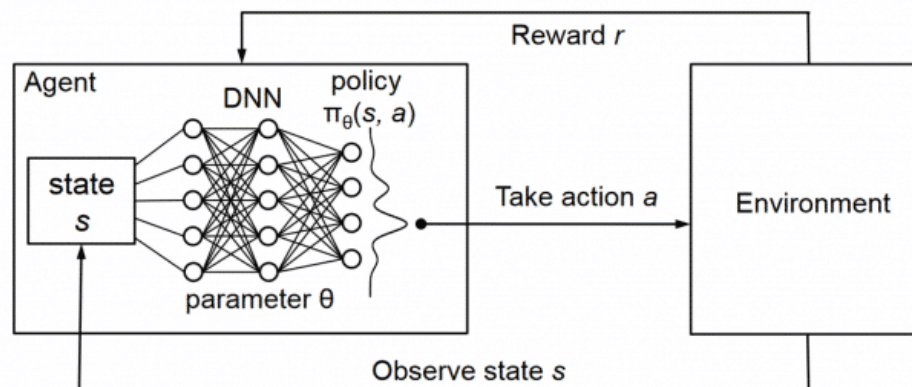
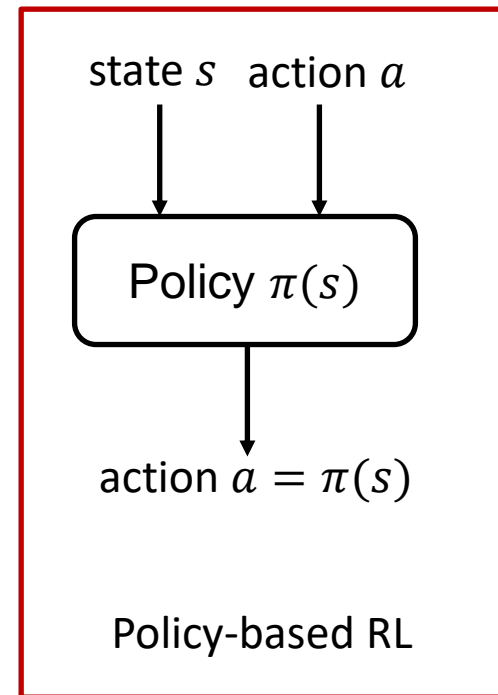
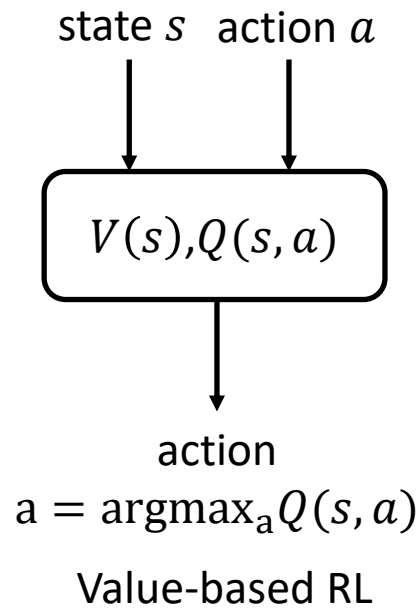
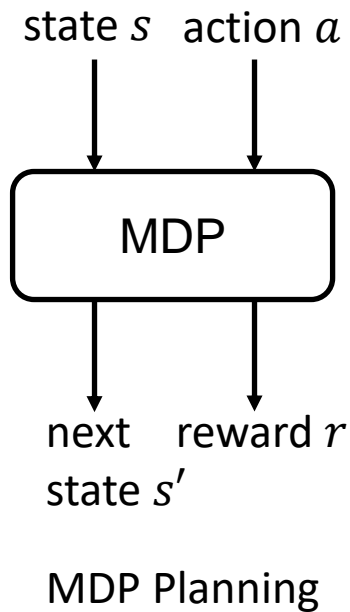


# L7.3 Policy-based RL

Zonghua Gu 2021

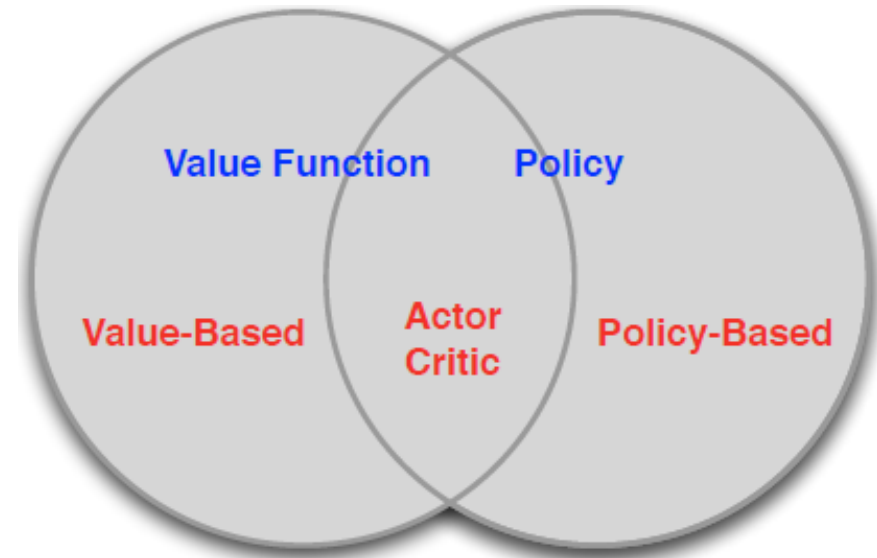


# Policy-based RL



# CH13 Policy Gradient Methods

- Value Based
  - Learnt Value Function
  - Implicit policy (e.g.  $\epsilon$ -greedy)
- Policy Based
  - No Value Function
  - Learnt Policy
- Actor-Critic
  - Learnt Value Function
  - Learnt Policy

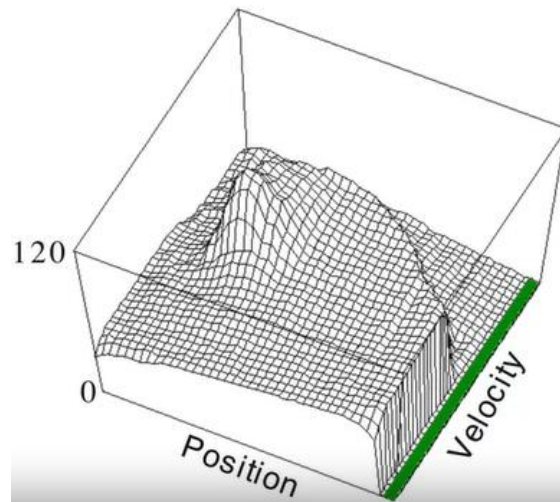
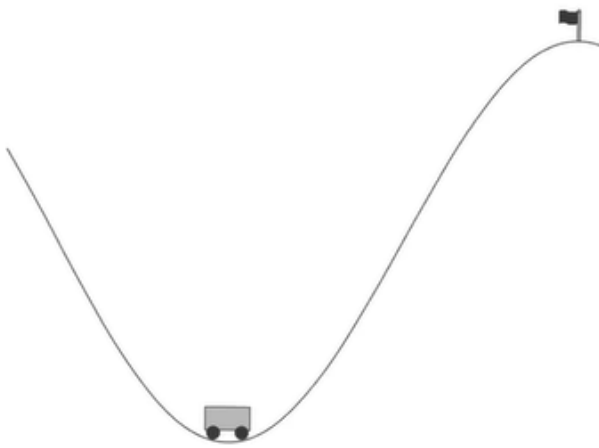


# Policy-based RL Pros and Cons

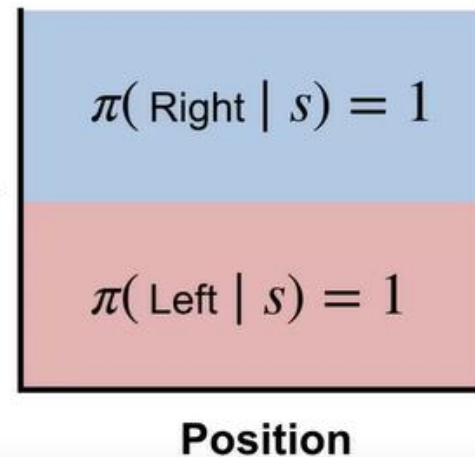
- Pros:
  - Effective in high-dimensional or continuous action space.
    - Value-based RL is only applicable to discrete action space; inefficient to discretize continuous actions for high-dim action space, as taking  $\operatorname{argmax}_a Q(s, a)$  may be expensive.
  - Can learn stochastic policies
    - Value-based RL learns a near-deterministic policy (greedy or  $\epsilon$ -greedy).
  - Policy typically converges faster than value functions.
- Disadvantages:
  - Typically converges to a local rather than global optimum.
  - Evaluating a policy is typically inefficient and high variance.

# Mountain Car Example

- Middle: a complex value function
- Right: a simple policy that works well: accelerate in the direction of current velocity.

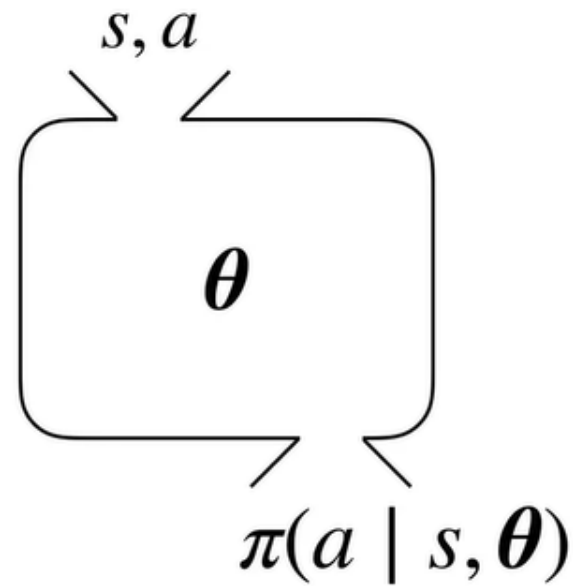
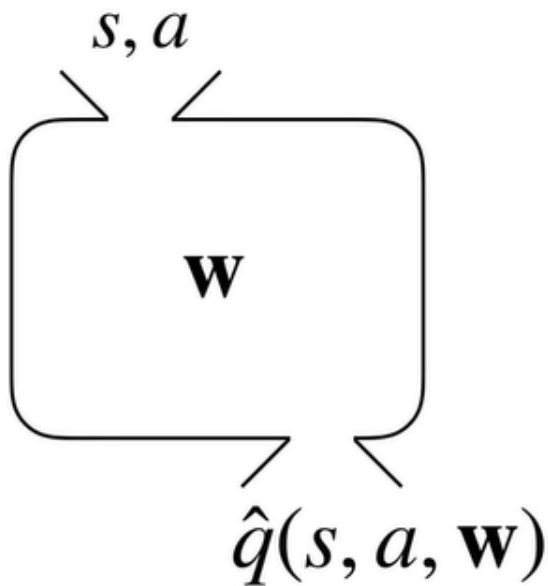


moving  
right  
**Velocity**  
moving  
left



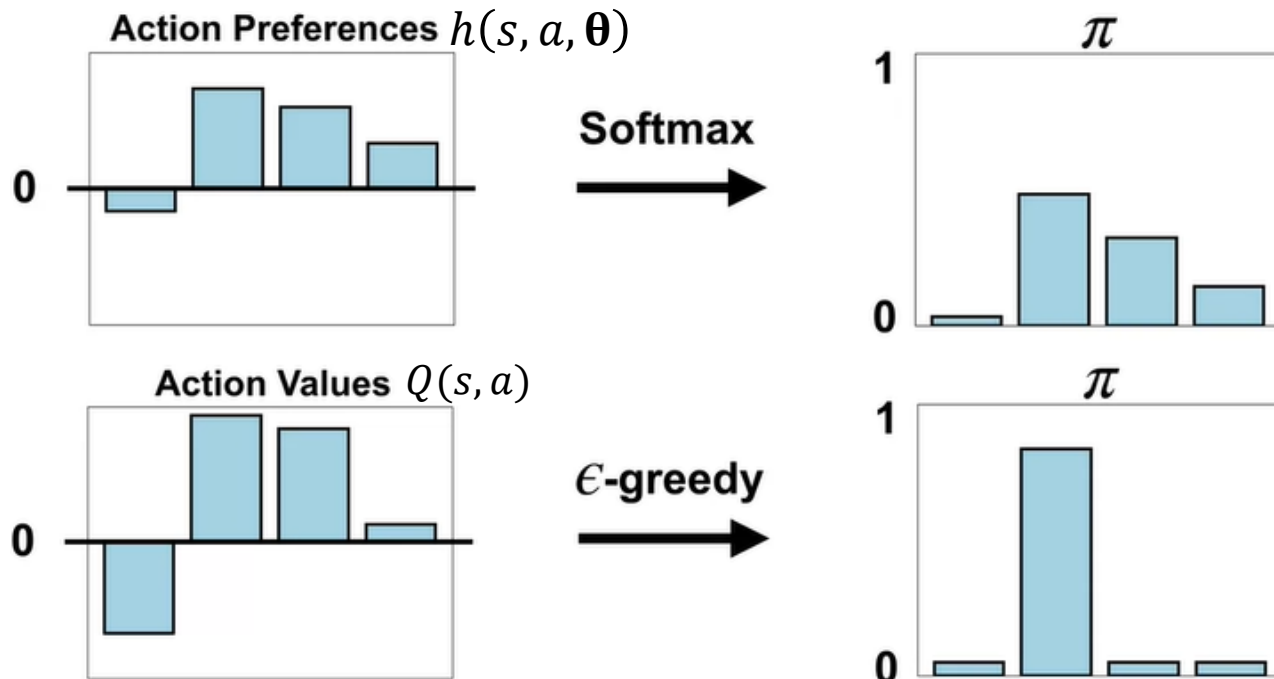
# Function Approximation for Action Value Function vs. Policy

- Value-based RL learns a function approximation for action value function  $\hat{q}(s, a, \mathbf{w})$ .
  - Deterministic policy  $a = \operatorname{argmax} \hat{q}(s, a, \mathbf{w})$  (may be  $\epsilon$ -greedy during training)
- Policy-based RL learns a function approximation for stochastic policy  $\pi(a|s, \boldsymbol{\theta})$ :
  - Probability that action  $a$  is taken in state  $s$ , with parameter  $\boldsymbol{\theta}$ . The actual action taken is sampled from the probability distribution  $A \sim \pi(a|s, \boldsymbol{\theta})$
  - Probability must be non-negative:  $\pi(a|s, \boldsymbol{\theta}) \geq 0, \forall a \in \mathcal{A} \wedge \forall s \in \mathcal{S}$
  - Probabilities must sum to 1:  $\sum_{a \in \mathcal{A}} \pi(a|s, \boldsymbol{\theta}) = 1, \forall s \in \mathcal{S}$



# SoftMax vs. $\epsilon$ -greedy for Discrete Actions

- SoftMax policy:  $\pi(a|s, \theta) \doteq \frac{e^{h(s,a,\theta)}}{\sum_{a' \in \mathcal{A}} e^{h(s,a',\theta)}}$ 
  - $h(s, a, \theta)$  is action preference, which may be a linear function  $\theta^T \mathbf{x}(s, a)$ , or the logit from the SoftMax layer of a DNN
  - A bad action with very negative  $h(s, a, \theta)$  will be very unlikely to be selected
  - Action probabilities change smoothly as a function of  $h(s, a, \theta)$
- $\epsilon$ -greedy: select the greedy action  $\underset{a}{\operatorname{argmax}} Q(s, a)$  with prob  $1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|}$ 
  - No distinction between policies that are not the optimal one; A bad action with very low  $Q(s, a)$  will be selected with equal prob as all other non-optimal policies
  - Action probabilities may change dramatically for an arbitrarily small change in  $Q(s, a)$ , if that change results in a different optimal action



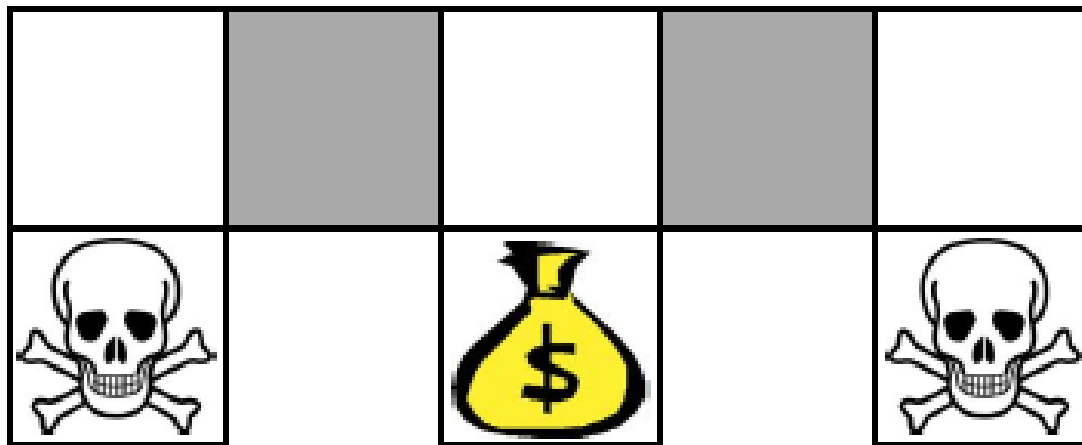
# Stochastic Policy vs. Deterministic Policy

- Example 1: two-player game of rock-paper-scissors
  - Scissors beat paper; paper beats rock; rock beats scissors
  - For iterated game, a deterministic policy is easily exploited by the opponent; a uniform random policy is optimal, and achieves Nash equilibrium
- Example 2: Aliased Grid World (POMDP)



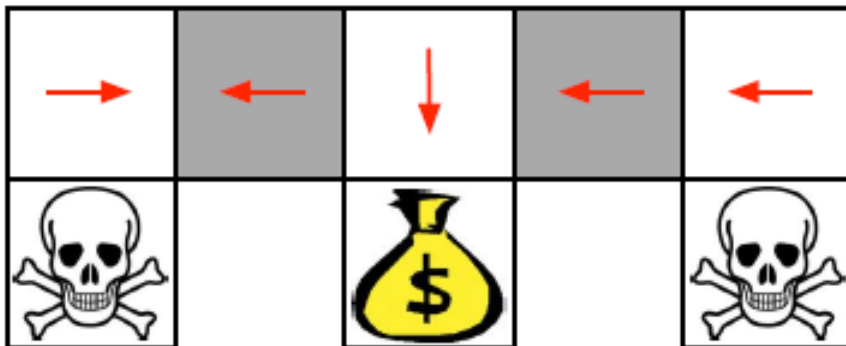
# Aliased Grid World (POMDP)

- Env has 3 terminal states: 2 with high positive reward and 1 with high negative reward. It is a Partially Observable MDP (POMDP): Agent cannot observe its position directly; it can only observe features of the following form (for all directions  $d_1, d_2, \dots \in (N, E, S, W)$ ):
  - $\phi(s) = \mathbf{1}(\text{wall to } d_1, d_2 \dots)$  (it can detect walls, e.g., w. Radar or Lidar)
    - $\mathbf{1}(x)$  is indicator function:  $\mathbf{1}(x) = 1$  if  $x = \text{true}$
  - Agent cannot differentiate between the 2 grey states
- Value-based RL learns a deterministic policy:  $a = \operatorname{argmax} \hat{q}(s, a, \mathbf{w}) = \operatorname{argmax} f(\phi(s), \mathbf{w})$
- Policy-based RL learns a stochastic policy:  $\pi(a|s, \boldsymbol{\theta}) = g(\phi(s), \boldsymbol{\theta})$

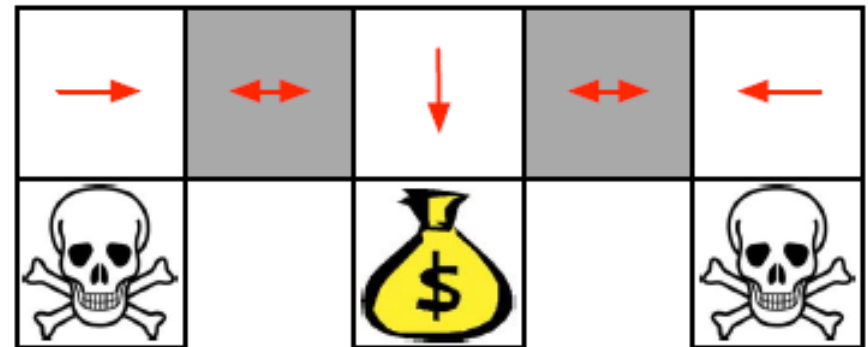


# Aliased Grid World (POMDP)

- Left: an optimal deterministic policy:
  - Either move  $W$  in both grey states (red arrows), or move  $E$  in both grey states; Either way, agent can get stuck and never reach the money
- Right: the optimal stochastic policy:
  - Randomly move  $E$  or  $W$  in grey states:  $\pi(\text{move } E | \text{wall to N and S}, \theta) = \pi(\text{move } W | \text{wall to N and S}, \theta) = 0.5$
  - Agent will likely reach the goal state quickly
- How about adding  $\epsilon$ -greedy on top of the opt det policy?
  - Agent may get into one of the 2 bad states, since each non-optima action is given equal probability in every state



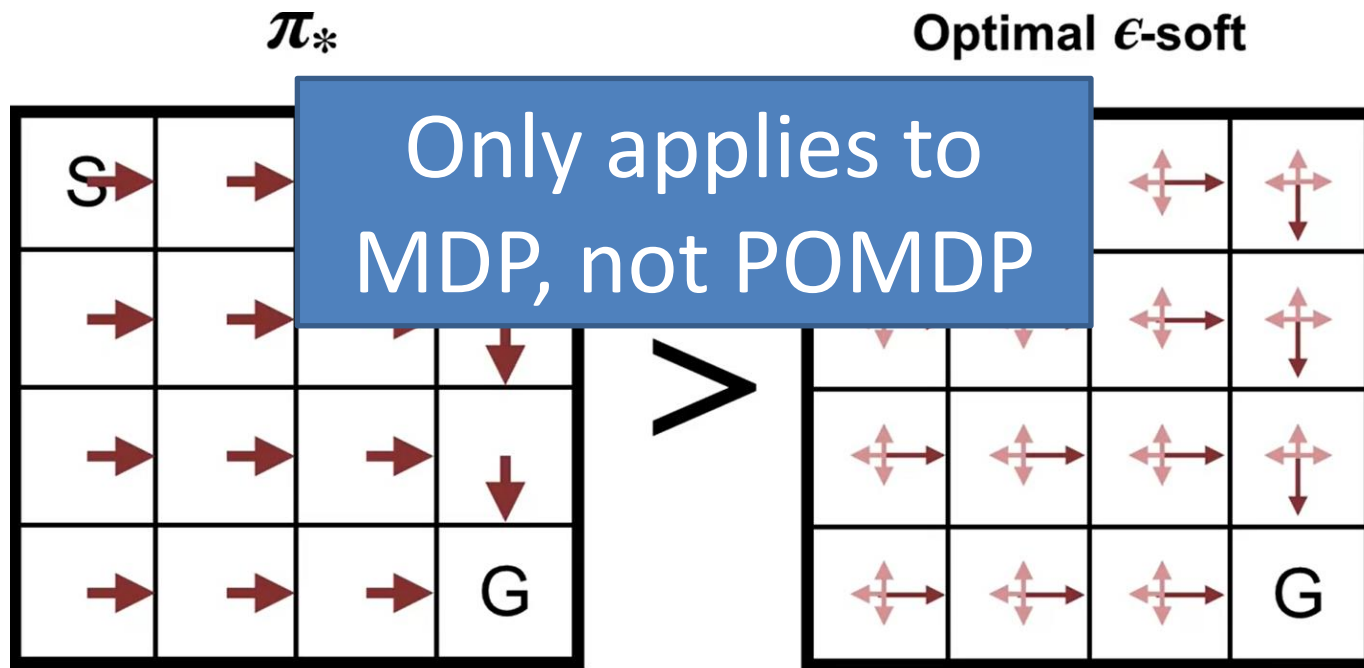
Opt det policy



Opt Sto policy

# Optimal $\epsilon$ -Soft Policy

- The optimal  $\epsilon$ -soft policy is the policy with the highest value in each state among all  $\epsilon$ -soft policies. It performs worse than the optimal greedy deterministic policy  $\pi_*$  in general.
- But it often performs reasonably well, and avoids exploring starts.

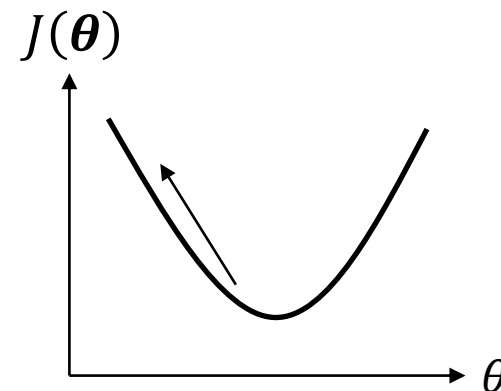
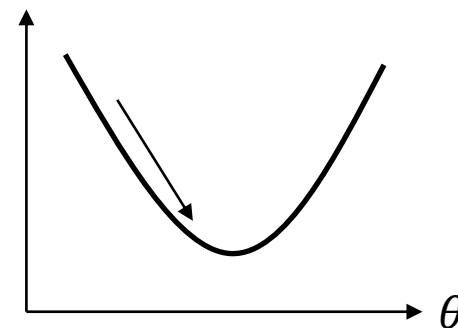


# Optimization Objective

- We consider the episodic case, and would like to optimize the expected value of the start state of each episode, with start state  $s_0$ :
- $J(\boldsymbol{\theta}) \doteq v_{\pi}(s_0)$
- Where  $v_{\pi}(s_0)$  is the true value function for policy  $\pi$ , parametrized by  $\boldsymbol{\theta}$ :  $\pi(a|s, \boldsymbol{\theta})$

# Model Training in Supervised Learning vs. Policy Gradient in RL

- SL: to solve  $\min_{\theta} \mathbb{E}_{(x,y) \sim D} \text{Loss}(x, y; \theta)$   
for model training: gradient **descent**  $\theta \leftarrow \theta - \alpha \nabla_{\theta} \text{Loss}(x, y; \theta)$ 
  - Update **model params**  $\theta$  by following the gradient **downhill**, in order to **decrease**  $\text{Loss}(x, y; \theta)$ . ( $\alpha$  is the Learning Rate)
- RL: to solve  $\max_{\theta} J(\theta)$  in Policy Gradient: gradient **ascent**  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ 
  - Update **policy model params**  $\theta$  by following the gradient **uphill**, in order to **increase**  $J(\theta)$
- We use  $\nabla$  as shorthand for  $\nabla_{\theta}$

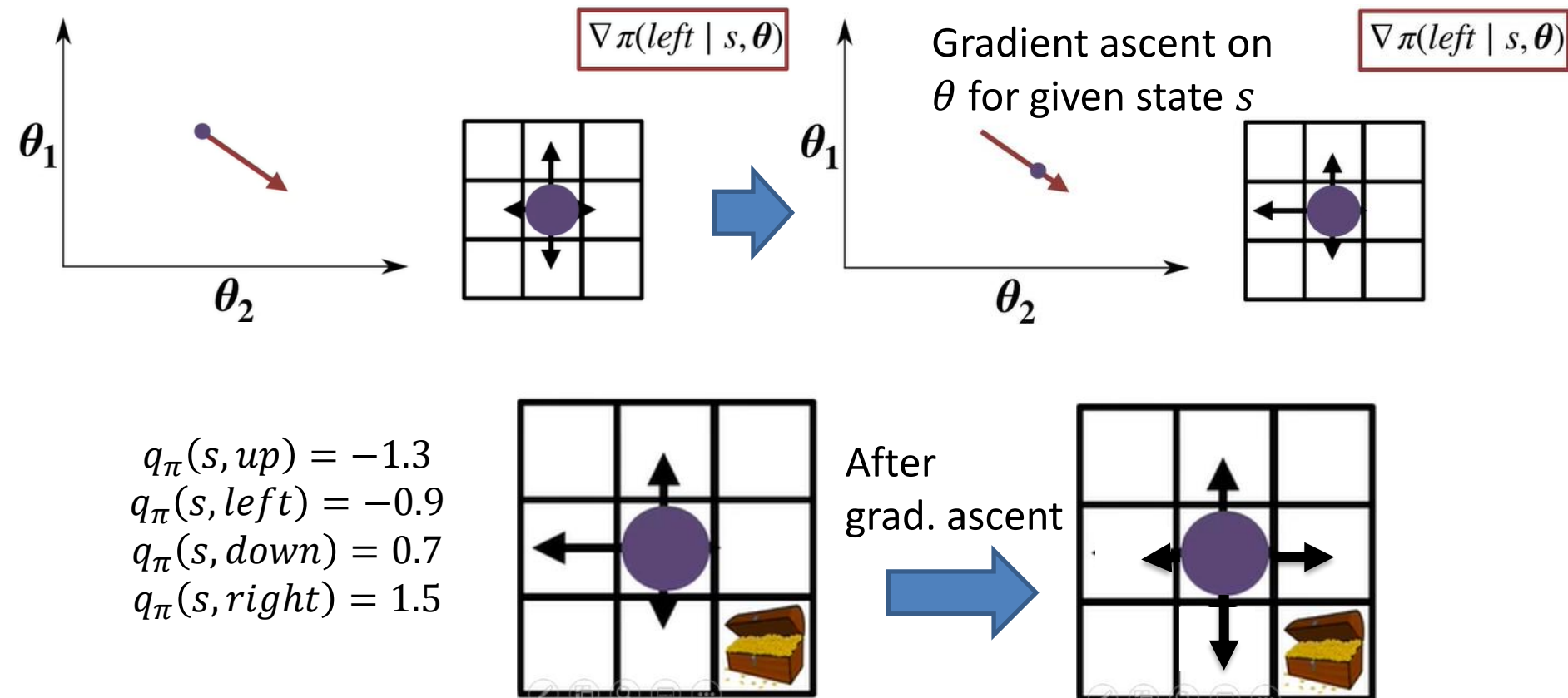


# Policy Gradient Theorem

- $\nabla J(\boldsymbol{\theta}) = \nabla v_{\pi}(s_0)$
- $= \nabla [\sum_a \pi(a|s, \boldsymbol{\theta}) q_{\pi}(s, a)]$
- $= \sum_a [\nabla \pi(a|s, \boldsymbol{\theta}) q_{\pi}(s, a) + \pi(a|s, \boldsymbol{\theta}) \nabla q_{\pi}(s, a)]$
- Policy Gradient Theorem (proof in RLBook p. 325; assuming discount factor  $\gamma = 1$ ):
- $\nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a \nabla \pi(a|s, \boldsymbol{\theta}) q_{\pi}(s, a)$ 
  - $\mu(s)$ : on-policy distribution under policy  $\pi$ ,  $\sum_{s \in \mathcal{S}} \mu(s) = 1$ .  
A larger  $\mu(s)$  denotes state  $s$  is visited more frequently, hence the policy gradient term for state  $s$  is given more weight (i.e., we care more about frequently-visited states than rarely-visited states)

# Policy Gradient Example

- The gradient  $\nabla \pi(\text{left} | s, \theta)$  tells us how to change the policy parameters  $\theta$  to make action *left* more likely to be selected in state  $s$ . By gradient ascent on  $\theta$ , we increase the probability for taking action *left* in state  $s$ .



# MC REINFORCE

- $\nabla J(\boldsymbol{\theta}) = \sum_s \mu(s) \sum_a \nabla \pi(a|s, \boldsymbol{\theta}) q_\pi(s, a)$
- $= \mathbb{E}_\pi [\sum_a \nabla \pi(a|S_t, \boldsymbol{\theta}) q_\pi(S_t, a)]$ 
  - Outer exp over policy  $\pi$ : average over all experienced states under policy  $\pi$ , i.e.,  $S_t \sim s$
- $= \mathbb{E}_\pi \left[ \sum_a \pi(a|S_t, \boldsymbol{\theta}) \frac{\nabla \pi(a|S_t, \boldsymbol{\theta})}{\pi(a|S_t, \boldsymbol{\theta})} q_\pi(S_t, a) \right]$
- $= \mathbb{E}_\pi [\sum_a \pi(a|S_t, \boldsymbol{\theta}) \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}) q_\pi(S_t, a)]$  (since  $\nabla \log f(x) = \frac{\nabla f(x)}{f(x)}$ )
- $= \mathbb{E}_\pi [\mathbb{E}_\pi \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}) q_\pi(S_t, A_t)]$ 
  - Inner exp over policy  $\pi$ : average over all experienced actions  $A_t$  from state  $S_t$  under policy  $\pi$ , i.e.,  $A_t \sim \pi(a|S_t, \boldsymbol{\theta})$
- $= \mathbb{E}_\pi [\nabla \log \pi(A_t|S_t, \boldsymbol{\theta}) q_\pi(S_t, A_t)]$ 
  - Outer and inner exp over policy  $\pi$  combined: execute policy  $\pi$  and average over all experienced states  $S_t$  and actions  $A_t$  from state  $S_t$ , i.e.,  $S_t \sim s, A_t \sim \pi(a|S_t, \boldsymbol{\theta})$
- $= \mathbb{E}_\pi [\nabla \log \pi(A_t|S_t, \boldsymbol{\theta}) G_t]$ 
  - One way to approximate  $q_\pi(S_t, A_t)$  is by taking expectation over the return  $G_t$  in every episode for taking action  $A_t$  in state  $S_t$  ( $\mathbb{E}_\pi [G_t|S_t, A_t] = q_\pi(S_t, A_t)$ )
- SGD update ( $\gamma = 1$ ):  $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}_t) G_t$
- SGD update ( $\gamma < 1$ ):  $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \gamma^t \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}_t) G_t$  (proof omitted)



# PG Variants

- Monte Carlo REINFORCE:
  - $\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\pi}[\log \pi(A_t|S_t, \boldsymbol{\theta}) G_t]; \boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}_t) G_t$
- Monte Carlo REINFORCE with a baseline of estimated state value function:
  - $\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\pi}[\log \pi(A_t|S_t, \boldsymbol{\theta}) (G_t - \hat{v}(S_t, \mathbf{w}))]; \boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}_t) (G_t - \hat{v}(S_t, \mathbf{w}))$
- Actor-Critic:
  - $\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\pi}[\log \pi(A_t|S_t, \boldsymbol{\theta}) \delta_t]$ , TD error  $\delta_t = R_{t+1} + \gamma \hat{v}_{\pi}(S_{t+1}, \mathbf{w}_t) - \hat{v}_{\pi}(S_t, \mathbf{w}_t)$
  - $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}_t) \delta_t$
- Q Actor-Critic:
  - $\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\pi}[\log \pi(A_t|S_t, \boldsymbol{\theta}) q_{\pi}(S_t, A_t)]; \boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}_t) q_{\pi}(S_t, A_t)$
- Advantage Actor-Critic:
  - $\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\pi}[\log \pi(A_t|S_t, \boldsymbol{\theta}) A_{\pi}(S_t, A_t)]$ ,  $A_{\pi}(S_t, A_t) = q_{\pi}(S_t, A_t) - v_{\pi}(S_t)$
  - $A_{\pi}(S_t, A_t)$  captures the advantage of taking action  $A_t$  in state  $S_t$  then follow policy  $\pi$ , compared always following policy  $\pi$  from state  $S_t$
  - $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \nabla \log \pi(A_t|S_t, \boldsymbol{\theta}_t) A_{\pi}(S_t, A_t)$

# MC REINFORCE Pseudo-Code

- At the end of each episode, for each timestep  $0 \leq t < T$ :
  - Calculate the return  $G_t$  from each timestep  $t$
  - Update policy params  $\theta$  with SGD

## REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for $\pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Algorithm parameter: step size  $\alpha > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (G_t)$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \theta)$$

# Actor-Critic Pseudo-Code

- For update to critic params  $\mathbf{w}$ , refer to L7.2 Value-based RL, p 75 “Semi-Gradient TD(0) for Estimating  $\hat{v} \approx v_\pi$ ”

## One-step Actor–Critic (episodic), for estimating $\pi_\theta \approx \pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Input: a differentiable state-value function parameterization  $\hat{v}(s, \mathbf{w})$

Parameters: step sizes  $\alpha^\theta > 0$ ,  $\alpha^\mathbf{w} > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

    Initialize  $S$  (first state of episode)

$I \leftarrow 1$

    Loop while  $S$  is not terminal (for each time step):

$A \sim \pi(\cdot|S, \theta)$

        Take action  $A$ , observe  $S', R$

$\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$       (if  $S'$  is terminal, then  $\hat{v}(S', \mathbf{w}) \doteq 0$ )

$\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \delta \nabla \hat{v}(S, \mathbf{w})$

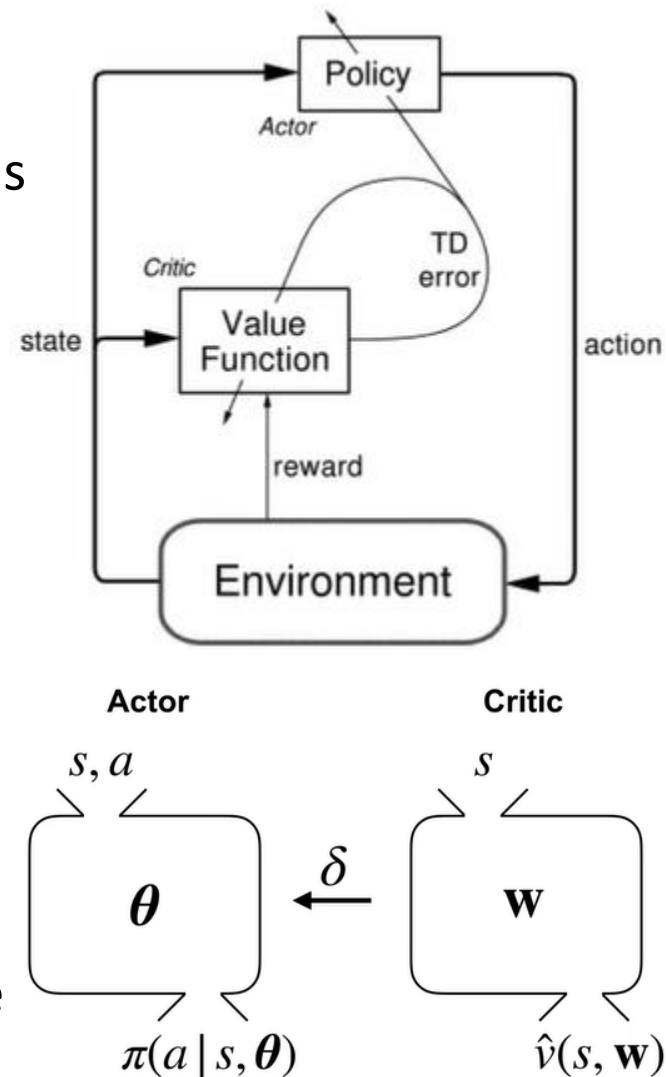
$\theta \leftarrow \theta + \alpha^\theta I \delta \nabla \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

$S \leftarrow S'$

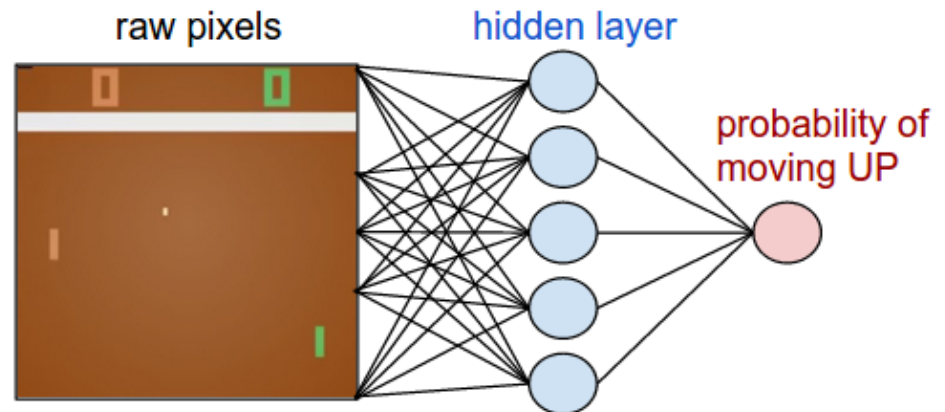
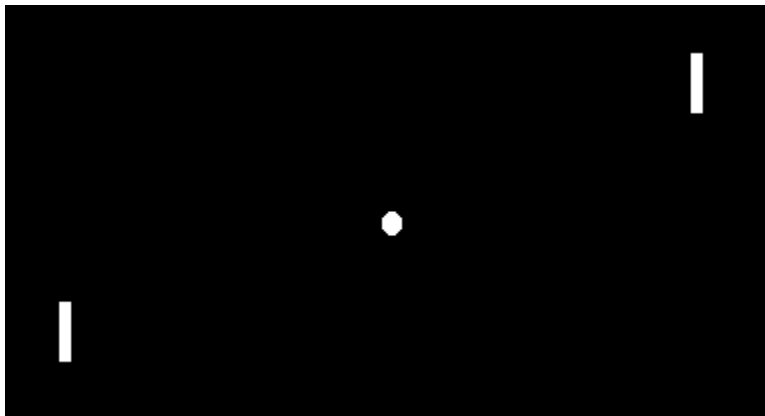
# Actor-Critic Explanations

- After each step of taking action  $A_t$  in state  $S_t$ :
- Critic computes TD error  $\delta_t = R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t)$ , and updates its params with semi-gradient TD(0)  $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha^w \delta_t \nabla_w \hat{v}(S_t, \mathbf{w}_t)$  (learning rate  $\alpha^w$ )
- Actor updates its params with Policy Gradient  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha^\theta \gamma^t \delta_t \nabla_{\boldsymbol{\theta}} \log \pi(A_t | S_t, \boldsymbol{\theta}_t)$ . If  $\delta_t > 0$ , then it means  $A_t$  resulted in a higher (one-step estimate) value than the expected  $\hat{v}(S_t, \mathbf{w}_t)$ , so probability of  $A_t$  in state  $S_t$  is increased; if  $\delta_t < 0$ , it is decreased (learning rate  $\alpha^\theta$ )
- Actor and Critic learn at the same time, constantly interacting. The actor is continually changing the policy params  $\boldsymbol{\theta}$  to exceed the critic's expectation, and the critic is constantly updating its value function params  $\mathbf{w}$  to evaluate the actor's changing policy.



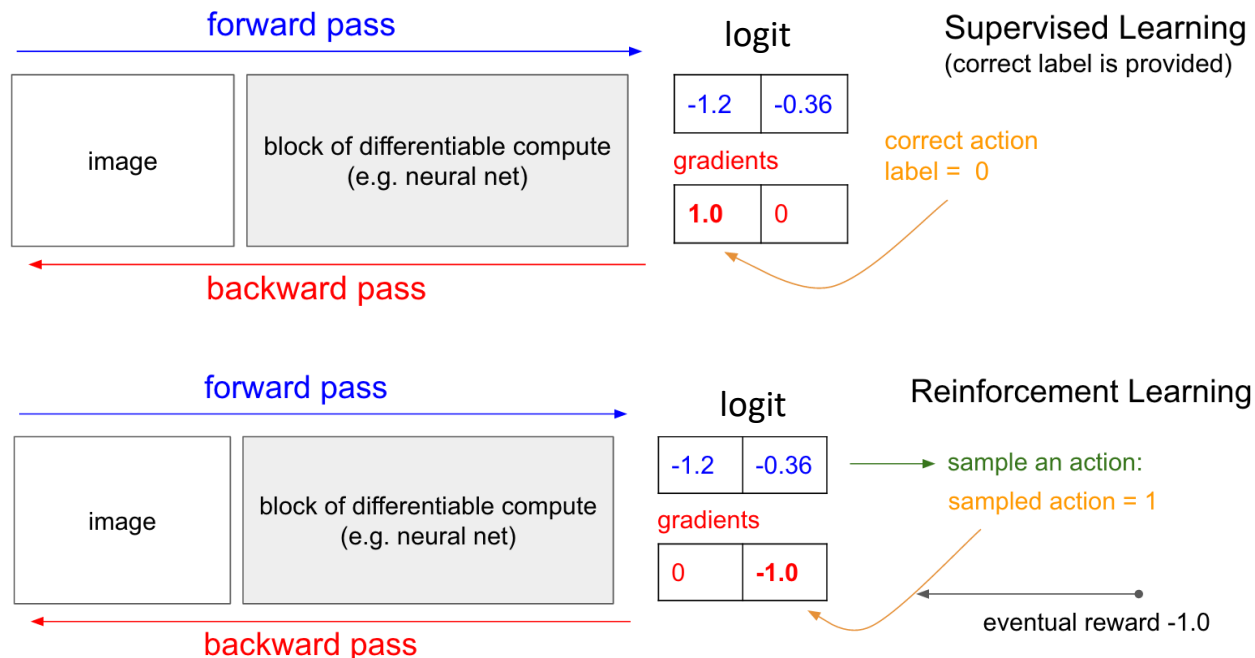
# Example: Game of Pong

- The agent plays one of the paddles (the other is controlled by a decent AI) and it has to bounce the ball past the other player. For each input image, agent decides if it wants to move the paddle UP or DOWN (2 discrete actions). At the end of the game the agent either wins (+1 reward) or loses (-1 reward)(sparse rewards).
  - (In practice we may stack multiple input images as input to the agent.)
- The agent implements a policy network, which maps from input image to two possible actions (UP or DOWN) with a stochastic SoftMax policy.



# SL vs. RL

- With SL: if the NN predicts  $y = UP$  for input  $x$ , and it is the correct ground truth label (e.g., the expert action in Imitation Learning), then gradient descent ( $\nabla_{\theta} \log p(y = UP|x)$ ) will make the NN more likely to predict  $y = UP$  for input  $x$
- With RL: if the NN predicts  $y = DOWN$  for input  $x$ , but we don't have the ground truth label in the middle of the game, so we must wait until the end of the game. If agent loses, then gradient descent ( $\nabla_{\theta} \log p(y = DOWN|x)$ ) will make the NN less likely to predict  $y = DOWN$  for input  $x$ 
  - Credit assignment problem: in an episode of many steps, which step contributed the most to the final outcome?



# MC REINFORCE for Pong

- $J(\theta) = \sum_s \mu_\pi(s) \sum_a \pi(a|s, \theta) q_\pi(s, a)$
- Consider agent in a given state  $S_t$  at time step  $t$ , and we want to select policy params  $\theta$  to maximize
- $v_\pi(S_t) = \sum_a \pi(a|S_t, \theta) q_\pi(S_t, a)$
- $= \pi(U|S_t, \theta) q_\pi(S_t, U) + \pi(D|S_t, \theta) q_\pi(S_t, D)$
- Taking derivative w.r.t policy params  $\theta$
- $\nabla_\theta v_\pi(S_t) = \sum_a \nabla_\theta \pi(a|S_t, \theta) q_\pi(S_t, a)$
- $= \nabla_\theta \pi(U|S_t, \theta) q_\pi(S_t, U) + \nabla_\theta \pi(D|S_t, \theta) q_\pi(S_t, D)$
- $= \pi(U|S_t, \theta) \nabla_\theta \log \pi(U|S_t, \theta) q_\pi(S_t, U) + \pi(D|S_t, \theta) \nabla_\theta \log \pi(D|S_t, \theta) q_\pi(S_t, D)$
- $= \mathbb{E}_\pi \nabla_\theta \log \pi(a|S_t, \theta) q_\pi(S_t, a)$
- $= \mathbb{E}_\pi \nabla_\theta \log \pi(a|S_t, \theta) G_t$
- Suppose  $q_\pi(S_t, U) > 0, q_\pi(S_t, D) < 0$ . We know  $\log \pi(a|S_t, \theta) \geq 0$ .
  - For the UP action in state  $S_t$ , the policy update  $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi(U|S_t, \theta) q_\pi(S_t, U)$  will push up  $\pi(U|S_t, \theta)$ , i.e, make it more likely to move UP in state  $S_t$
  - For the DOWN action in state  $S_t$ , the policy update  $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi(D|S_t, \theta) q_\pi(S_t, D)$  will push down  $\pi(D|S_t, \theta)$ , i.e, make it less likely to move DOWN in state  $S_t$

## MC REINFORCE

$$\nabla J(\theta) = \sum_s \mu_\pi(s) \sum_a \nabla \pi(a|s, \theta) q_\pi(s, a)$$

$$= \mathbb{E}_\pi [\sum_a \nabla \pi(a|S_t, \theta) q_\pi(S_t, a)]$$

– Outer exp over policy  $\pi$ : average over all experienced states under policy  $\pi$ , i.e.,  $S_t \sim_s$

$$= \mathbb{E}_\pi \left[ \sum_a \pi(a|S_t, \theta) \frac{\nabla \pi(a|S_t, \theta)}{\pi(a|S_t, \theta)} q_\pi(S_t, a) \right]$$

$$= \mathbb{E}_\pi [\sum_a \pi(a|S_t, \theta) \nabla \log \pi(a|S_t, \theta) q_\pi(S_t, a)] \text{ (since } \nabla \log f(x) = \frac{\nabla f(x)}{f(x)})$$

$$= \mathbb{E}_\pi [\mathbb{E}_\pi \nabla \log \pi(A_t|S_t, \theta) q_\pi(S_t, A_t)]$$

– Inner exp over policy  $\pi$ : average over all experienced actions  $A_t$  from state  $S_t$  under policy  $\pi$ , i.e.,  $A_t \sim \pi(a|S_t, \theta)$

$$= \mathbb{E}_\pi [\nabla \log \pi(A_t|S_t, \theta) q_\pi(S_t, A_t)]$$

– Outer and inner exp over policy  $\pi$  combined: execute policy  $\pi$  and average over all experienced states  $S_t$  and actions  $A_t$  from state  $S_t$ , i.e.,  $S_t \sim_s, A_t \sim \pi(a|S_t, \theta)$

$$= \mathbb{E}_\pi [\nabla \log \pi(A_t|S_t, \theta) G_t]$$

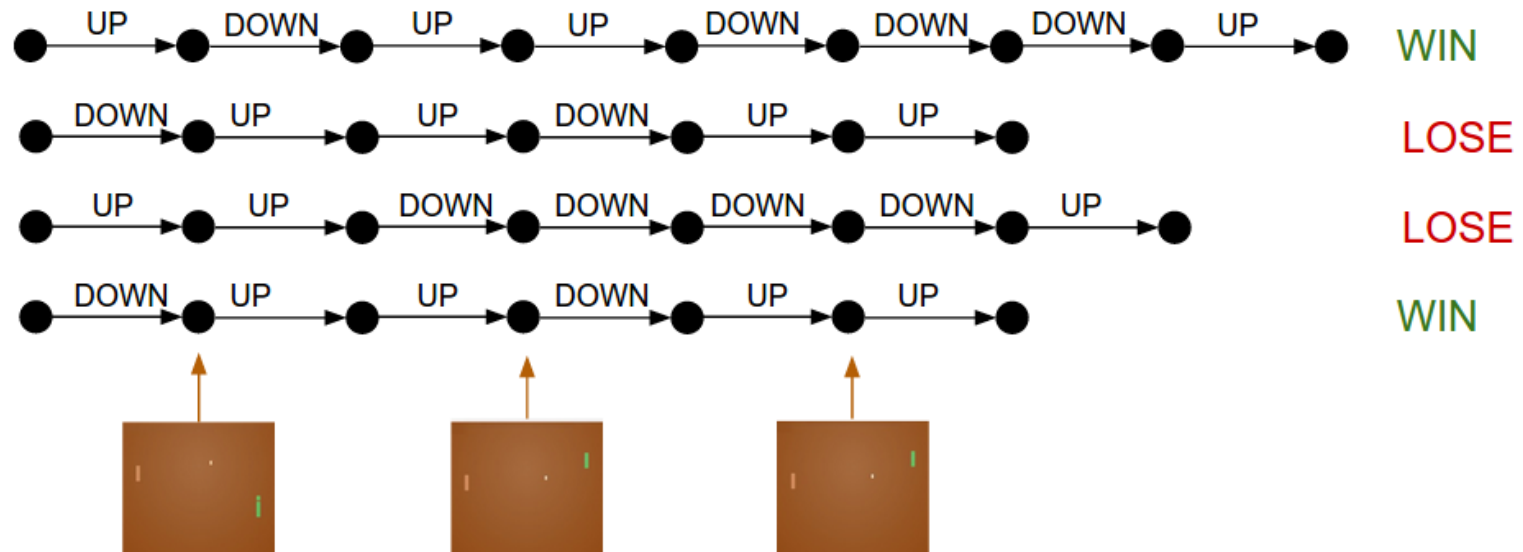
– One way to approximate  $q_\pi(S_t, A_t)$  is by taking expectation over the return  $G_t$  in every episode for taking action  $A_t$  in state  $S_t$  ( $\mathbb{E}_\pi [G_t|S_t, A_t] = q_\pi(S_t, A_t)$ )

SGD update ( $\gamma = 1$ ):  $\theta_{t+1} \leftarrow \theta_t + \alpha \nabla \log \pi(A_t|S_t, \theta) G_t$

SGD update ( $\gamma < 1$ ):  $\theta_{t+1} \leftarrow \theta_t + \alpha \gamma^t \nabla \log \pi(A_t|S_t, \theta) G_t$  (proof omitted)

# MC REINFORCE in Action

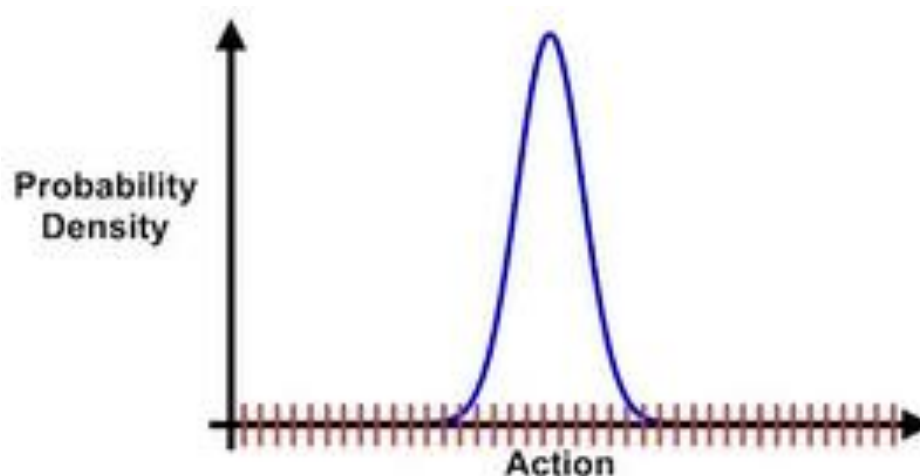
- Agent plays 4 rollouts (episodes), and won 2 episodes and lost 2. Assume that each episode lasts 200 steps, so agent made 200 decisions of UP or DOWN in each episode. We can compute  $G_t$  at timestep  $t$  for each episode with MC Policy Evaluation (recall L7.2 Value-based RL)
- For each of the 2 **winning** episodes ( $v_\pi(S_t) > 0$ ), NN params w.  $\theta \leftarrow \theta + \alpha \nabla \log \pi(A_t | S_t, \theta) q_\pi(S, A_t)$  are updated to **encourage** all taken actions in the 200 steps (**push up**  $\pi(A_t | S_t, \theta)$ ).
- For each of the 2 **losing** episodes ( $v_\pi(S_t) < 0$ ), NN params w.  $\theta \leftarrow \theta + \alpha \nabla \log \pi(A_t | S_t, \theta) q_\pi(S, A_t)$  are updated to **discourage** all taken actions in the 200 steps (**push down**  $\pi(A_t | S_t, \theta)$ ).
- The NN will now become slightly more likely to repeat actions that worked, and slightly less likely to repeat actions that didn't work.





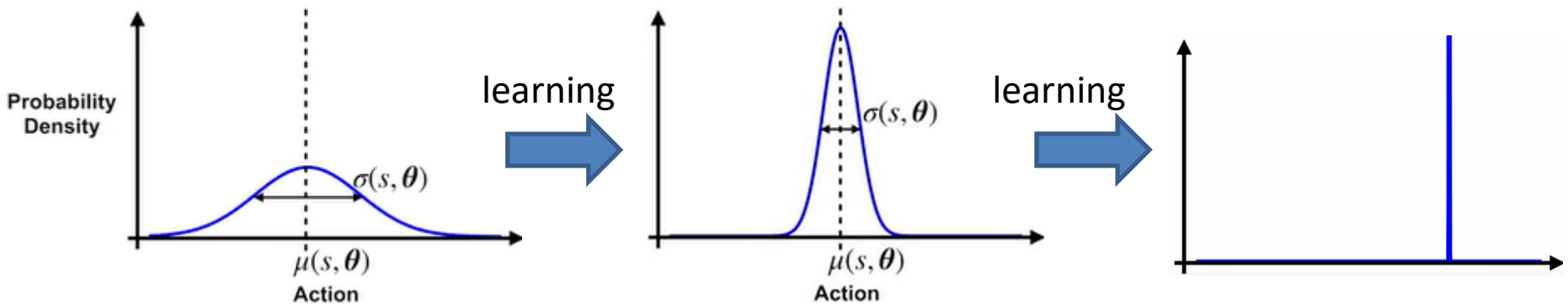
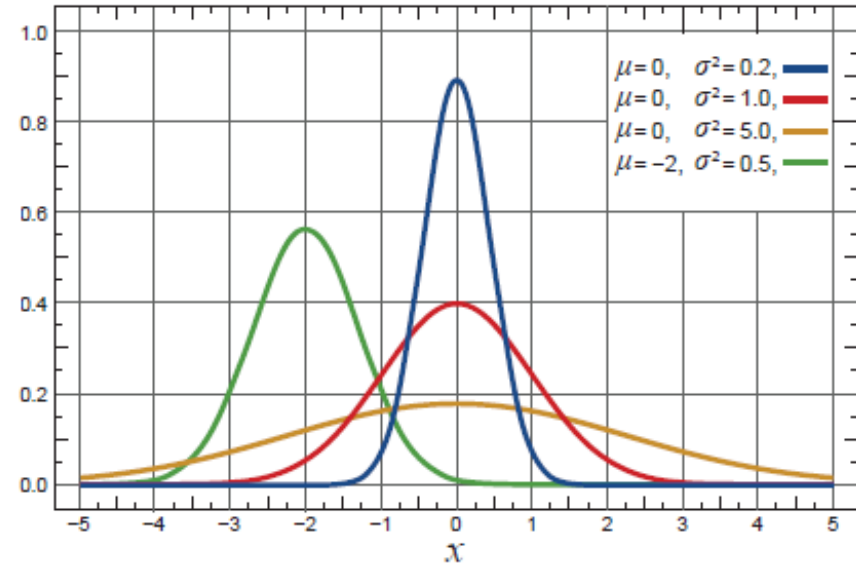
# Advantages of Continuous Actions

- It might not be straightforward to choose a proper discrete set of actions
- Continuous actions allow us to generalize over actions
  - If an action is good, its neighboring actions are also likely to be good
  - Discrete actions lack generalization: each action is independent of others, including its neighbors (similar to value functions for discrete states)



# Gaussian Policy for Continuous Actions

- Gaussian Policy  $\pi(a|s, \theta) \doteq \frac{1}{\sigma(s, \theta)\sqrt{2\pi}} \exp\left(-\frac{(a-\mu(s, \theta))^2}{2\sigma(s, \theta)^2}\right)$ 
  - Mean  $\mu(s, \theta)$  is the most likely action
  - Variance  $\sigma(s, \theta)^2$  controls the degree of exploration.



Policy variance initially large,  
more exploration

Variance gradually reduced during learning w. PG,  
converging towards deterministic policy  $a = \mu(s, \theta)$