

# 6.7900 Machine Learning (Fall 2023)

**Lecture 17:  
Reinforcement Learning Cont'd  
(supporting slides)  
Shen Shen**

# Outline

- Value-based RL
  - (Tabular) Q-learning
- Policy-based RL
  - What does the policy gradient do?
  - Policy gradient derivation
  - Policy gradient estimates
  - Variance reduction
    - Constant Baselines
    - Temporal structure
    - Actor-critic intro

# References

- More RL-flavored presentation:
  - Reinforcement Learning: An Introduction, Sutton and Barto; The MIT Press, 2018.  
Chapter 6 and 13
- Seminal papers referenced on slides.
- Some slides adapted from: Philip Isola, Pieter Abbeel, and Andrej Karpathy

# Reinforcement Learning

Unknown Model

Multi-Armed Bandits

Reinforcement Learning

Known Model

Stochastic Optimization

Markov Decision Process

Actions Don't Impact State

Actions Change State

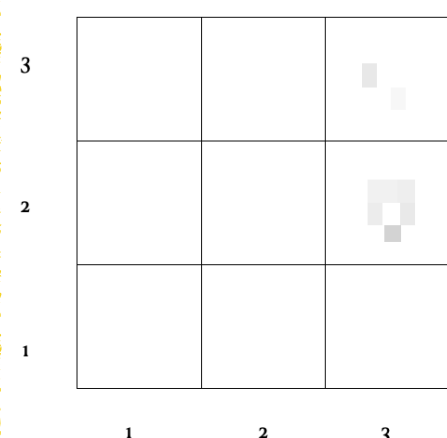
## Markov Decision Process RL

- $\mathcal{S}$ : a state space which contains all possible states  $s$  of the system.
- $\mathcal{A}$ : an action space which contains all possible actions  $a$  an agent can take.
- $P(s'|s, a)$ : the probability of transition from state  $s$  to  $s'$  if action  $a$  is taken.
- $R(s, a)$ : a function that takes in the (state and action) and returns a real-valued reward.

Sometimes, also:

- $s_0$ : initial state.
- Objective version (may involve a  $\gamma \in [0,1]$ : discount factor (details later), and/or  $T$ : horizon. Details later).

## Example: Grid World (in RL)



(Almost deterministic) Transitions:

- Normally, actions take us deterministically to the “intended” state. E.g., in state (1,1), action “North” gets us to state (1,2)
- If an action would take us out of this world, stay put
- In state (3,2), action “North” leads to two possible next state:
  - chance ends in (3,3)
  - chance ends in (2,3)

State space: 9 cells

Actions space: {North, South, East, West}

Discount  $\gamma = 0.9$

Deterministic Rewards:

- State (3,3), any action gets reward
- State (3,2), any action gets reward
- Any other (state, action) pairs get reward 0

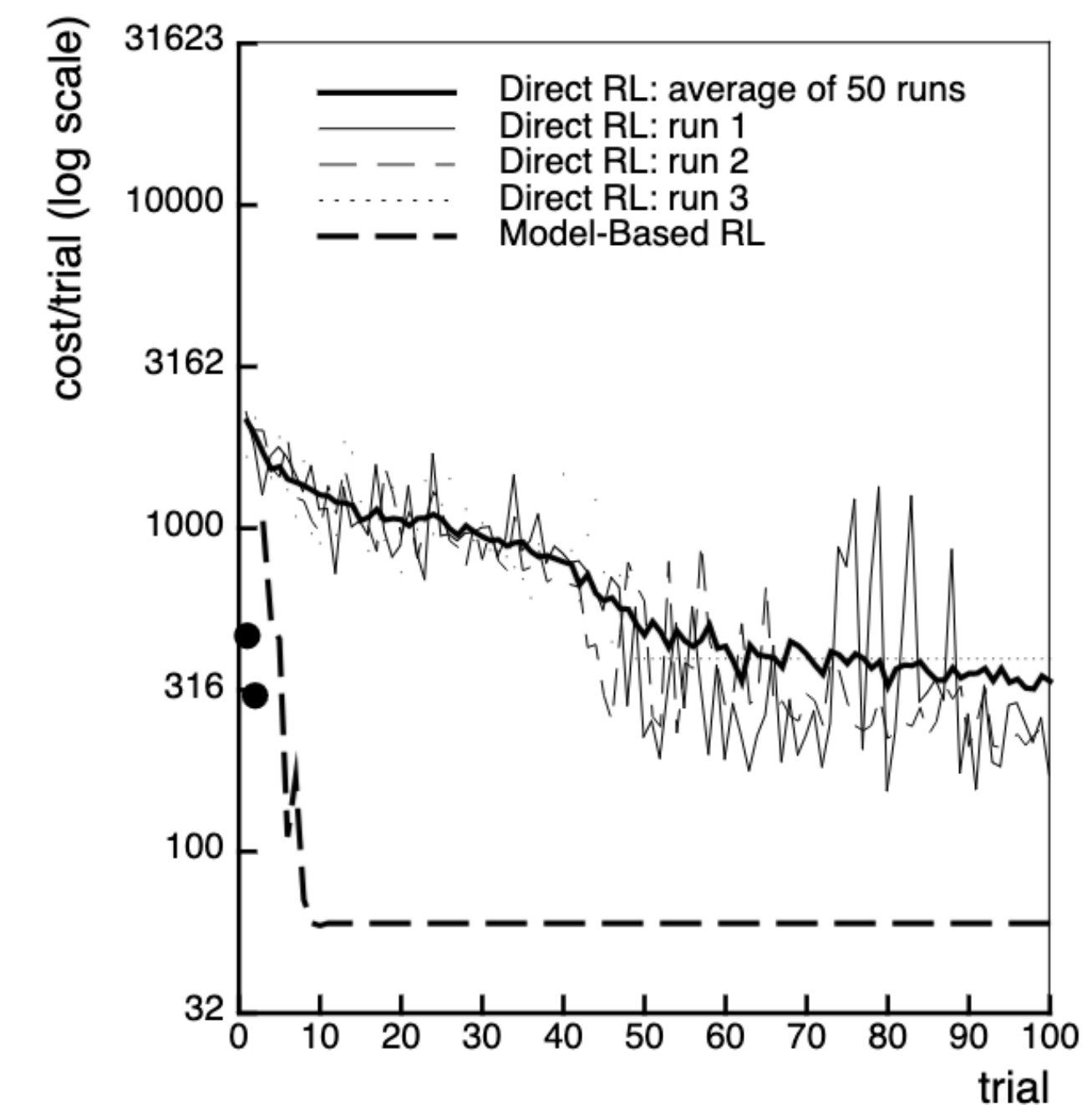
## RL — MDP Goal

- Find a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , such that:

$$\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) | s_0 = s] \text{ is maximized for all } s_0$$

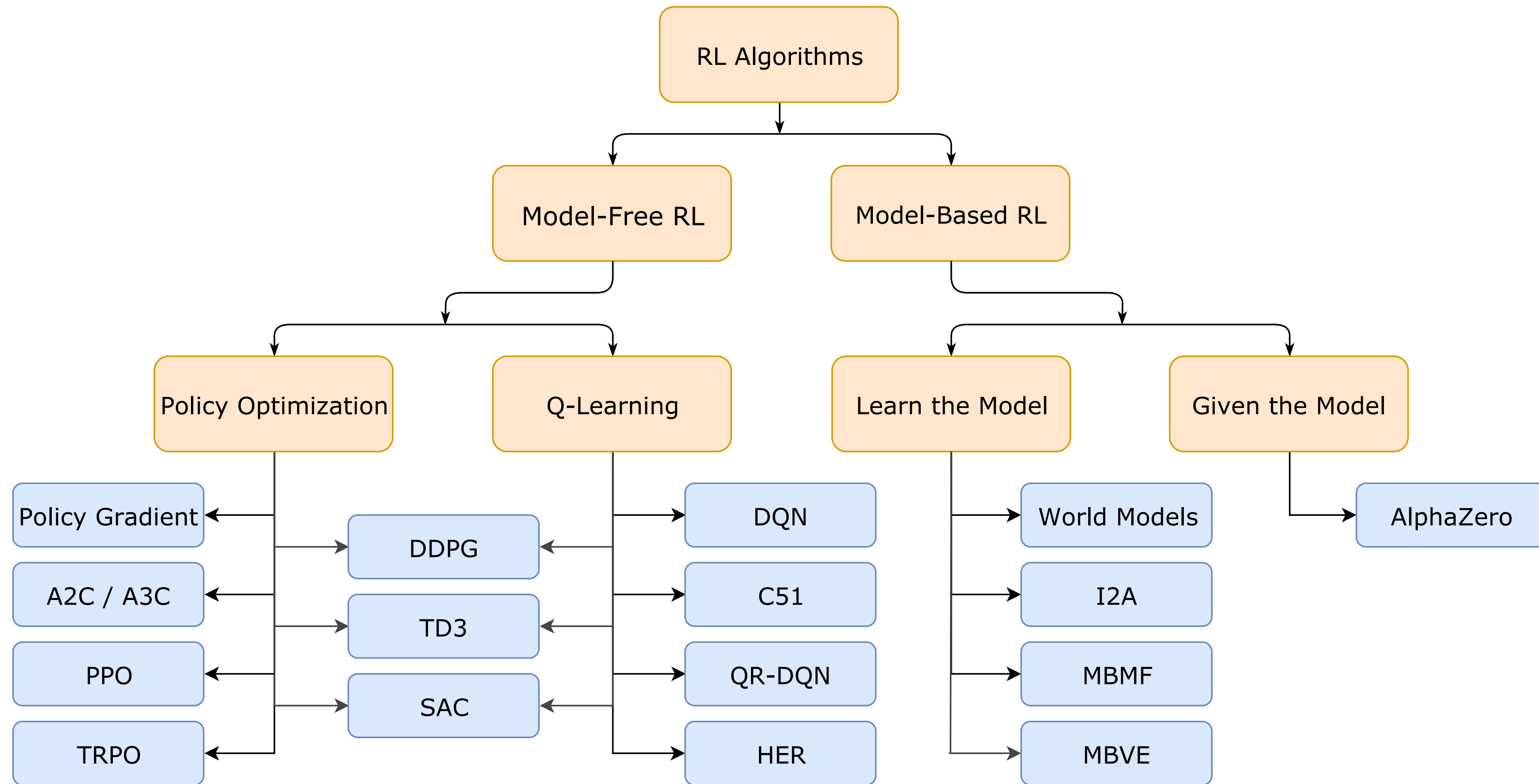
# Model-Based RL

- Collect trajectories to estimate the transition and rewards model (system identification in control)
  - $\hat{P}(s' | s, a) = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=1}^{L_k-1} 1(s_{k,t} = s, a_{k,t} = a, s_{k,t+1} = s')$
  - $\hat{R}(s, a) = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{t=1}^{L_k-1} 1(s_{k,t} = s, a_{k,t} = a) r_{t,k}$
  - Where  $1(\cdot)$  is indicator function and  $N(s, a)$  is the count of trajectories starting from  $(s, a)$
- Then solve the estimated MDP
- Typically more sample efficient than the so-called model-free RL
- Inherit limitations of MDP exact methods, e.g., can be very computationally expensive



[Atkeson and Santamaría, 96]

# A Glance of RL Algorithms



# Q Learning

- Recall that using Q-value iteration, we update our estimate of Q via
$$Q_{\text{new}}(s, a) \leftarrow \mathbb{E}[R(s, a)] + \gamma \sum_{s'} p(s' | s, a) \max_{a'} Q(s', a')$$
- Without access to  $P$  and  $R$ , how would we be able to use this?
- One idea is to sample a state and action pair  $(s, a)$ , simulate, observe  $s', r$ , and then 
$$Q_{\text{new}}(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')$$
- But this is too “current sample dependent” — assumes the observed  $r$  is the only possible reward, assumes the observed  $s'$  is the only possible next state.
- So, instead, “smooth” the update with a step-size  $\alpha$ 
$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a'))$$
 target

# (Tabular) Q Learning

Q-LEARNING( $\mathcal{S}, \mathcal{A}, s_0, \gamma, \alpha, \epsilon$ )

```
1   $Q(s, a) = 0$  for  $s \in \mathcal{S}, a \in \mathcal{A}$ 
2   $s = s_0$  // (e.g.,  $s_0$  can be drawn randomly from  $\mathcal{S}$ )
3  while True:
4       $a = \text{select\_action}(s, Q)$ 
5       $r, s' = \text{execute}(a)$ 
6       $Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a'))$ 
7       $s \leftarrow s'$ 
8      if  $|Q - Q_{\text{new}}| < \epsilon$ : // (or, if reached some max iteration)
9          return  $Q_{\text{new}}$ 
10  $Q \leftarrow Q_{\text{new}}$ 
```

target



# Q-learning Comments

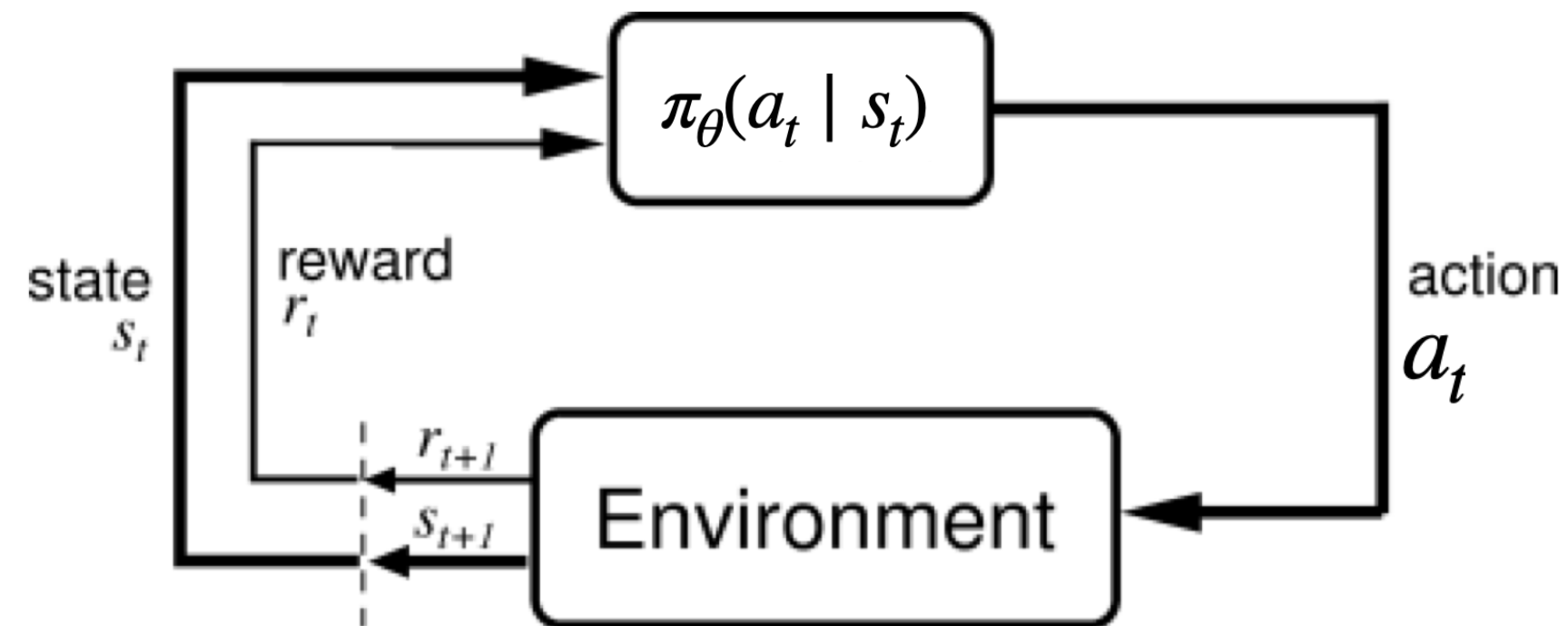
- Face the same exploration versus exploitation dilemma as in bandits (due to unknown model)
- selection\_action in line4 often uses epsilon-greedy; many other options available
- Rearranging terms in line 6, the update can also be interpreted via temporal-difference (TD) error:  $Q_{new}(s, a) \leftarrow Q(s, a) + \alpha(\text{target} - Q(s, a))$  TD error
- In TD-error form, the update looks quite like SGD.
- Closely connects to Fitted Q-learning (coming up in future lecture).

Q-LEARNING( $\mathcal{S}, \mathcal{A}, s_0, \gamma, \alpha, \epsilon$ )

```
1  Q(s, a) = 0 for s ∈ S, a ∈ A
2  s = s_0 // (e.g., s_0 can be drawn randomly from S)
3  while True:
4      a = select_action(s, Q)
5      r, s' = execute(a)
6      Q_new(s, a) ← (1 - α)Q(s, a) + α(r + γ max_{a'} Q(s', a'))
7      s ← s'
8      if |Q - Q_new| < ε: // (or, if reached some max iteration)
9          return Q_new
10 Q ← Q_new
```

- Can converge to true  $Q^*$  if:
  - All states and actions visited infinity often
  - Step-size  $\alpha$  are annealed (i.e. if  $\alpha_k, k$  being the iteration number of line 6, satisfy:  
 $\sum_{k=1}^{\infty} \alpha_k = \infty$  and  $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$ )

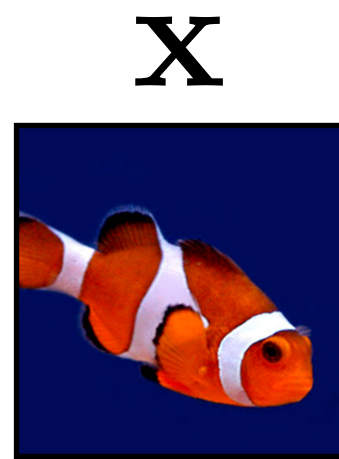
# Policy Optimization



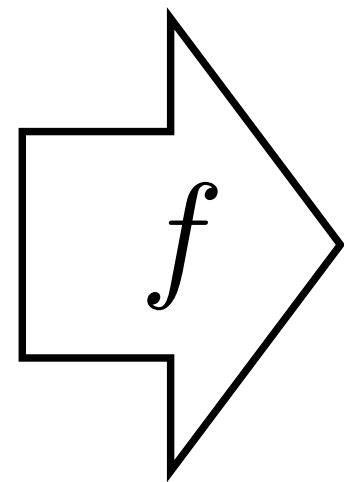
- Parameterize policy by  $\theta$  and directly try  $\max_{\theta} \mathbb{E} \left[ \sum \gamma^t R (s_t, a_t) \mid \pi_{\theta} \right]$
- Stochastic policy class  $\pi_{\theta}(a | s)$  : probability of action  $a$  in state  $s$ 
  - Discrete  $\mathcal{A}$ : e.g.  $\pi_{\theta}(a | s)$  softmax
  - Continuous  $\mathcal{A}$ : e.g.  $\pi_{\theta}(a | s)$  Gaussian with mean/variance parameterized by  $\theta$
  - Smooths out the optimization problem
  - Also encourages exploration

# Why Policy Optimization

- Often  $\pi$  can be simpler than  $Q$  or  $V$ 
  - e.g. lots of  $\pi$  are roughly good
- $V(s)$ : doesn't prescribe actions
  - $\pi^*(s) = \arg \max_a \left[ \mathbb{E}[R(s, a)] + \gamma \sum_{s'} p(s' | s, a) V^*(s') \right]$
  - Would still need world model (and compute one-step Bellman update)
- $Q$ : need to be able to efficiently solve  $\arg \max_a Q(s, a)$ 
  - $\pi^*(s) = \arg \max_a Q^*(s, a)$
  - Can be challenging for continuous / high-dimensional action spaces
- Maybe makes sense to directly optimize policy end-to-end
- So how do we do this?

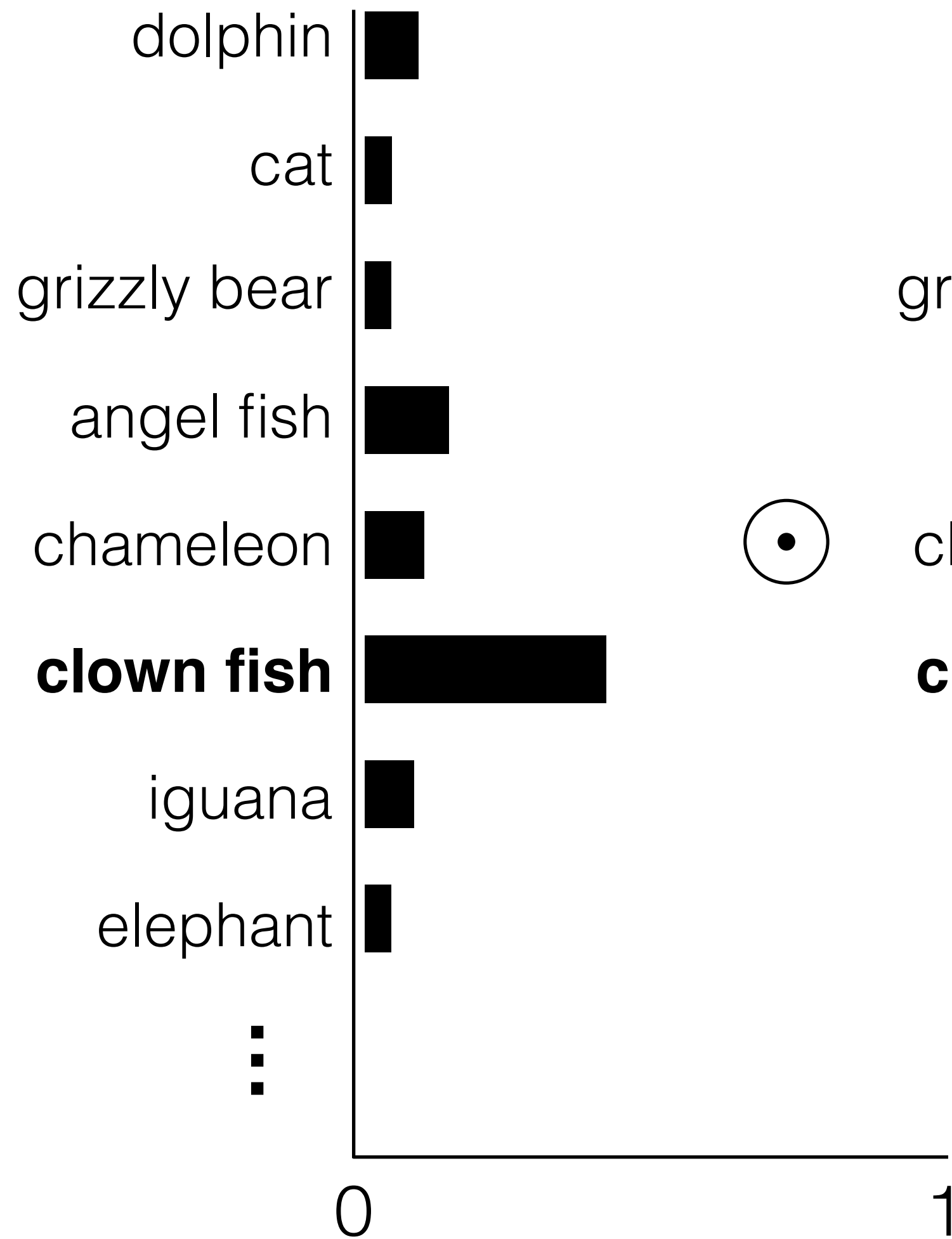


**X**

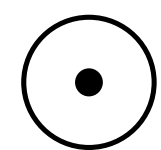


Prediction  $\hat{y}$

$$f_{\theta} : X \rightarrow \mathbb{R}^K$$

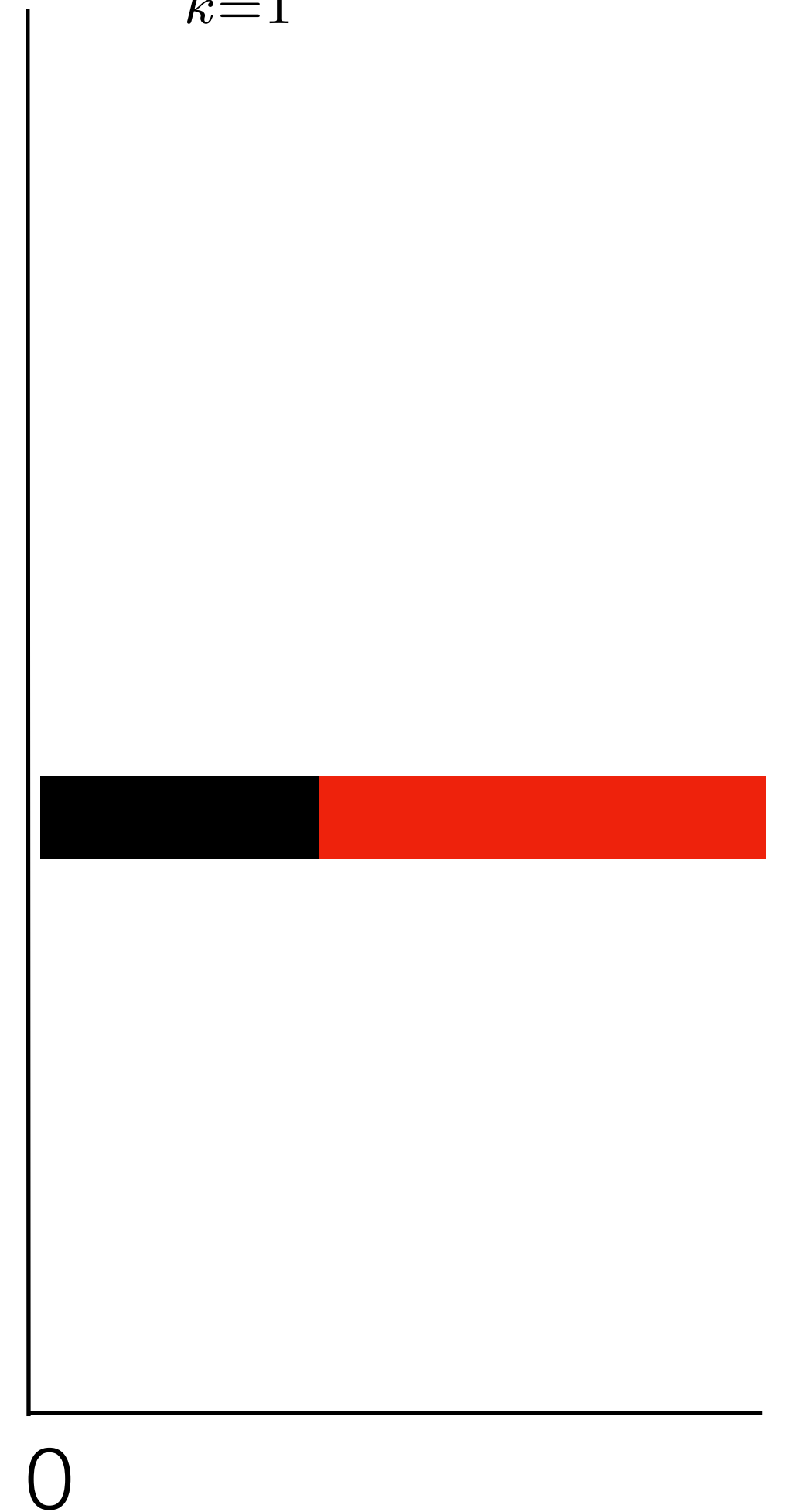


Ground truth label  $y$

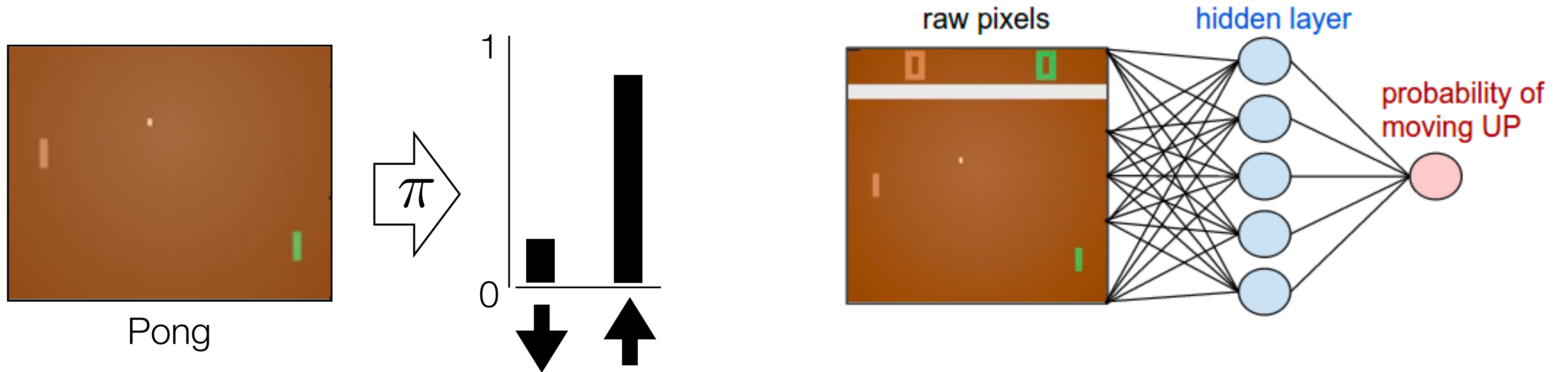


Loss

$$-\sum_{k=1}^K y_k \log \hat{y}_k$$

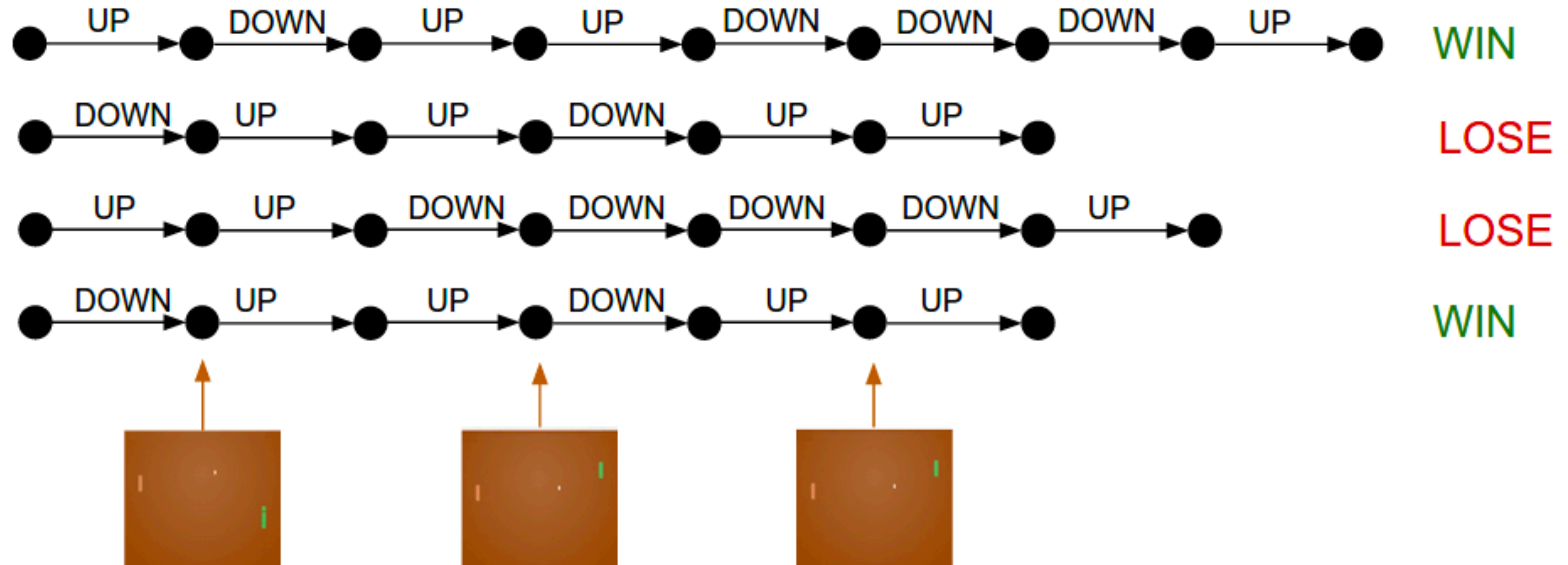


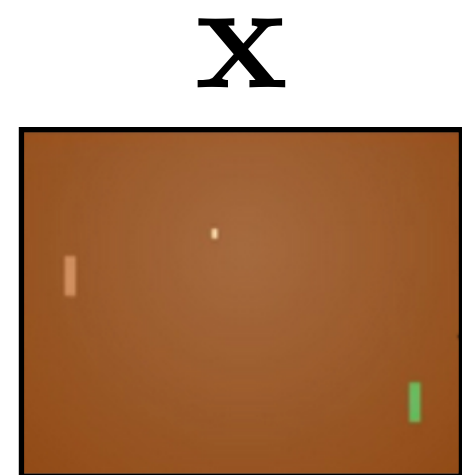
- If explicit “good” state-action pair is given, also supervised learning.
- Behavior cloning or imitation learning.
- But what if no explicit guide?



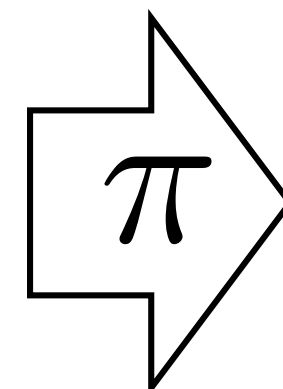
[Adapted from Andrej Karpathy: <http://karpathy.github.io/2016/05/31/r1/>]

**Policy gradients:** Run a policy for a while. See what actions led to good return. Increase their likelihood.





$\mathbf{x}$

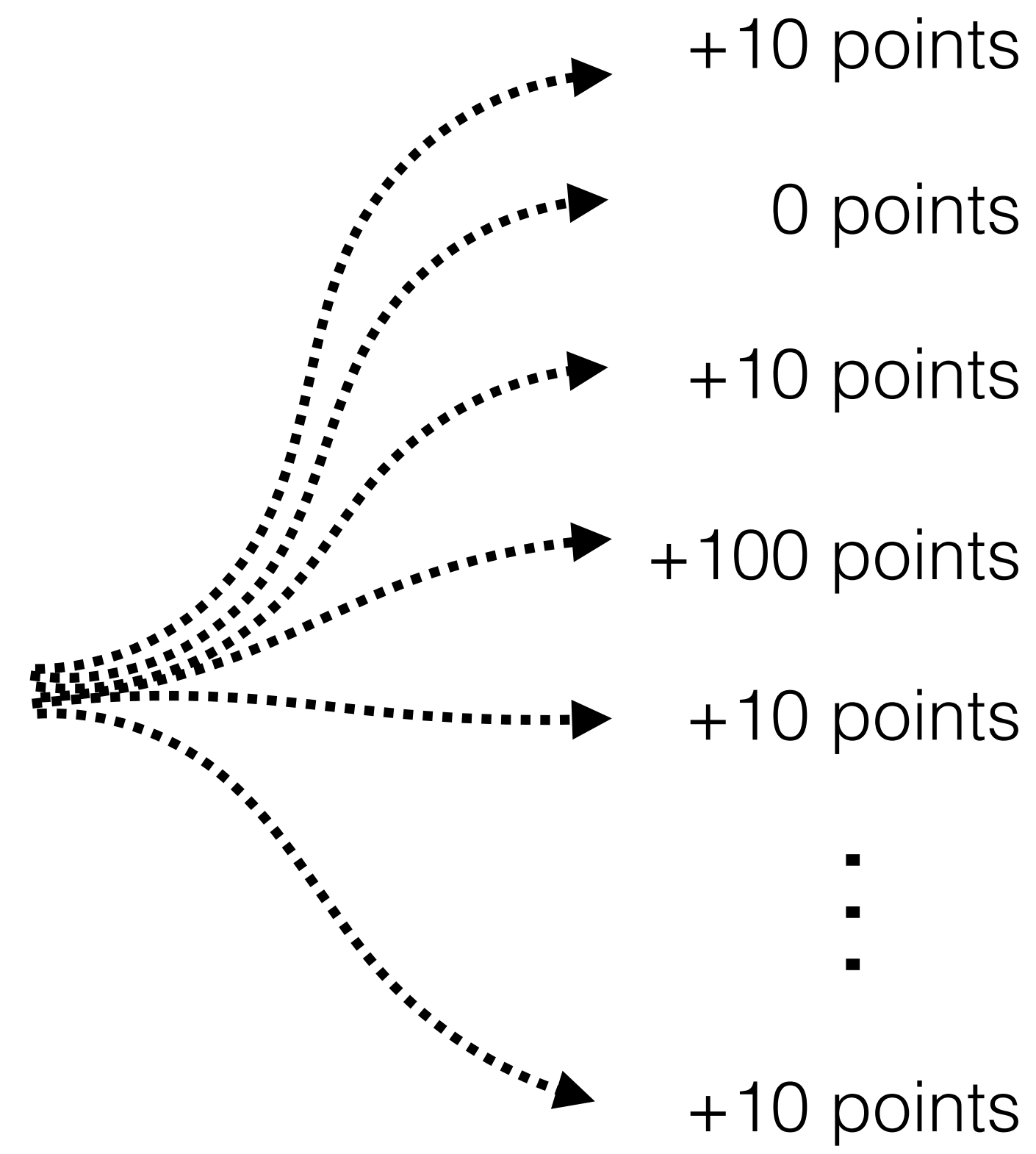


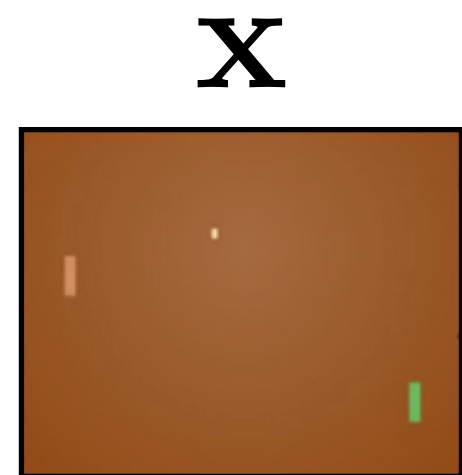
Policy output



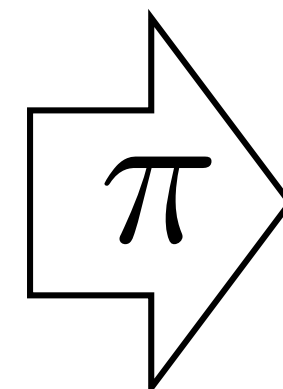
Action

Up





$\mathbf{x}$

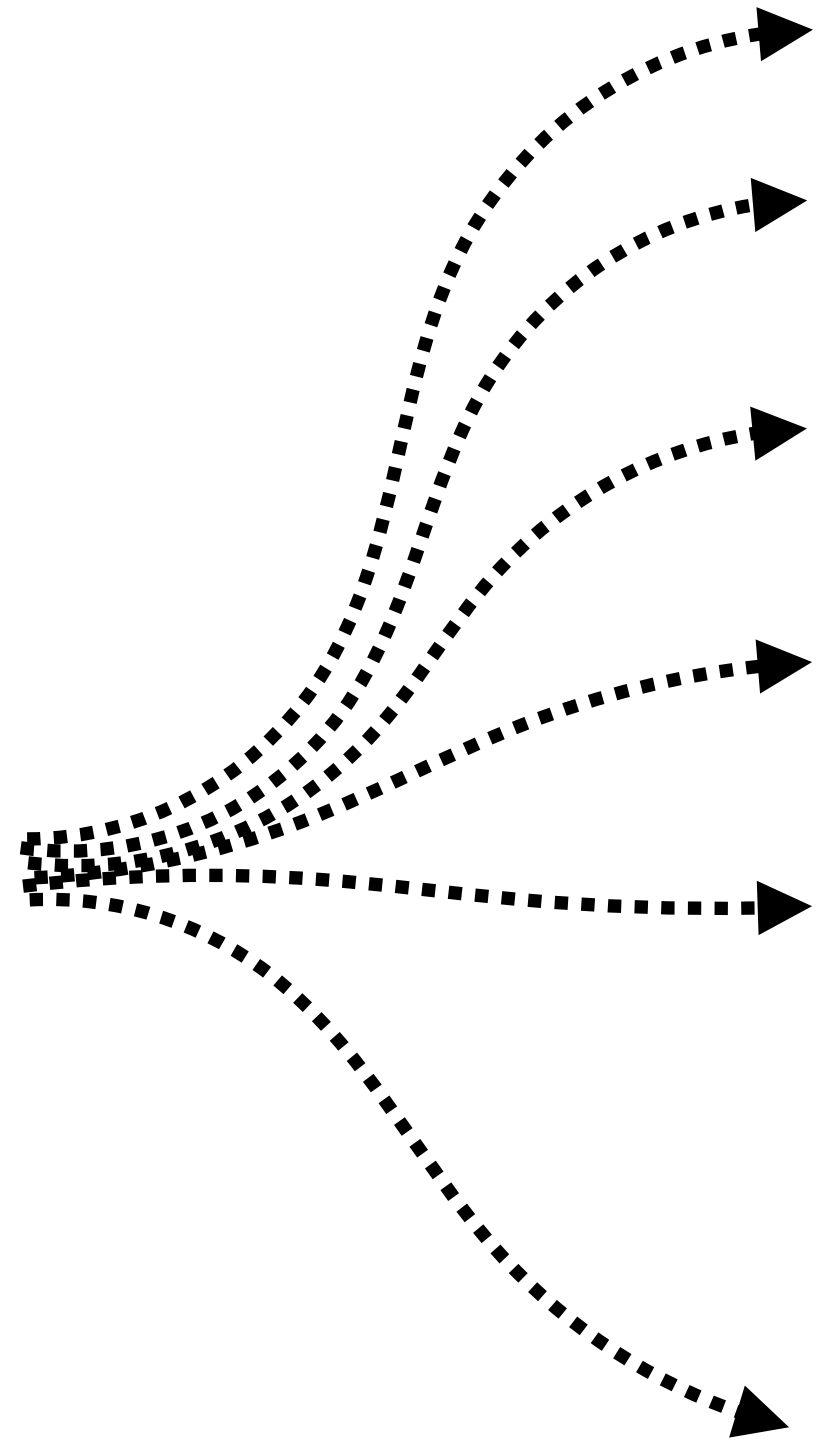


Policy output



Action

Down

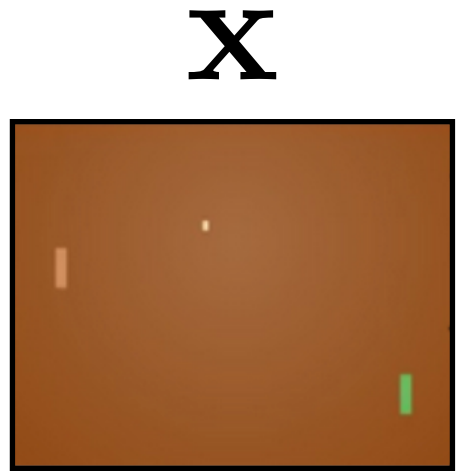


Eventual return

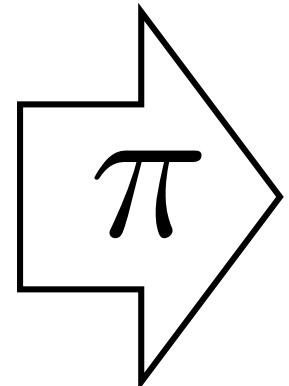
- 0 points
- 0 points
- +10 points
- 0 points
- 0 points
- ⋮
- +10 points



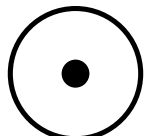
**Approximated via lots of sampling**



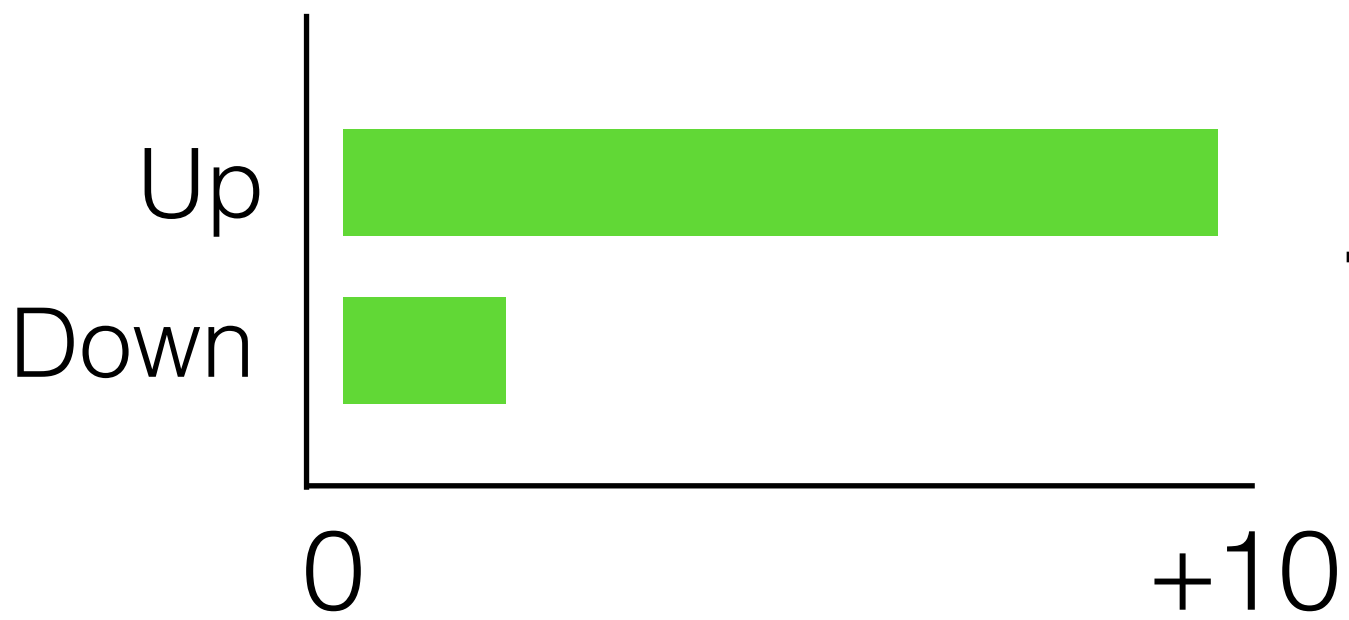
**X**



Policy output



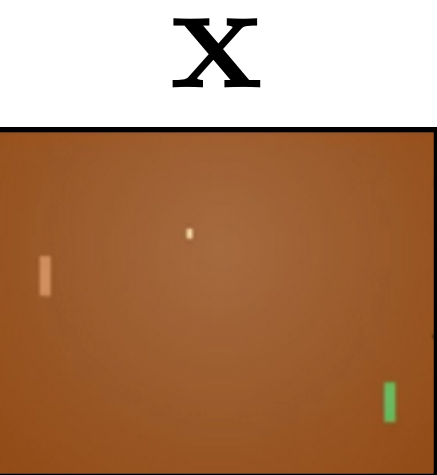
Average return after taking each action



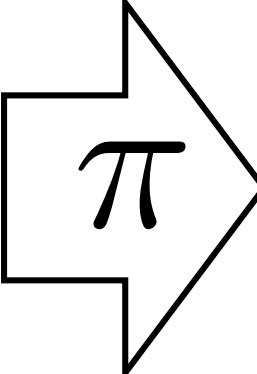
**+6**

Expected return

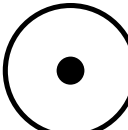
Approximated via lots of sampling



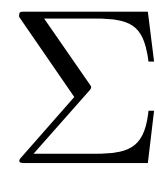
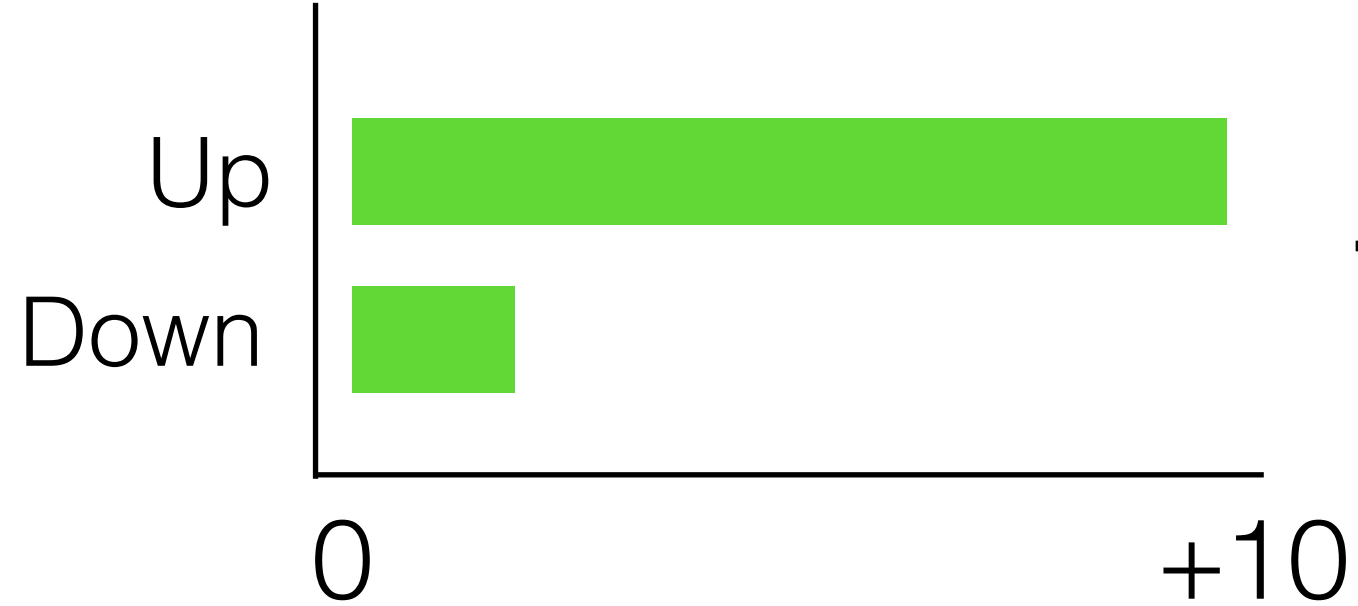
x



Policy output



Average return after taking each action



Expected return

+6

How is this gradient update done though (as we don't have the world model Pong)?

# Likelihood Ratio Policy Gradient

▸ We overload notation:

- Let  $\tau$  denote a state-action sequence:  $\tau = s_0, a_0, s_1, a_1, \dots$

- Let  $R(\tau)$  denote the sum of discounted rewards on  $\tau$ :  $R(\tau) = \sum_t \gamma^t R(s_t, a_t)$

- W.l.o.g. assume  $R(\tau)$  is deterministic in  $\tau$

- Let  $P(\tau; \theta)$  denote the probability of trajectory  $\tau$  induced by  $\pi_\theta$

- Let  $U(\theta)$  denote the objective:  $U(\theta) = \mathbb{E}[\sum_t \gamma^t R(s_t, a_t) \mid \pi_\theta]$

▸ Our goal is to find  $\theta$ :  $\max_\theta U(\theta) = \max_\theta \sum_\tau P(\tau; \theta) R(\tau)$

# Likelihood Ratio Policy Gradient

Taking the gradient w.r.t.  $\theta$  gives

$$\begin{aligned}\nabla_{\theta} U(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau)\end{aligned}$$

$$U(\theta) = \sum_{\tau} P(\tau; \theta) R(\tau)$$

$$\text{But } P(\tau; \theta) = \prod_{t=0} \underbrace{P(s_{t+1} | s_t, a_t)}_{\text{transition}} \cdot \underbrace{\pi_{\theta}(a_t | s_t)}_{\text{policy}}$$

Use identity

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

# Likelihood Ratio Policy Gradient

$$\nabla_{\theta} U(\theta) = \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau)$$

Approximate with the empirical estimate for  $m$  sample traj. under policy  $\pi_{\theta}$

$$\nabla_{\theta} U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

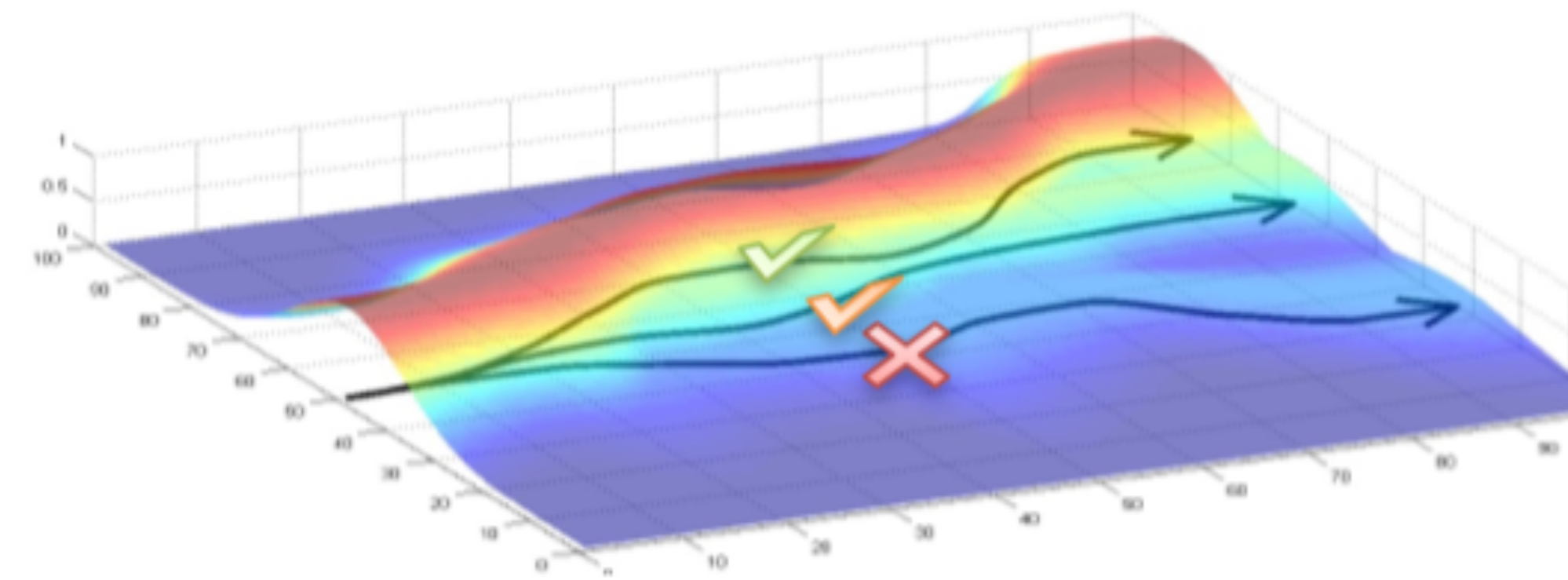
Valid even when:

- Reward function discontinuous and/or unknown
- Discrete state and/or action spaces

# Likelihood Ratio Gradient

$$\nabla_{\theta} U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

- Checks out with our intuition that:
  - Increase likelihood of trajectory with big reward
  - Decrease prob of trajectory with negative reward



- How do we evaluate  $\nabla_{\theta} \log P(\tau^{(i)}; \theta)$  though?

Didn't we say we don't know the transition?

$$P(\tau; \theta) = \prod_{t=0} \underbrace{P(s_{t+1} | s_t, a_t)}_{\text{transition}} \cdot \underbrace{\pi_{\theta}(a_t | s_t)}_{\text{policy}}$$

# Decompose a trajectory

$$\begin{aligned}\nabla_{\theta} \log P(\tau; \theta) &= \nabla_{\theta} \log \left[ \prod_{t=0} \underbrace{P(s_{t+1} | s_t, a_t)}_{\text{transition}} \cdot \underbrace{\pi_{\theta}(a_t | s_t)}_{\text{policy}} \right] \\ &= \nabla_{\theta} \left[ \sum_{t=0} \log P(s_{t+1} | s_t, a_t) + \sum_{t=0} \log \pi_{\theta}(a_t | s_t) \right] \\ &= \nabla_{\theta} \sum_{t=0} \log \pi_{\theta}(a_t | s_t) \\ &= \sum_{t=0} \underbrace{\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)}_{\text{no transition model required,}}\end{aligned}$$

# Likelihood Ratio Gradient - Summary

- The following expression provides us with an unbiased estimate of the gradient, and we can compute it without access to the world model:

$$\hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

- Here

$$\nabla_{\theta} \log P(\tau; \theta) = \sum_{t=0} \underbrace{\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)}_{\text{no need of dynamics}}$$

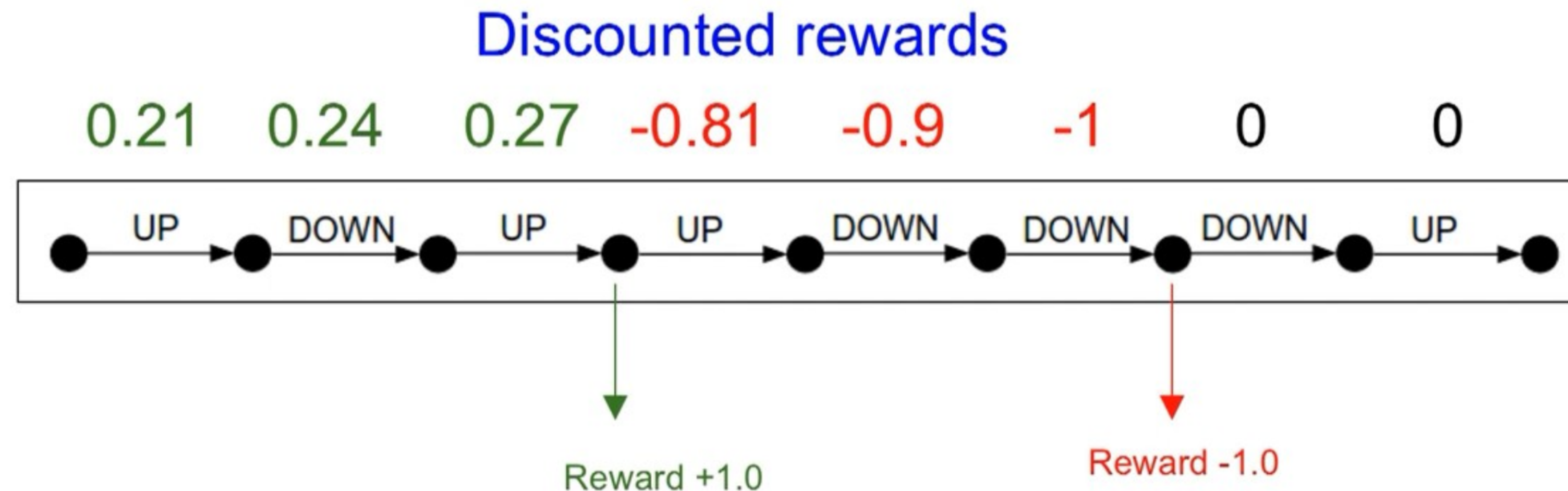
- Unbiased estimator  $E[\hat{g}] = \nabla_{\theta} U(\theta)$ , but very noisy.



# Variance Reduction - Discount

Blame each action assuming that its effects have exponentially decaying impact into the future.

---



- In the extreme, if discount of 0, almost no variance at all.
- So discount can be both a problem definition, or a hyper-parameter.

# Variance Reduction - Baseline

- Sample estimate, unbiased but can be very noisy

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

- Can we keep unbiasedness but reduce variance? Yes!
- Subtract an appropriate baseline can keep the unbiasedness

$$\begin{aligned} & \mathbb{E} [\nabla_{\theta} \log P(\tau; \theta) b] \\ &= \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) b \\ &= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} b \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) b \\ &= \nabla_{\theta} \left( \sum_{\tau} P(\tau) b \right) = b \nabla_{\theta} \left( \sum_{\tau} P(\tau) \right) = b \times 0 \end{aligned}$$
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) \left( R(\tau^{(i)}) - b \right)$$

# Variance-reduction Baselines

▸ Constant  $b = \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$

▸ Optimal constant baseline:  $b = \frac{\sum_i \left( \nabla_{\theta} \log P(\tau^{(i)}; \theta) \right)^2 R(\tau^{(i)})}{\sum_i \left( \nabla_{\theta} \log P(\tau^{(i)}; \theta) \right)^2}$

[Greensmith, Bartlett, Baxter, JMLR 2004 for variance reduction techniques.]

▸ Estimated state-dependent value functions:  $b(s_t) = \hat{V}^{\pi}(s_t)$

- I.e.,  $\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) \left( R(\tau^{(i)}) - \hat{V}^{\pi}(s) \right)$  Advantage

- We'll discuss methods on how to estimate  $\hat{V}^{\pi}(s_t)$  later.

- This kind of "value" baseline very roughly gets us to actor-critic methods.

# Variance Reduction - Temporal Structure

- Current gradient estimate:

$$\begin{aligned}\hat{g} &= \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) (R(\tau^{(i)}) - b) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \right) \left( \sum_{t=0}^{H-1} R(s_t^{(i)}, a_t^{(i)}) - b \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ \left( \sum_{k=0}^{t-1} R(s_k^{(i)}, a_k^{(i)}) \right) + \left( \sum_{k=t}^{H-1} R(s_k^{(i)}, a_k^{(i)}) - b \right) \right] \right)\end{aligned}$$

- Removing terms that don't depend on current action can lower variance:

$$\frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left( \sum_{k=t}^{H-1} R(s_k^{(i)}, a_k^{(i)}) - b(s_t^{(i)}) \right)$$

[Policy Gradient Theorem: Sutton et al 1999;  
GPOMDP: Bartlett & Baxter, 2001; Survey: Peters &  
Schaal, 2006]

# Estimation of $V^\pi$ (coming up later)

- State-dependent expected return:  $b(s_t) = \mathbb{E} [r_t + r_{t+1} + r_{t+2} + \dots + r_{H-1}] = V^\pi(s_t)$ 
  - Increase the prob of action proportionally to how much its returns are better than the expected return under the current policy
- Can't exactly solve for  $V^\pi$ ; again need to estimate. How?
  - Either collect  $\tau_1, \dots, \tau_m$  and regress against empirical return:

$$\phi_{i+1} \leftarrow \arg \min_{\phi} \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{H-1} \left( V_{\phi}^{\pi}(s_t^{(i)}) - \left( \sum_{k=t}^{H-1} R(s_k^{(i)}, u_k^{(i)}) \right) \right)^2$$

- Or similar to fitted Q-learning, do fitted V-learning:

$$\phi_{i+1} \leftarrow \min_{\phi} \sum_{(s,u,s',r)} \left\| r + V_{\phi}^{\pi}(s') - V_{\phi}(s) \right\|_2^2$$

---

**Algorithm 1** “Vanilla” policy gradient algorithm

---

Initialize policy parameter  $\theta$ , baseline  $b$

**for** iteration=1, 2, ... **do**

Collect a set of trajectories by executing the current policy

At each timestep in each trajectory, compute

the *return*  $R_t = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$ , and

the *advantage estimate*  $\hat{A}_t = R_t - b(s_t)$ .

Re-fit the baseline, by minimizing  $\|b(s_t) - R_t\|^2$ ,

summed over all trajectories and timesteps.

Update the policy, using a policy gradient estimate  $\hat{g}$ ,

which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$

**end for**

---

**Thanks!**

**Questions?**