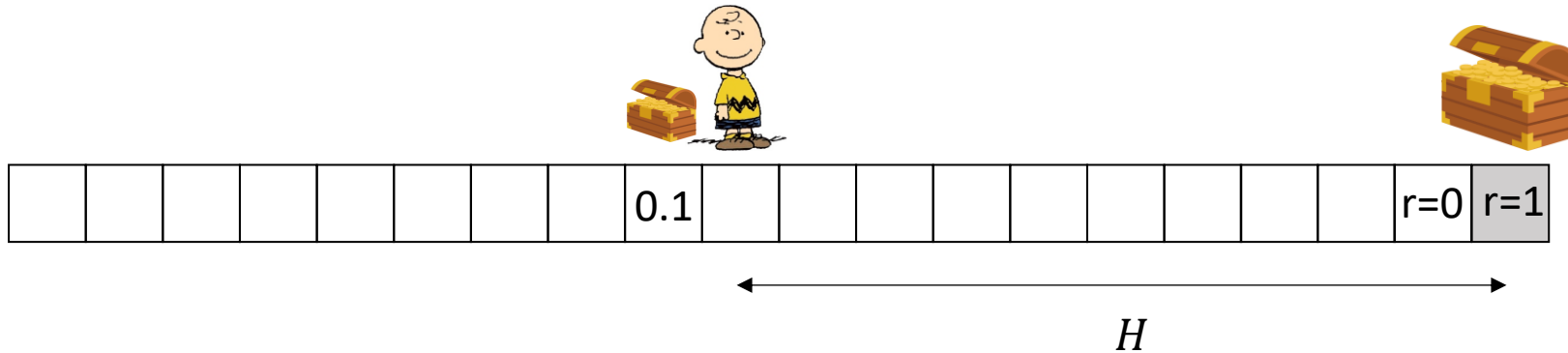


# Exploration in MDPs

Chen-Yu Wei

# State-Space Exploration in MDPs



Environment:

- Fixed-horizon MDP with episode length  $H$
- Initial state at 0
- A single rewarding state at state  $H$
- Actions: Go LEFT or RIGHT

Suppose we perform DQN with  $\epsilon$ -greedy with random initialization

$\Rightarrow$  On average, we need  $2^H$  episodes to see the reward

# Regret Analysis for MDPs?

- We have done regret analysis for several bandit algorithms:

- Regression oracle + ( $\epsilon$ -greedy or inverse gap weighting)
- UCB
- EXP3

$$\sum_{(s,a,s') \in \mathcal{B}} \left( Q_\theta(s,a) - R(s,a) - \sqrt{\frac{1}{t}} \left[ \max_{a'} \sum_{s'} Q_\theta(s',a') \right] \right)^2$$

- We did not really establish regret bounds for MDPs. We only argued:

- Approximate value iteration: under the assumption that the data in replay buffer is exploratory *small*
- Approximate policy iteration: monotonically improvement

$$Q^{n+1}(s,a) \geq Q^n(s,a)$$

$$\max_{s,a} \left| Q(s,a) - R(s,a) - \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ \max_{a'} Q(s',a') \right] \right| \leq \epsilon$$

↓

$$\max_{s,a} |Q(s,a) - Q^*(s,a)| \leq \frac{\epsilon}{1-\gamma}$$

# Regret Analysis for MDPs?

$$\mathbb{E}_{s \sim \rho}[V^{\pi^*}(s)] - \mathbb{E}_{s \sim \rho}[V^{\pi}(s)]$$

$$= \sum_{s,a} d_{\rho}^{\pi}(s,a) (V^*(s) - Q^*(s,a))$$

For VI-based algorithm (approximating  $Q^*$ )

Approximating  $Q^*(s,a)$  requires the replay buffer to cover **wide range of** state-actions.

$$= \sum_{s,a} d_{\rho}^{\pi^*}(s,a) (Q^{\pi}(s,a) - V^{\pi}(s))$$

$$\boxed{(s,a) \sim d_{\rho}^{\pi^*}(s,a)}$$

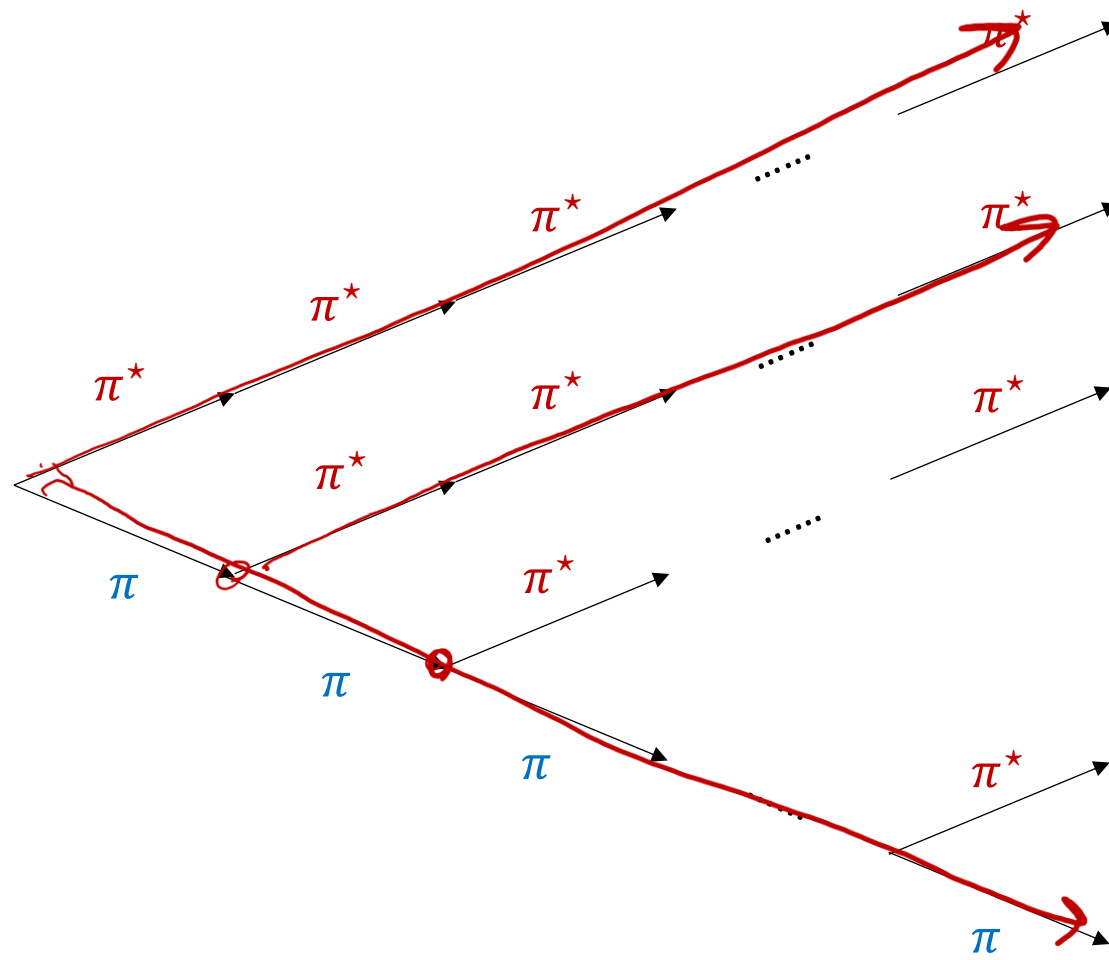
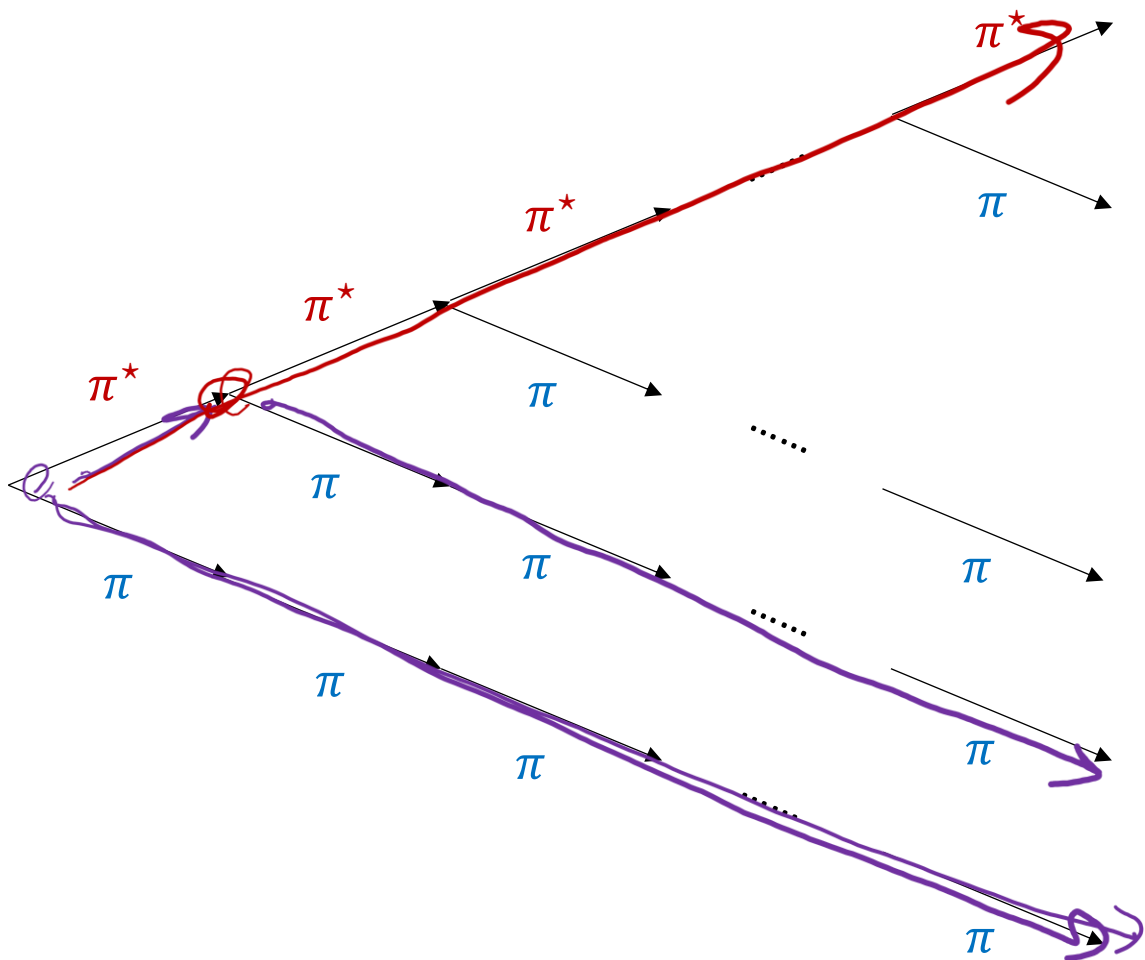
make  $Q^{\pi}(s,a) = \sum_{a'} \pi(a'|s) Q^{\pi}(s,a')$

For PI-based algorithm (approximating  $Q^{\pi}$ )

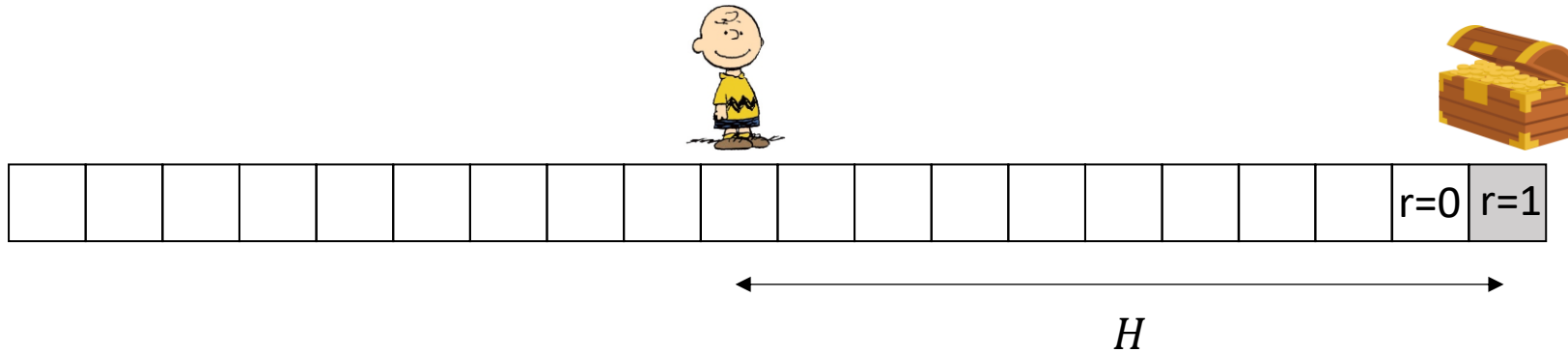
Approximating  $Q^{\pi}(s,a)$  only requires state-actions generated from current policy

But...

$$\sum_{h=1}^H \sum_{s,a} d_{\rho,h}^{\pi^*}(s) (\pi'_h(a|s) - \pi_h(a|s)) Q_h^{\pi}(s,a) = \sum_{h=1}^H \sum_{s,a} d_{\rho,h}^{\pi}(s) (\pi'_h(a|s) - \pi_h(a|s)) Q_h^{\pi^*}(s,a)$$



# Regret Analysis for MDPs?



$$\sum_{s,a} d_{\rho}^{\pi}(s,a) (V^{\star}(s) - Q^{\star}(s,a))$$

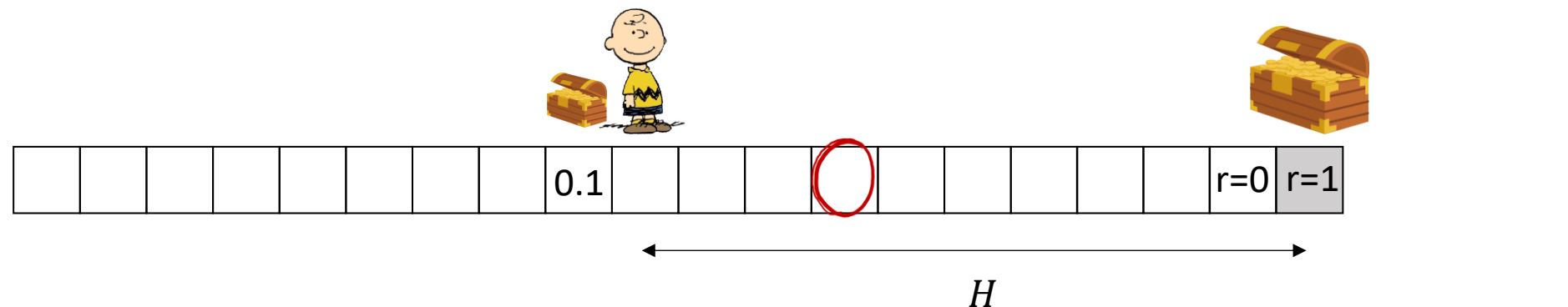
$$\sum_{s,a} d_{\rho}^{\pi^{\star}}(s,a) (Q^{\pi}(s,a) - V^{\pi}(s))$$

PI-based algorithm only tries to make  $\sum_{s,a} d_{\rho}^{\pi^k}(s,a) (Q^{\pi}(s,a) - V^{\pi}(s))$  small.

It can only quickly find optimal policy when  $d_{\rho}^{\pi^k} \approx d_{\rho}^{\pi^{\star}}$

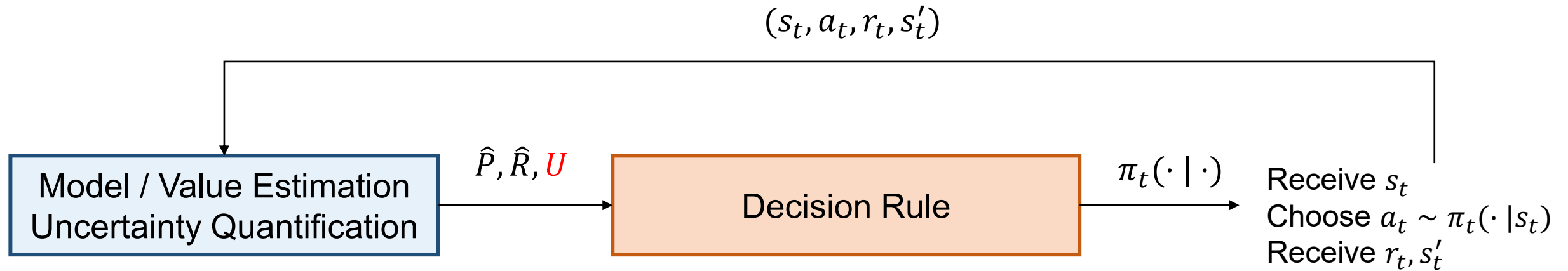
# Insufficiency of algorithms we have discussed for MDPs

- Lack of **exploration over the state space** (we need **deep exploration**)
- This issue is particularly critical if
  - Local reward does not provide any information
  - Local reward provide misleading information



- Solution
  - Try to make the data (i.e., state-action) distribution close to  $d^{\pi^*}$
  - Try to visit as many states as possible (by quantifying the learner's **uncertainty** about a state)

# Exploration via Uncertainty Quantification





# Exploration Bonus for Bandits (Optimism Principle)

- We have discussed this idea for action exploration – UCB.

## Upper Confidence Bound

$$a_t = \operatorname{argmax}_a \hat{R}_t(a) + \sqrt{\frac{2 \log(2/\delta)}{N_t(a)}}$$

$\hat{R}_t(a)$  = the empirical mean of arm  $a$  up to time  $t - 1$ .

$N_t(a)$  = the number of times we draw arm  $a$  up to time  $t - 1$ .

$$\hat{R}_t(a) + \sqrt{\frac{2 \log(2/\delta)}{N_t(a)}} \geq R(a) \quad \text{w.h.p.}$$

# Exploration Bonus for MDPs

## UCB Value Iteration (UCBVI)

For episode  $1, 2, \dots, T$ :

$$\tilde{Q}_{H+1}(s, a) = 0 \quad \forall s, a$$

For step  $H, H-1, \dots, 1$ :

$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + H \sqrt{\frac{2S \log(2/\delta)}{N_t(s, a)}} \quad \forall s, a$$

Receive  $s_1 \sim \rho$

For step  $1, 2, \dots, H$ :

Take action  $a_h = \operatorname{argmax}_a \tilde{Q}_h(s_h, a)$

Receive  $r_h = R(s_h, a_h) + \text{noise}$ ,  $s_{h+1} \sim P(\cdot | s_h, a_h)$

$N_t(s, a, s')$  ← # times we visit  $s, a$  and see next state  $s'$

$N_t(s, a)$  ← # times we visit  $s, a$

$\hat{R}(a) + \sqrt{\frac{1}{N_t(a)}}$

$\tilde{Q}_h(s, a) \geq Q_h^*(s, a)$

$$-1 \leq R(s,a) \leq 1$$

$$|\hat{R}_t(s,a) - R(s,a)| \lesssim \sqrt{\frac{1}{N_t(s,a)}}$$

$$\|\hat{P}(\cdot|s,a) - P(\cdot|s,a)\|_1 \lesssim \sqrt{\frac{S}{N_t(s,a)}}$$

# Exploration Bonus for MDPs

$$\tilde{Q}_h(s,a) \triangleq \hat{R}(s,a) + \sum_{s'} \hat{P}(s'|s,a) \max_{a'} \tilde{Q}_{h+1}(s',a') + H \sqrt{\frac{2S \log(2/\delta)}{N_t(s,a)}} \quad \forall s,a$$

$$Q_h^*(s,a) = R(s,a) + \sum_{s'} P(s'|s,a) \max_{a'} Q_{h+1}^*(s',a')$$

$$\tilde{Q}_h(s,a) \geq Q_h^*(s,a)$$

$$\tilde{Q}_h(s,a) - Q_h^*(s,a)$$

$$+ b_t(s,a)$$

$$= \hat{R}(s,a) - R(s,a) + \sum_{s'} \left( \hat{P}(s'|s,a) \max_{a'} \tilde{Q}_{h+1}(s',a') - P(s'|s,a) \max_{a'} Q_{h+1}^*(s',a') \right)$$

$$\geq \hat{R}(s,a) - R(s,a) + \sum_{s'} \left( \hat{P}(s'|s,a) \max_{a'} Q_{h+1}^*(s',a') - P(s'|s,a) \max_{a'} Q_{h+1}^*(s',a') \right) + b_t(s,a)$$

This holds as long as

$$|\hat{R}(s,a) - R(s,a)|$$

$$+ \sum_{s'} |\hat{P}(s'|s,a) - P(s'|s,a)| \max_{a'} Q_{h+1}^*(s',a') \leq b_t(s,a)$$

$$\geq 0$$

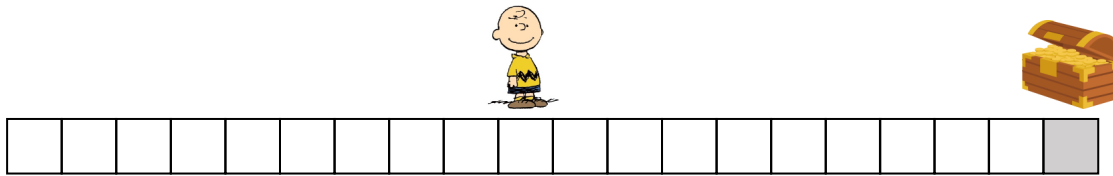
# Exploration Bonus for MDPs

## Theorem. Regret Bound of UCBVI

Proven in HW4

UCBVI ensures with high probability,

$$\text{Regret} = \sum_{t=1}^T (V^*(s_{t,1}) - V^{\pi_t}(s_{t,1})) \lesssim HS\sqrt{AT}.$$



Improving the required number of episodes from  $2^H$  to  $\text{poly}(H)$

Jaksch, Ortner, Auer. Near-Optimal Regret Bounds for Reinforcement Learning. 2010.

Azar, Osband, Munos. Minimax Regret Bounds for Reinforcement Learning. 2017.

# Thompson Sampling (Posterior Sampling)

## Bayesian interpretation:

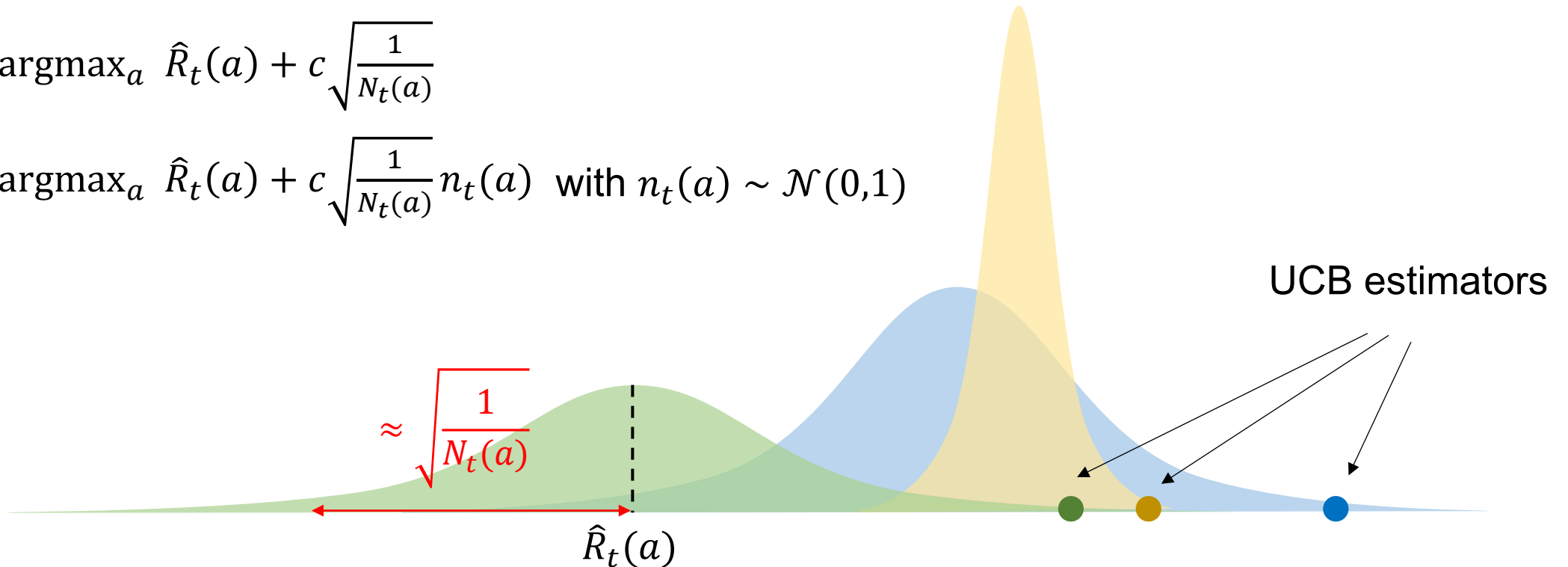
Assume the reward mean  $(\theta(1), \dots, \theta(A))$  is drawn from a Gaussian distribution (prior distribution).

Then the **posterior distribution** is

$$P(\theta(a)|\mathcal{H}_t) = \mathcal{N}\left(\hat{R}_t(a), \frac{1}{N_t(a)}\right)$$

$$\text{UCB: } a_t \approx \operatorname{argmax}_a \hat{R}_t(a) + c \sqrt{\frac{1}{N_t(a)}}$$

$$\text{TS: } a_t \approx \operatorname{argmax}_a \hat{R}_t(a) + c \sqrt{\frac{1}{N_t(a)}} n_t(a) \text{ with } n_t(a) \sim \mathcal{N}(0,1)$$



# Randomized Exploration for MDPs

## Randomized Value Iteration

For episode  $1, 2, \dots, T$ :

$$\tilde{Q}_{H+1}(s, a) = 0 \quad \forall s, a$$

For step  $H, H - 1, \dots, 1$ :

$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + H \sqrt{\frac{2S \log(2/\delta)}{N_t(s, a)}} \underbrace{n_t(s, a)}_{\sim \mathcal{N}(0,1)}$$

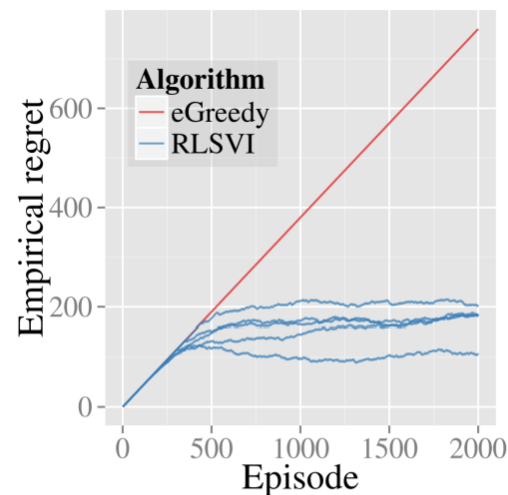
Receive  $s_1 \sim \rho$

For step  $1, 2, \dots, H$ :

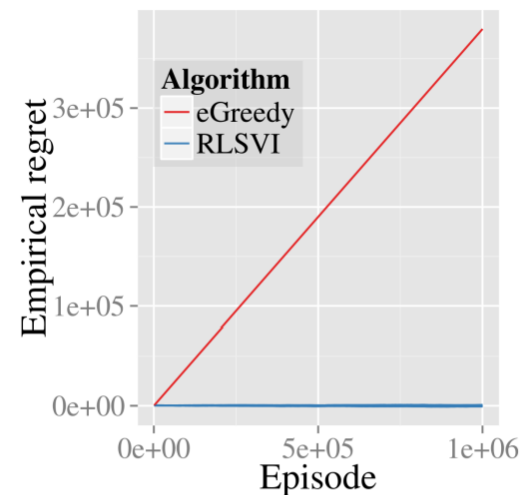
Take action  $a_h = \operatorname{argmax}_a \tilde{Q}_h(s_h, a)$

Receive  $r_h = R(s_h, a_h) + \text{noise}$ ,  $s_{h+1} \sim P(\cdot | s_h, a_h)$

# Randomized Exploration for MDPs



(a) First 2000 episodes



(b) First  $10^6$  episodes

Figure 2. Efficient exploration on a 50-chain

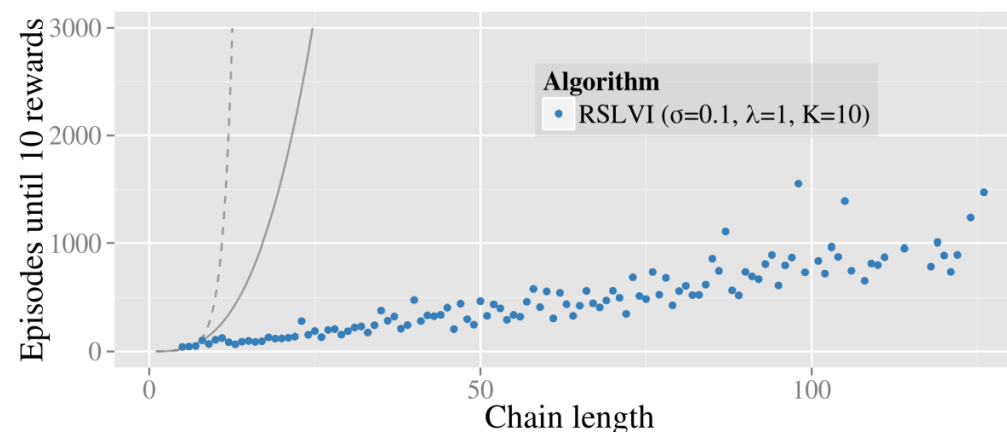
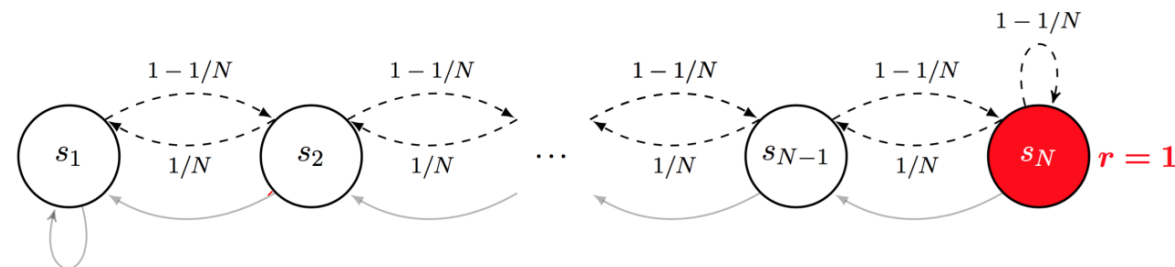


Figure 3. RLSVI learning time against chain length.

# Recap: Exploration in Finite-State Finite-Action MDPs

Find exploration bonus  $B(s, a)$  such that

$$|\hat{R}(s, a) - R(s, a)| \leq B(s, a) \quad \text{reward uncertainty}$$

$$|\mathbb{E}_{s' \sim \hat{P}(\cdot|s, a)}[V(s')] - \mathbb{E}_{s' \sim P(\cdot|s, a)}[V(s')]| \leq B(s, a) \quad \text{transition uncertainty}$$

$\uparrow$   
 $V^*(s)$

$$\mathbb{E}_{a_i \sim \pi_{\text{old}}(\cdot|s_i)} \left[ \frac{\pi_{\theta}(a_i|s_i)}{\pi_{\text{old}}(a_i|s_i)} B(s_i, a_i) \right]$$

Then perform VI (e.g. DQN) over the reward  $r(s, a) + \alpha B(s, a)$

or PI (e.g. PPO, PG) with the **reward estimator**

$$\frac{\pi_{\theta}(a_i|s_i)}{\pi_{\theta_{\text{old}}}(a_i|s_i)} (Q^{\pi_{\text{old}}}(s_i, a_i) - b(s_i)) + \alpha \sum_a \pi_{\theta}(a|s_i) B(s_i, a)$$

$+ \alpha B(s_i, a_i)$

$$\sum_a \pi_{\text{old}}(a|s_i) \cdot \frac{\pi_{\theta}(a|s_i)}{\pi_{\text{old}}(a|s_i)} B(s_i, a)$$

Formalized in HW4

$$= \sum_a \pi_{\theta}(a|s_i) B(s_i, a)$$



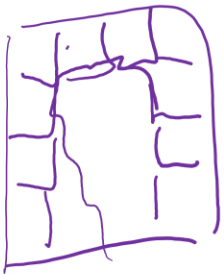
# Common Approaches of Exploration

- Optimistic Exploration
  - Upper Confidence Bound
- Randomized Exploration
  - Thompson Sampling (Posterior Sampling)
- Information-Directed Exploration
  - Information-Directed Sampling

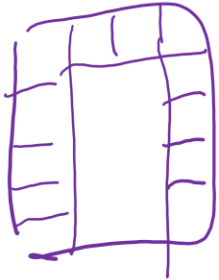
$$|\hat{R}(a) - R(a)|$$

$$0.02 \times \boxed{\phantom{0.02}} \leq \frac{1}{\text{gap}} = \frac{1}{0.02}$$

$$O(1)$$



① Navigation



$$\begin{matrix} 0.51 & 0.49 \\ \circ & \circ \\ \text{Ber}(0.51) & \text{Ber}(0.49) \end{matrix} \quad \circ \quad \textcircled{0}$$

World 1

$$\begin{matrix} 0.49 & 0.51 \\ \circ & \circ \\ \text{Ber}(0.49) & \text{Ber}(0.51) \end{matrix} \quad \circ \quad \textcircled{0.01}$$

World 2

# **Exploration in Large State Spaces**

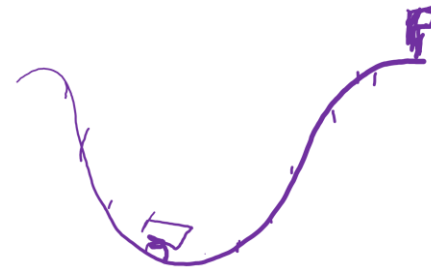
# UCB / TS with State(-Action) Discretization

HW4 Task

Partition the state-action space into a finite number of groups

Then instead of counting the #visits to individual state-action, we only count the #visits to each group

$g(s, a)$ : the group  $(s, a)$  belongs to



$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + c \cdot \frac{1}{\sqrt{N_t(g(s, a))}}$$

$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + c \cdot \frac{\mathcal{N}(0, 1)}{\sqrt{N_t(g(s, a))}}$$

# UCB / TS with State(-Action) Features

$$\phi(s,a) = e_{g(s,a)} = \begin{pmatrix} 0 \\ \vdots \\ 1 \end{pmatrix}$$

$$\phi(s,a) = e_{s,a} \in \mathbb{R}^{S \times A}$$

$$\phi(s,a)^T \Lambda^{-1} \phi(s,a) = \frac{1}{N(s,a)}$$

Suppose for any  $(s, a)$ , we have access to a feature vector  $\phi(s, a) \in \mathbb{R}^d$ .

Then instead of counting the #visits to every state-action, we can evaluate the **novelty of the feature**.

$$\Lambda_t = \sum_{i < t} \sum_{h=1}^H \phi(s_{ih}, a_{ih}) \phi(s_{ih}, a_{ih})^T$$

$$\Lambda = \sum_{(s,a) \in \mathcal{D}} \phi(s,a) \phi(s,a)^T$$

$$\phi(s,a)^T \Lambda^{-1} \phi(s,a)$$

$\frac{1}{N(g(s,a))}$

$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + c \cdot \sqrt{\phi(s, a) \Lambda_t^{-1} \phi(s, a)}$$

Jin et al. Provably efficient reinforcement learning with linear function approximation. 2019.

$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + c \cdot \mathcal{N}(0, \phi(s, a) \Lambda_t^{-1} \phi(s, a))$$

Zanette et al. Frequentist Regret Bounds for Randomized Least-Squares Value Iteration. 2019.

# Ideas for Exploration

Ideas from UCB:

1.  $\tilde{R}(s, a) = \hat{R}(s, a) + \frac{1}{\sqrt{N(s, a)}}$  where  $N(s, a) \approx$  Amount of prior visit to  $(s, a)$
2.  $\tilde{R}(s, a) = \hat{R}(s, a) + \underline{e(s, a)}$  where  $e(s, a) \approx$  Prediction error on  $\hat{R}(s, a)$  and  $\hat{P}(\cdot | s, a)$

$$|\hat{R}(s, a) - R(s, a)| \approx \frac{1}{\sqrt{N(s, a)}}$$

Ideas from TS:

3.  $\tilde{R}(s, a) = \hat{R}(s, a) +$  noise whose variance scales with the uncertainty of  $\hat{R}(s, a)$  and  $\hat{P}(\cdot | s, a)$

Ideas from Information-directed Sampling:

4.  $\tilde{R}(s, a) = \hat{R}(s, a) + \lambda \underbrace{\text{KL}(\mathcal{P}(\cdot | \mathcal{H}_t, s, a, s'), \mathcal{P}(\cdot | \mathcal{H}_t))}_{\text{Information gain}}$

After these modifications, just perform standard RL algorithm over  $\tilde{R}$ .

# **1. Bonus from Prediction Error**

# Bonus from Prediction Error

Ideally, we would like to quantify  $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1$  and set it as bonus.

However, modeling transition is not always easy. Sometimes, what we can do is just **predicting the next state** and measure  $\|\hat{s}'(s, a) - s'(s, a)\|$

There are some issues if we naively do this:

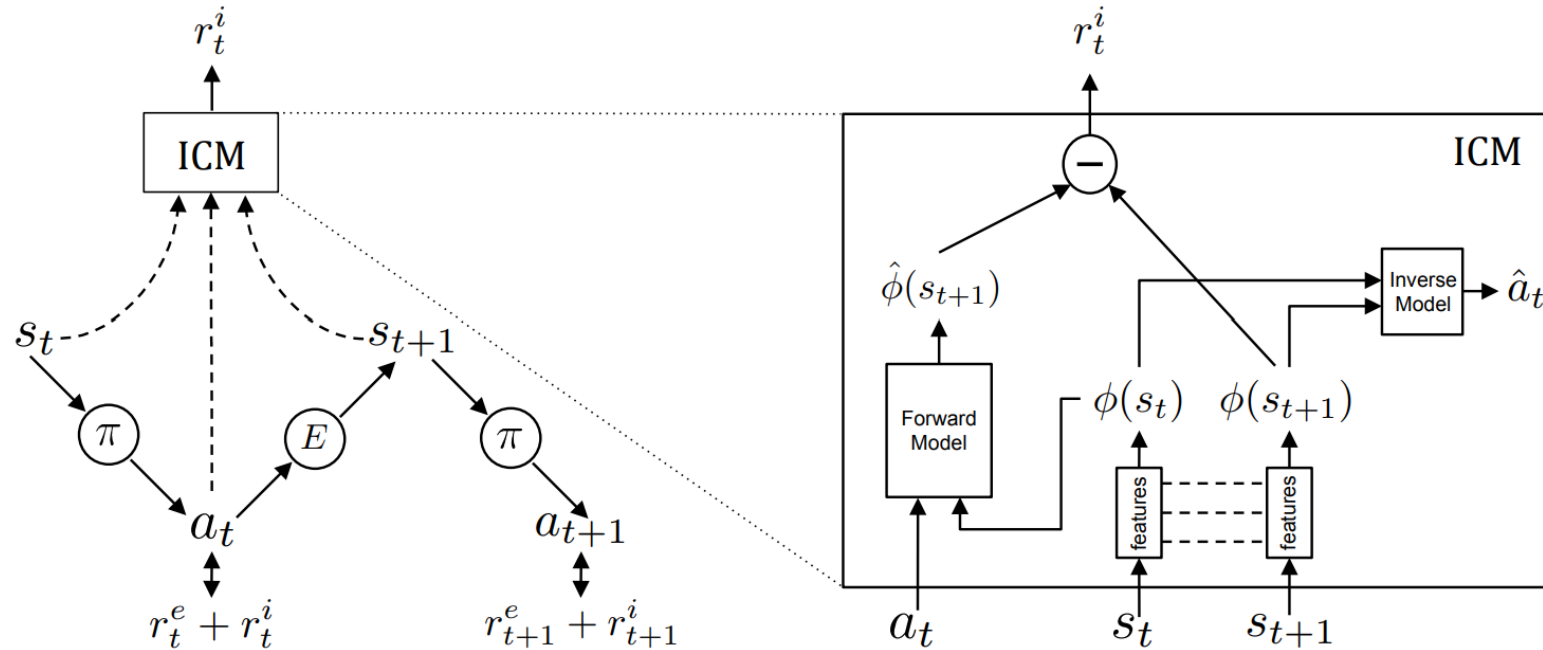
1. For environments with stochastic transitions, we will never have small prediction error for the next state.
2. For many environments, some part of the state is uncontrollable by the learner (e.g., movement of the clouds in the background).

# Bonus from Prediction Error

In some special cases, there exists a deterministic latent-state MDPs.



# Intrinsic Curiosity Module (ICM)



Task 1: Given  $s_t$  and  $s_{t+1}$ , predict  $a_t$ : make the feature  $\phi$  capture action-related dynamics

Task 2: Given  $\phi(s_t)$  and  $a_t$ , predict  $\phi(s_{t+1})$

# Random Network Distillation (RND)

HW4 Task

let's say we have some **target** function  $f^*(\mathbf{s}, \mathbf{a})$

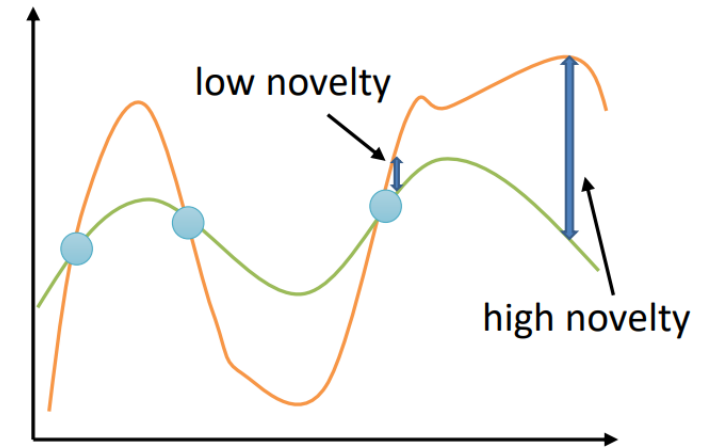
given our buffer  $\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i)\}$ , fit  $\hat{f}_\theta(\mathbf{s}, \mathbf{a})$

use  $\mathcal{E}(\mathbf{s}, \mathbf{a}) = \|\hat{f}_\theta(\mathbf{s}, \mathbf{a}) - f^*(\mathbf{s}, \mathbf{a})\|^2$  as bonus

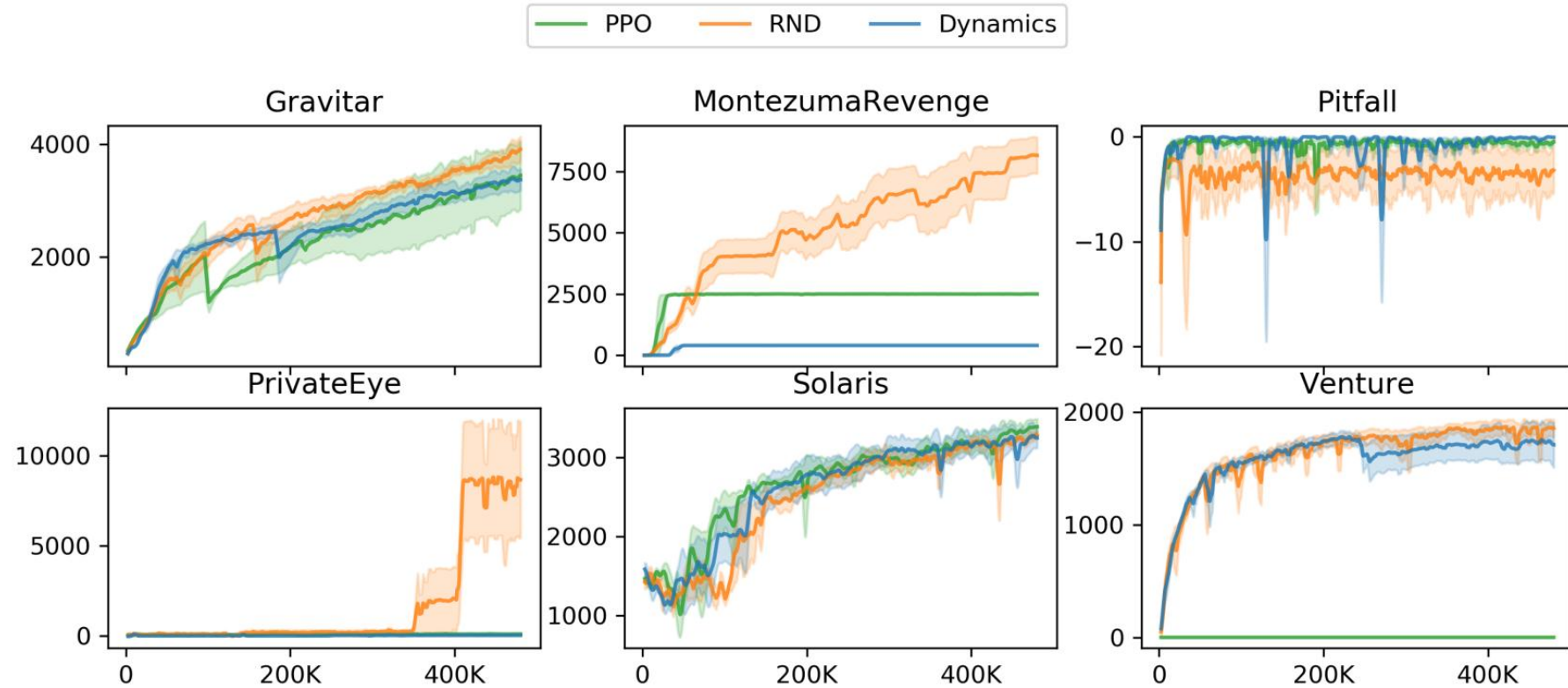
what should we use for  $f^*(\mathbf{s}, \mathbf{a})$ ?

one common choice: set  $f^*(\mathbf{s}, \mathbf{a}) = \mathbf{s}'$  – i.e., next state prediction

even simpler:  $f^*(\mathbf{s}, \mathbf{a}) = f_\phi(\mathbf{s}, \mathbf{a})$ , where  $\phi$  is a *random* parameter vector



# Random Network Distillation



## **2. Thompson Sampling**

# Recall: Randomized Value Iteration

## Randomized Value Iteration

For episode  $1, 2, \dots, T$ :

$$\tilde{Q}_{H+1}(s, a) = 0 \quad \forall s, a$$

For step  $H, H - 1, \dots, 1$ :

$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + H \sqrt{\frac{2 \log(2/\delta)}{N_t(s, a)}} \underbrace{n_t(s, a)}_{\sim \mathcal{N}(0,1)}$$

Receive  $s_1 \sim \rho$

For step  $1, 2, \dots, H$ :

Take action  $a_h = \operatorname{argmax}_a \tilde{Q}_h(s_h, a)$

Receive  $r_h = R(s_h, a_h) + \text{noise}$ ,  $s_{h+1} \sim P(\cdot | s_h, a_h)$

# Recall: Randomized Value Iteration

$$\tilde{Q}_h(s, a) \triangleq \hat{R}(s, a) + \sum_{s'} \hat{P}(s'|s, a) \max_{a'} \tilde{Q}_{h+1}(s', a') + n_t(s, a)$$

Adapting this idea to DQN:

$$\theta = \operatorname{argmin}_{\theta} \sum_{(s, a, r, s') \in \mathcal{B}} \left( r + \max_{a'} Q_{\bar{\theta}}(s', a') + n_t(s, a) - Q_{\theta}(s, a) \right)^2 \quad (*)$$

Notice that different noise gives different  $\theta$ .

**Direct generalization from Randomized VI** (not easy to implement)

$\Theta$  = Space of  $\theta$ 's

In each episode, sample a  $\theta \in \Theta$  with the distribution following (\*),  
and execute  $\pi(s) = \operatorname{argmax}_a Q_{\theta}(s, a)$

# Bootstrapped DQN

Osband et al. Deep Exploration via Bootstrapped DQN. 2016.

Osband et al. Randomized Prior Functions for Deep Reinforcement Learning. 2018.

Randomly initialize  $K$  instances of DQN  $\theta_1, \dots, \theta_K$   
(each  $\theta_i$  has their own target network  $\bar{\theta}_i$  and replay buffer  $\mathcal{B}_i$ ).

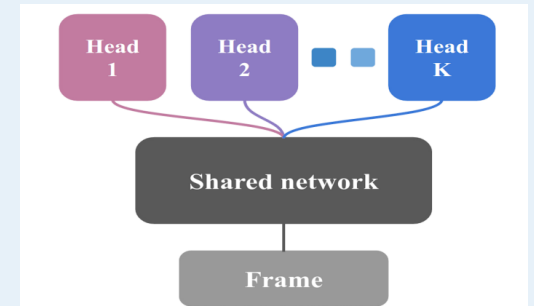
For each episode:

Randomly sample  $i \sim \text{Unif}\{1, 2, \dots, K\}$

Execute  $\pi(s) = \max_a Q_{\theta_i}(s, a)$  in the whole episode.

Randomly place the obtained  $(s, a, r, s')$  in some/all replay buffers.

Update all DQN parameters.



(a) Shared network architecture

# Bootstrapped DQN

Osband et al. Deep Exploration via Bootstrapped DQN. 2016.

Osband et al. Randomized Prior Functions for Deep Reinforcement Learning. 2018.

Some intuitions:

- The random initialization makes  $Q_{\theta_1}(s, a), \dots, Q_{\theta_K}(s, a)$  all very different. We can view them as associated with different initial noise  $n_1(s, a)$ .
- Over the course of training, for  $(s, a)$ 's that are more often visited, their effective magnitude of  $n_t(s, a)$  decreases (because we train those DQNs without adding more noise).
- For  $(s, a)$ 's that are not often visited, their effective magnitude of  $n_t(s, a)$  remains high.
- **Why does this perform deep exploration?** For a particular state  $s$ , if  $\max_a Q_{\theta_i}(s, a)$  is initialized high but has not been visited many times before, the training of  $\theta_i$  will propagate this high value to other state and encourage the learner to reach  $s$  from other states.

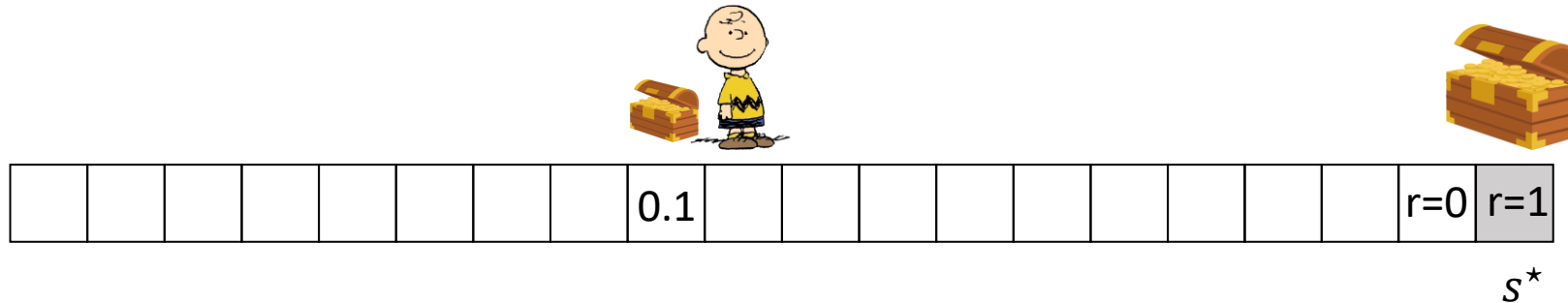


# Bootstrapped DQN

Osband et al. Deep Exploration via Bootstrapped DQN. 2016.

Osband et al. Randomized Prior Functions for Deep Reinforcement Learning. 2018.

- In the toy example, as long as **one of the  $K$  DQNs** initializes  $s^*$  (or some states close to it) with a high value, then it can help the learner explore to  $s^*$ .
- In this example, roughly we need  $K = O(\text{number of states})$  to achieve this effect.



# Bootstrapped DQN

Osband et al. Deep Exploration via Bootstrapped DQN. 2016.

Osband et al. Randomized Prior Functions for Deep Reinforcement Learning. 2018.

## "Deep Sea" Exploration

- Stylized "chain" domain testing "deep exploration":
  - State =  $N \times N$  grid, observations 1-hot.
  - Start in top left cell, fall one row each step.
  - Actions {0, 1} map to left/right in each cell.
  - "left" has reward = 0, "right" has reward =  $-0.1/N$
  - ... but if you make it to bottom right you get +1.
- Only one policy (out of more than  $2^N$ ) positive return.
- $\epsilon$ -greedy / Boltzmann / policy gradient / are useless.

