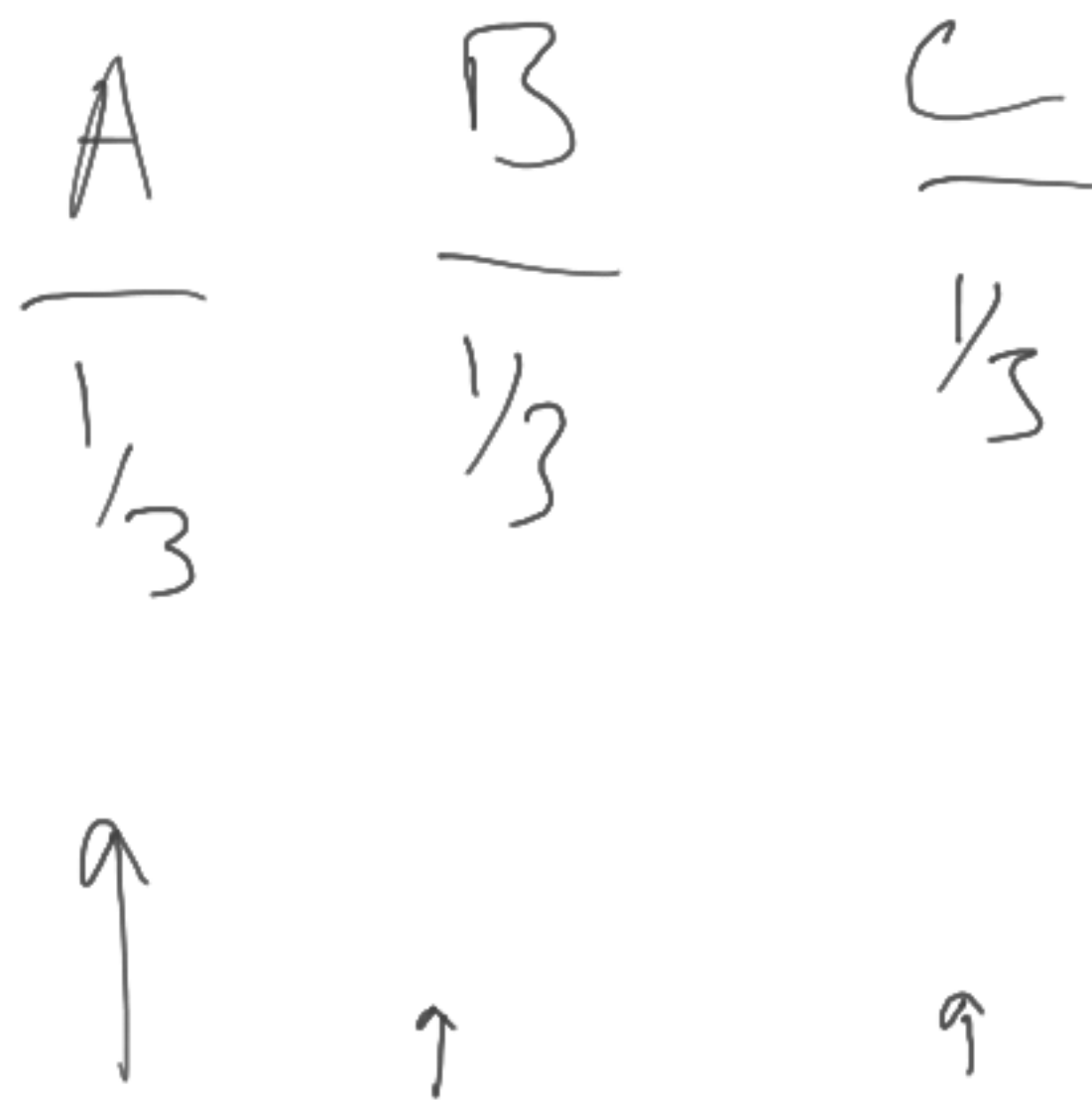
$\gamma = 0.5$

$\times J(\text{circles})$

$$1 + .5(1) + .5^2(1) + \dots = 1 + \frac{1}{2} + \frac{1}{4} + \dots = 2$$

$J(\text{go to end}) = 1 + .5(1) + .5^2(1) + \dots + .5^3(1000)$

*(Note:  $\gamma^t$  is under the first term and  $R_t$  is under the last term in the second equation)*



## Softmax

