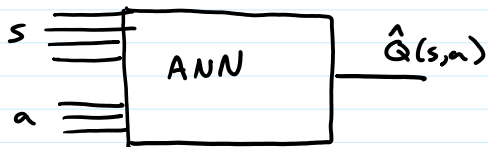


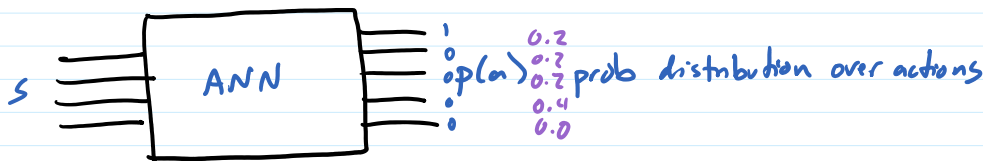
# Deep Q Learning



two networks: learning target

play using learning observe  $(s_t, a_t, s'_t, r_t) \rightarrow$  add to replay database  
 compute error  $\hat{Q}_t(s_t, a_t) - (r_t + \gamma \max_{a'} \hat{Q}_t(s'_t, a'))$   
 adjust learning accordingly  $\hookrightarrow (s_t, a_t, s'_t, r_t)$  sampled from replay database  
 periodically copy learning to target

# REINFORCE: learn policy directly



objective: maximize  $E \left[ \sum_{t=1}^T r_t \right]$

$(s_1, a_1) (s_2, a_2) (s_3, a_3) \dots$

over all possible trajectories  $\tau$  where probability of  $\tau$  is given by ANN and its weights  $\theta$

objective as fun of weights on neural network

$$J(\theta) = \int \underbrace{\pi_{\theta}(\tau)}_{\text{prob ANN w/ weights } \theta \text{ gives trajectory } \tau} \cdot \underbrace{r(\tau)}_{\text{sum of reward at each step}} d\tau$$

gradient w.r.t. weights  $\theta$

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} \pi_{\theta}(\tau) \cdot r(\tau) d\tau$$

how to adjust weights  $\theta$  to increase  $J(\theta)$  most? so in dir of steepest ascent

$$= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) d\tau$$

$$\begin{aligned} & \pi_{\theta}(\tau) \cdot \nabla_{\theta} \log \pi_{\theta}(\tau) \\ &= \cancel{\pi_{\theta}(\tau)} \cdot \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\cancel{\pi_{\theta}(\tau)}} \end{aligned}$$

$$= E \left[ \nabla_{\theta} \log \pi_{\theta}(\tau) \cdot r(\tau) \right]$$

$$= \cancel{\pi_{\theta}(\tau)} \cdot \frac{(\nabla_{\theta} \pi_{\theta}(\tau))}{\cancel{\pi_{\theta}(\tau)}}$$

$$= E \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \cdot \left( \sum_{t=1}^T r(s_t, a_t) \right) \right]$$

↑  
 still over  
 all trajectories  $\tau$   
 with prop  $\pi_{\theta}(\tau)$

$\pi_{\theta}(\tau)$  = prob of trajectory  $\tau$

$$= \pi_{\theta}((s_1, a_1), (s_2, a_2), \dots)$$

prob ANN w/ weights  $\theta$  chooses  $a_t$  in  $s_t$

estimate expected value by running policy  $N$  times

$$= p(s_1) \cdot \pi_{\theta}(a_1 | s_1) \cdot p(s_2 | s_1, a_1) \cdot \pi_{\theta}(a_2 | s_2) \dots$$

prob of transition in MDP

$$\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \left( \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \cdot \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

output of ANN

$$= p(s_1) \cdot \prod_{i=1}^T \pi_{\theta}(a_i | s_i) \cdot p(s_{t+1} | s_t, a_t)$$

$$\log \pi_{\theta}(\tau) = \log$$

$$= \log p(s_1) + \sum_{i=1}^T (\log \pi_{\theta}(a_i | s_i) + \log p(s_{t+1} | s_t, a_t))$$

$$\nabla_{\theta} \log \pi_{\theta}(\tau) = \nabla_{\theta} \log p(s_1) + \sum_{i=1}^T (\log \pi_{\theta}(a_i | s_i) + \log p(s_{t+1} | s_t, a_t))$$

↔ not functions of  $\theta$  ↔

## REINFORCE

$$\tau_i = (s_{i,1}, a_{i,1}) (s_{i,2}, a_{i,2}) \dots$$

get sample trajectory by running policy

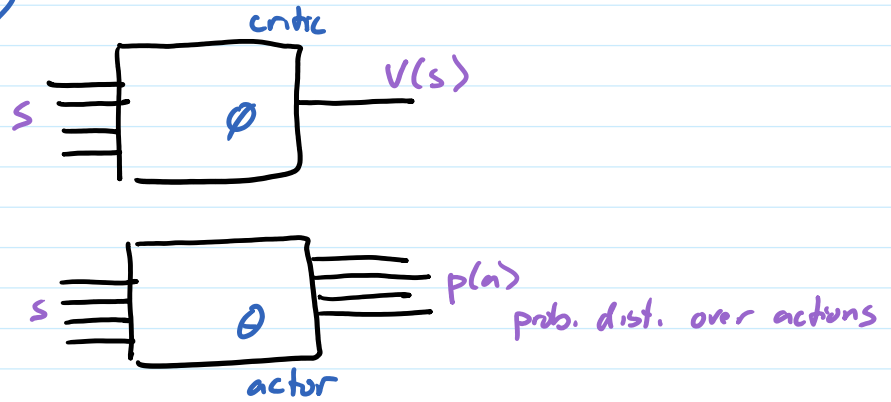
compute estimate of  $\nabla_{\theta} J(\theta)$

$$\text{update } \theta : \theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

↑  
 learning rate

(A2C)

Two networks



get sample trajectory using actor

fit critic to observed rewards

evaluate advantage

↳ difference between critic's expectation and actual result

critic predicts  
observed

$$\hat{V}_\phi(s_t) - (r_t + \hat{V}_\phi(s_{t+1}))$$

compute gradient  $\nabla_\theta J(\theta) = \sum_i \nabla_\theta \log \pi_\theta(a_i | s_i) \cdot \hat{A}(s_i, a_i)$

do gradient ascent on actor  $\theta \leftarrow \theta + \alpha \cdot \nabla_\theta J(\theta)$