

# Generalization in Large scale MDPs

Wen Sun

CS 6789: Foundations of Reinforcement Learning

# Recap on Bellman Error and Bellman Operator

Two types of Bellman error of  $f(s, a) (\approx Q^*)$

# Recap on Bellman Error and Bellman Operator

Two types of Bellman error of  $f(s, a)$  ( $\approx Q^*$ )

$$BE_Q(s, a) = f(s, a) - \left( r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

# Recap on Bellman Error and Bellman Operator

Two types of Bellman error of  $f(s, a) (\approx Q^*)$

$$BE_Q(s, a) = f(s, a) - \left( r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

$$V_f(s) = \arg \max_a f(s, a), \pi_f(s) = \arg \max_a f(s, a)$$

# Recap on Bellman Error and Bellman Operator

Two types of Bellman error of  $f(s, a)$  ( $\approx Q^*$ )

$$BE_Q(s, a) = f(s, a) - \left( r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

$$V_f(s) = \arg \max_a f(s, a), \pi_f(s) = \arg \max_a f(s, a)$$

$$BE_V(s) = V_f(s) - r(s, \pi_f(s)) - \mathbb{E}_{s' \sim P_h(s, \pi_f(s))} V_f(s')$$

# Recap on Bellman Error and Bellman Operator

Two types of Bellman error of  $f(s, a)$  ( $\approx Q^*$ )

$$BE_Q(s, a) = f(s, a) - \left( r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

$$V_f(s) = \arg \max_a f(s, a), \pi_f(s) = \arg \max_a f(s, a)$$

$$BE_V(s) = V_f(s) - r(s, \pi_f(s)) - \mathbb{E}_{s' \sim P_h(s, \pi_f(s))} V_f(s')$$

If  $BE(s, a) \neq 0$ , then  $f \neq Q^*$

# Notations

Probability of  $\pi$  visiting  $(s, a)$  at time step  $h$ :  $d_h^\pi(s, a)$

# Question for Today

We have seen tabular MDP and linear MDP, is there a **more general framework** that captures these two, and potentially many more, where efficient learning is possible?



# Question for Today

We have seen tabular MDP and linear MDP, is there a **more general framework** that captures these two, and potentially many more, where efficient learning is possible?

In other words, what structural conditions permit RL generalization, provably?

# Outline for Today

1. Bellman rank Definitions

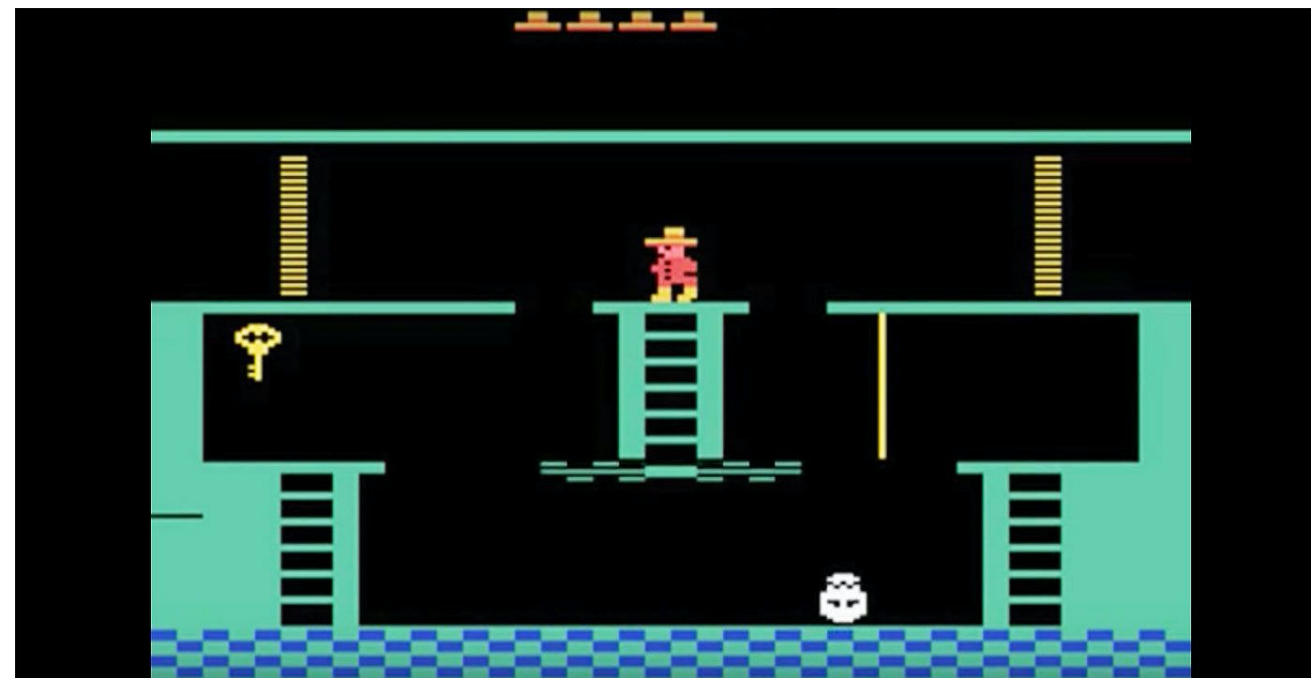
2. Examples that are captured by the Bellman rank framework

# Setting

Finite horizon episodic MDP  $\{ \{S_h\}_{h=0}^H, \{A_h\}_{h=0}^{H-1}, H, s_0, r, P \}$

State space  $S_h$  is extremely large:

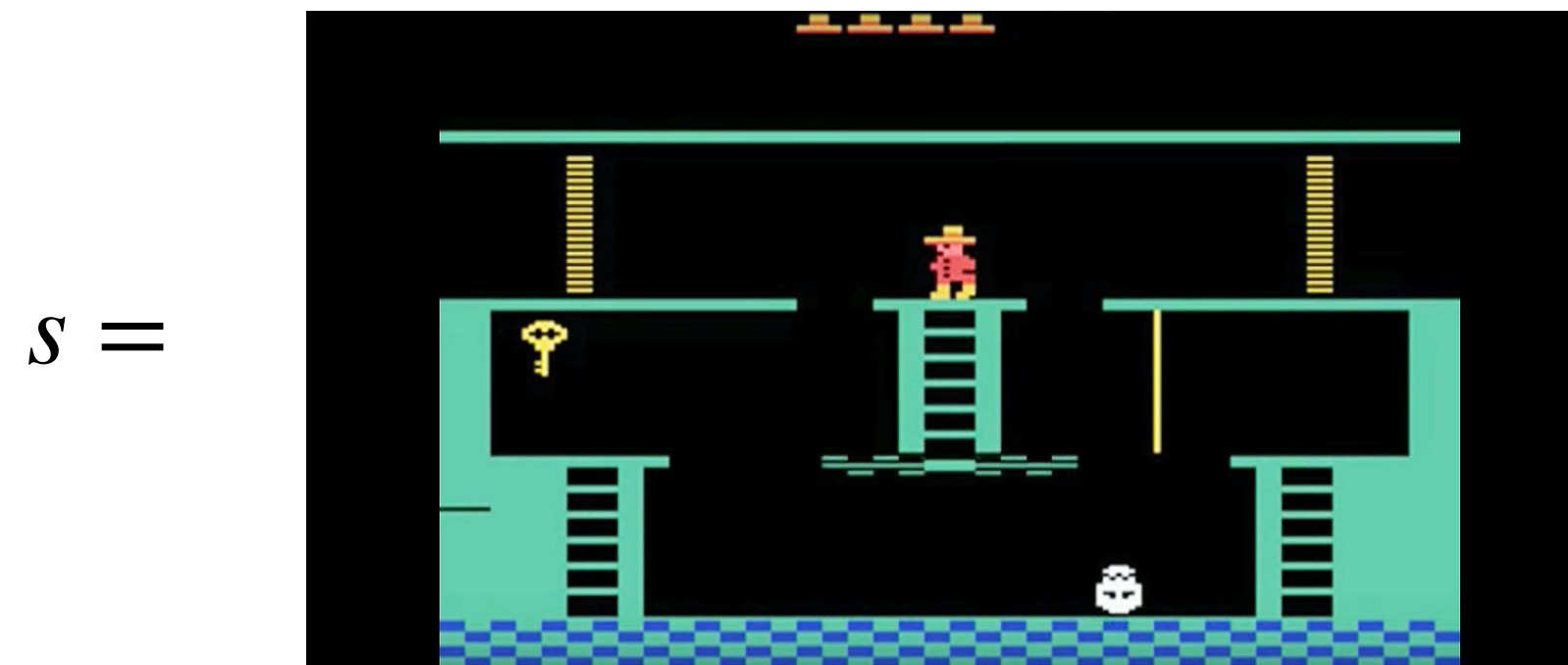
$S =$



# Setting

Finite horizon episodic MDP  $\{ \{S_h\}_{h=0}^H, \{A_h\}_{h=0}^{H-1}, H, s_0, r, P \}$

State space  $S_h$  is extremely large:



Not acceptable:  $\text{poly}(|S|)$

Need to generalize via (nonlinear) function approximation

# Let's set up function class in RL setting

We will consider **Q function class**

$$\mathcal{F} \subset S \times A \mapsto [0, H]$$

# Let's set up function class in RL setting

We will consider **Q function class**

$$\mathcal{F} \subset S \times A \mapsto [0, H]$$

**Realizability** assumption:

$$Q^* \in \mathcal{F}$$

# Let's set up function class in RL setting

We will consider **Q function class**

$$\mathcal{F} \subset S \times A \mapsto [0, H]$$

**Realizability** assumption:

$$Q^* \in \mathcal{F}$$

Define **policy class**:  $\Pi = \{ \pi : \pi(s) = \arg \max_{a \in A} f(s, a), \forall s \in S \mid f \in \mathcal{F} \}$

# Let's set up function class in RL setting

We will consider **Q function class**

$$\mathcal{F} \subset S \times A \mapsto [0, H]$$

**Realizability** assumption:

$$Q^* \in \mathcal{F}$$

Define **policy class**:  $\Pi = \{ \pi : \pi(s) = \arg \max_{a \in A} f(s, a), \forall s \in S \mid f \in \mathcal{F} \}$

Define **value function class**:  $\mathcal{V} = \{ V_f : V_f(s) = \max_a f(s, a) \mid f \in \mathcal{F} \}$



## **Learning Goal:**

We will do PAC in this lecture rather than regret.

## Learning Goal:

We will do PAC in this lecture rather than regret.

Given approximation error  $\epsilon$  and failure prob  $\delta$ ,  
can we learn  $\epsilon$  *near optimal policy* (i.e.,  $V^{\hat{\pi}} \geq V^* - \epsilon$ ) in # of samples scaling  
*poly* with all relevant parameters (*here, we need poly in  $\ln(|\mathcal{F}|)$* )

## How to check if a Q-approximator is good?

We define **average** Bellman error of a Q-estimate  $g$  below:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

## How to check if a Q-approximator is good?

We define **average** Bellman error of a Q-estimate  $g$  below:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

$f$ : defines roll-in distribution over  $s_h, a_h$ .

## How to check if a Q-approximator is good?

We define **average** Bellman error of a Q-estimate  $g$  below:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

$f$ : defines roll-in distribution over  $s_h, a_h$ .

We know that  $\mathcal{E}(Q^*; f, h) = 0, \forall f$

## How to check if a Q-approximator is good?

We define **average** Bellman error of a Q-estimate  $g$  below:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

$f$ : defines roll-in distribution over  $s_h, a_h$ .

We know that  $\mathcal{E}(Q^*; f, h) = 0, \forall f$

Hence, any  $g$  such that  $\mathcal{E}(g; f, h) \neq 0$ , is an incorrect  $Q^*$  approximator

## How to check if a Q-approximator is good?

We can define **average** Bellman error wrt the V-function induced by  $g$  as well:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

## How to check if a Q-approximator is good?

We can define **average** Bellman error wrt the V-function induced by  $g$  as well:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

Again we have  $\mathcal{E}(Q^*; f, h) = 0, \forall f$



## How to check if a Q-approximator is good?

We can define **average** Bellman error wrt the V-function induced by  $g$  as well:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

Again we have  $\mathcal{E}(Q^\star; f, h) = 0, \forall f$

( because:  $V_{Q^\star}(s) - r(s, \pi_{Q^\star}(s)) - \mathbb{E}_{s' \sim P_h(\cdot | s, \pi_{Q^\star}(s))} V_{Q^\star}(s') = 0$  )

## How to check if a Q-approximator is good?

We can define **average** Bellman error wrt the V-function induced by  $g$  as well:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

Again we have  $\mathcal{E}(Q^\star; f, h) = 0, \forall f$

( because:  $V_{Q^\star}(s) - r(s, \pi_{Q^\star}(s)) - \mathbb{E}_{s' \sim P_h(\cdot | s, \pi_{Q^\star}(s))} V_{Q^\star}(s') = 0$  )

Hence, any  $g$  such that  $\mathcal{E}(g; \pi, h) \neq 0$ , is an incorrect  $Q^\star$  approximator

# The Q / V-Bellman rank

$$\forall h : \mathcal{E}_h \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$$

	$g$	$f$			
$\pi_f$	$\mathcal{E}_{g;f,h}$	$\mathcal{E}_{f;f,h}$			

# The Q / V-Bellman rank

$$\forall h : \mathcal{E}_h \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$$

	$g$	$f$			
$\pi_f$	$\mathcal{E}_{g;f,h}$	$\mathcal{E}_{f;f,h}$			

Rank of this Matrix is defined as Bellman Rank

## The Q / V-Bellman rank

In other words, there are two mappings  $W_h : \mathcal{F} \mapsto \mathbb{R}^d$ ,  $X_h : \mathcal{F} \mapsto \mathbb{R}^d$  (d = Bellman-rank)

$$\forall f, g \in \mathcal{F} : \mathcal{E}(g; f, h) = \langle W_h(g), X_h(f) \rangle$$

Note, we just assume the existence of  $W, X$ , but they are unknown

# Outline for Today

 1. Bellman rank Definitions

2. Examples that are captured by the Bellman rank framework

# The Linear Bellman Completion Model

Given feature  $\phi$ , take any linear function  $\theta^\top \phi(s, a)$ :

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s, a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

# The Linear Bellman Completion Model

Given feature  $\phi$ , take any linear function  $\theta^\top \phi(s, a)$ :

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s, a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

**Claim: it has Q-Bellman rank d**



# The Linear Bellman Completion Model

Given feature  $\phi$ , take any linear function  $\theta^\top \phi(s, a)$ :

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s, a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

**Claim: it has Q-Bellman rank d**

$\forall g(s, a) := \theta^\top \phi(s, a)$ , we have:

# The Linear Bellman Completion Model

Given feature  $\phi$ , take any linear function  $\theta^\top \phi(s, a)$ :

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s, a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

**Claim: it has Q-Bellman rank d**

$\forall g(s, a) := \theta^\top \phi(s, a)$ , we have:

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} \theta^\top \phi(s_{h+1}, a) \right] \right]$$

# The Linear Bellman Completion Model

Given feature  $\phi$ , take any linear function  $\theta^\top \phi(s, a)$ :

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s, a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

**Claim: it has Q-Bellman rank d**

$\forall g(s, a) := \theta^\top \phi(s, a)$ , we have:

$$\begin{aligned} \mathcal{E}(g; f, h) &= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} \theta^\top \phi(s_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - \mathcal{T}_h(\theta)^\top \phi(s_h, a_h) \right] \end{aligned}$$

# The Linear Bellman Completion Model

Given feature  $\phi$ , take any linear function  $\theta^\top \phi(s, a)$ :

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s, a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

**Claim: it has Q-Bellman rank d**

$\forall g(s, a) := \theta^\top \phi(s, a)$ , we have:

$$\begin{aligned} \mathcal{E}(g; f, h) &= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} \theta^\top \phi(s_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - \mathcal{T}_h(\theta)^\top \phi(s_h, a_h) \right] \\ &= \left\langle \theta - \mathcal{T}_h(\theta), \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} [\phi(s_h, a_h)] \right\rangle \end{aligned}$$

# The Linear Bellman Completion Model

Given feature  $\phi$ , take any linear function  $\theta^\top \phi(s, a)$ :

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s, a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

**Claim: it has Q-Bellman rank d**

$\forall g(s, a) := \theta^\top \phi(s, a)$ , we have:

$$\begin{aligned} \mathcal{E}(g; f, h) &= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} \theta^\top \phi(s_{h+1}, a) \right] \right] \\ &= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - \mathcal{T}_h(\theta)^\top \phi(s_h, a_h) \right] \\ &= \left\langle \theta - \mathcal{T}_h(\theta), \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} [\phi(s_h, a_h)] \right\rangle \end{aligned}$$

Note linear Bell-completion captures tabular / linear mdp already

## The Linear $Q^*$ & $V^*$ model:

Assume  $Q^*(s, a) = (w^*)^\top \phi(s, a)$ ,  $V^*(s) = (\theta^*)^\top \psi(s)$ ,  $\forall s, a$

## The Linear $Q^*$ & $V^*$ model:

Assume  $Q^*(s, a) = (w^*)^\top \phi(s, a)$ ,  $V^*(s) = (\theta^*)^\top \psi(s)$ ,  $\forall s, a$

**Claim: it has Q-Bellman rank  $2d$**

## The Linear $Q^*$ & $V^*$ model:

Assume  $Q^*(s, a) = (w^*)^\top \phi(s, a)$ ,  $V^*(s) = (\theta^*)^\top \psi(s)$ ,  $\forall s, a$

**Claim: it has Q-Bellman rank  $2d$**

$$\mathcal{F}_h = \left\{ (w, \theta) : \max_a w^\top \phi(s, a) = \theta^\top \psi(s), \forall s \right\}$$



## The Linear $Q^*$ & $V^*$ model:

Assume  $Q^*(s, a) = (w^*)^\top \phi(s, a)$ ,  $V^*(s) = (\theta^*)^\top \psi(s)$ ,  $\forall s, a$

**Claim: it has Q-Bellman rank  $2d$**

$$\mathcal{F}_h = \left\{ (w, \theta) : \max_a w^\top \phi(s, a) = \theta^\top \psi(s), \forall s \right\}$$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \theta^\top \psi(s_{h+1}) \right] \right]$$

## The Linear $Q^*$ & $V^*$ model:

Assume  $Q^*(s, a) = (w^*)^\top \phi(s, a)$ ,  $V^*(s) = (\theta^*)^\top \psi(s)$ ,  $\forall s, a$

**Claim: it has Q-Bellman rank 2d**

$$\mathcal{F}_h = \left\{ (w, \theta) : \max_a w^\top \phi(s, a) = \theta^\top \psi(s), \forall s \right\}$$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \theta^\top \psi(s_{h+1}) \right] \right]$$

$$= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - (w^*)^\top \phi(s_h, a_h) + \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ (\theta^*)^\top \psi(s_{h+1}) \right] - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \theta^\top \psi(s_{h+1}) \right] \right]$$

## The Linear $Q^*$ & $V^*$ model:

Assume  $Q^*(s, a) = (w^*)^\top \phi(s, a)$ ,  $V^*(s) = (\theta^*)^\top \psi(s)$ ,  $\forall s, a$

**Claim: it has Q-Bellman rank  $2d$**

$$\mathcal{F}_h = \left\{ (w, \theta) : \max_a w^\top \phi(s, a) = \theta^\top \psi(s), \forall s \right\}$$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} [\theta^\top \psi(s_{h+1})] \right]$$

$$= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - (w^*)^\top \phi(s_h, a_h) + \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} [(\theta^*)^\top \psi(s_{h+1})] - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} [\theta^\top \psi(s_{h+1})] \right]$$

$$= \left\langle \begin{bmatrix} w - w^* \\ \theta - \theta^* \end{bmatrix}, \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \begin{bmatrix} \phi(s_h, a_h) \\ -\mathbb{E}_{s' \sim P_h(s_h, a_h)} [\psi(s')] \end{bmatrix} \right\rangle$$

## The Linear $Q^*$ & $V^*$ model:

Assume  $Q^*(s, a) = (w^*)^\top \phi(s, a)$ ,  $V^*(s) = (\theta^*)^\top \psi(s)$ ,  $\forall s, a$

**Claim: it has Q-Bellman rank  $2d$**

$$\mathcal{F}_h = \left\{ (w, \theta) : \max_a w^\top \phi(s, a) = \theta^\top \psi(s), \forall s \right\}$$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} [\theta^\top \psi(s_{h+1})] \right]$$

$$= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - (w^*)^\top \phi(s_h, a_h) + \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} [(\theta^*)^\top \psi(s_{h+1})] - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} [\theta^\top \psi(s_{h+1})] \right]$$

$$= \left\langle \begin{bmatrix} w - w^* \\ \theta - \theta^* \end{bmatrix}, \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \begin{bmatrix} \phi(s_h, a_h) \\ -\mathbb{E}_{s' \sim P_h(s_h, a_h)} [\psi(s')] \end{bmatrix} \right\rangle$$

As we will see, linear  $Q^*$  &  $V^*$  is learnable, and recall linear  $Q^*$  is not...

# $Q^*$ - state abstraction

We have a small latent state space  $Z$ , and a **known** mapping  $\xi$  from state  $s$  to  $z$

$$Q^*(s_1, a) = Q^*(s_2, a), \forall a, \text{ if } \xi(s_1) = \xi(s_2)$$

# $Q^*$ - state abstraction

We have a small latent state space  $Z$ , and a **known** mapping  $\xi$  from state  $s$  to  $z$

$$Q^*(s_1, a) = Q^*(s_2, a), \forall a, \text{ if } \xi(s_1) = \xi(s_2)$$

**Claim: this model has Q-Bellman rank  $|Z| |A| + |Z|$**

We can show that this model is captured by linear  $Q^*$  &  $V^*$

# Low-rank MDP

$$P_h(s' | s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a) \quad (\text{neither } \mu^\star \text{ nor } \phi^\star \text{ is known})$$

# Low-rank MDP

$$P_h(s' | s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a) \quad (\text{neither } \mu^\star \text{ nor } \phi^\star \text{ is known})$$

**Claim:** this model has V-Bellman rank  $d$



# Low-rank MDP

$$P_h(s' | s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a) \quad (\text{neither } \mu^\star \text{ nor } \phi^\star \text{ is known})$$

**Claim:** this model has V-Bellman rank  $d$

Define representation class  $\Phi$ , with  $\phi^\star \in \Phi$

$$\mathcal{F}_h = \{\theta^\top \phi(\cdot, \cdot) : \|\theta\|_2 \leq W, \phi \in \Phi\}$$

# Low-rank MDP

$$P_h(s' | s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a) \quad (\text{neither } \mu^\star \text{ nor } \phi^\star \text{ is known})$$

**Claim: this model has V-Bellman rank  $d$**

Define representation class  $\Phi$ , with  $\phi^\star \in \Phi$

$$\mathcal{F}_h = \{\theta^\top \phi(\cdot, \cdot) : \|\theta\|_2 \leq W, \phi \in \Phi\}$$

$$\mathbb{E}_{s_h \sim d_h^{\pi_g}} \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right]$$

# Low-rank MDP

$$P_h(s' | s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a) \quad (\text{neither } \mu^\star \text{ nor } \phi^\star \text{ is known})$$

**Claim: this model has V-Bellman rank  $d$**

Define representation class  $\Phi$ , with  $\phi^\star \in \Phi$

$$\mathcal{F}_h = \{\theta^\top \phi(\cdot, \cdot) : \|\theta\|_2 \leq W, \phi \in \Phi\}$$

$$\begin{aligned} & \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] \\ &= \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} \mathbb{E}_{s_h \sim P_{h-1}(\cdot | \tilde{s}, \tilde{a})} \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] \end{aligned}$$

# Low-rank MDP

$$P_h(s' | s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a) \quad (\text{neither } \mu^\star \text{ nor } \phi^\star \text{ is known})$$

**Claim: this model has V-Bellman rank  $d$**

Define representation class  $\Phi$ , with  $\phi^\star \in \Phi$

$$\mathcal{F}_h = \{\theta^\top \phi(\cdot, \cdot) : \|\theta\|_2 \leq W, \phi \in \Phi\}$$

$$\begin{aligned} & \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] \\ &= \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} \mathbb{E}_{s_h \sim P_{h-1}(\cdot | \tilde{s}, \tilde{a})} \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] \\ &= \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} \int_{s_h} \mu_{h-1}^\star(s_h)^\top \phi_{h-1}^\star(\tilde{s}, \tilde{a}) \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] d(s_h) \end{aligned}$$

# Low-rank MDP

$$P_h(s' | s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a) \quad (\text{neither } \mu^\star \text{ nor } \phi^\star \text{ is known})$$

**Claim: this model has V-Bellman rank  $d$**

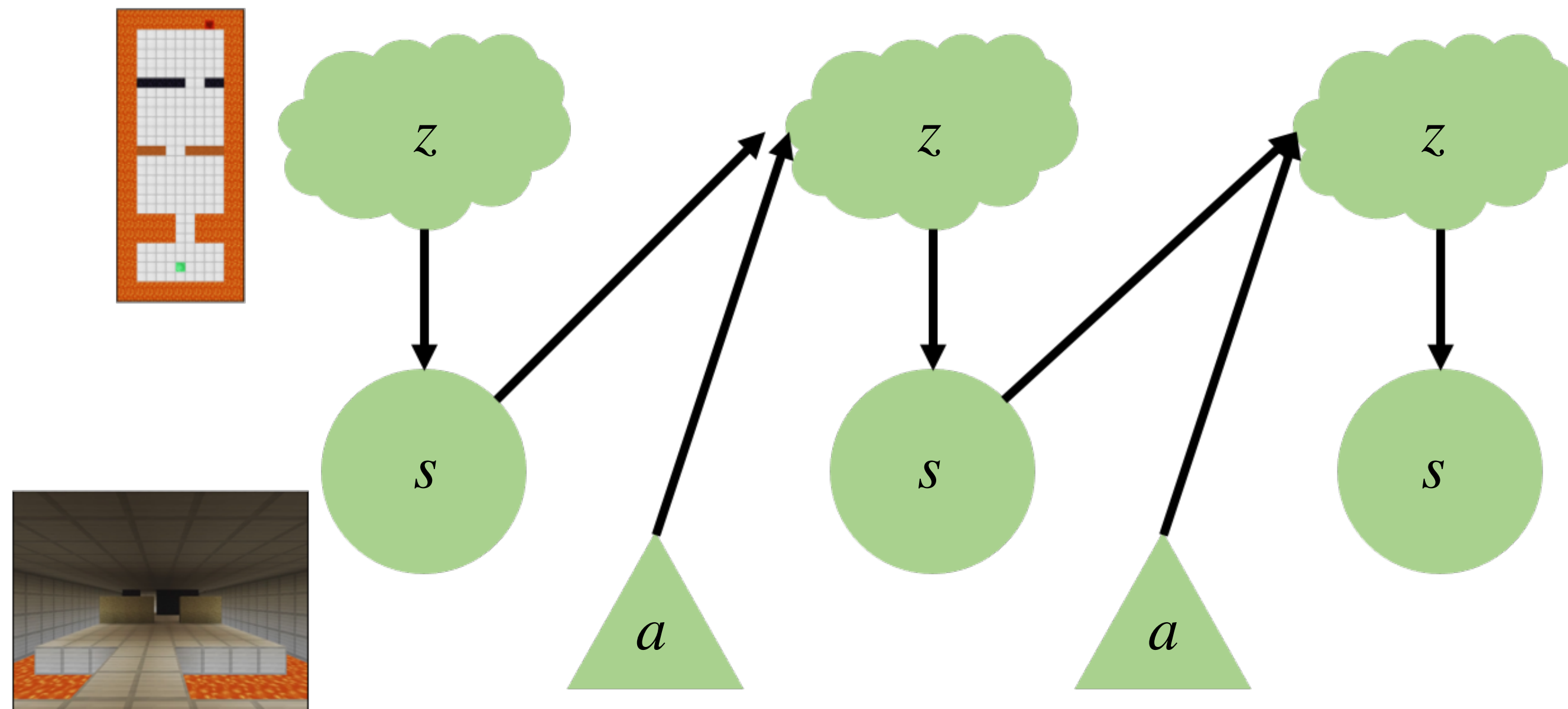
Define representation class  $\Phi$ , with  $\phi^\star \in \Phi$

$$\mathcal{F}_h = \{\theta^\top \phi(\cdot, \cdot) : \|\theta\|_2 \leq W, \phi \in \Phi\}$$

$$\begin{aligned} & \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] \\ &= \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} \mathbb{E}_{s_h \sim P_{h-1}(\cdot | \tilde{s}, \tilde{a})} \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] \\ &= \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} \int_{s_h} \mu_{h-1}^\star(s_h)^\top \phi_{h-1}^\star(\tilde{s}, \tilde{a}) \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] d(s_h) \\ &= \left\langle \int_{s_h} \mu_{h-1}^\star(s_h) \left[ V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, \pi_g(s_h))} [V_g(s_{h+1})] \right] d(s_h), \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} [\phi_{h-1}^\star(\tilde{s}, \tilde{a})] \right\rangle \end{aligned}$$

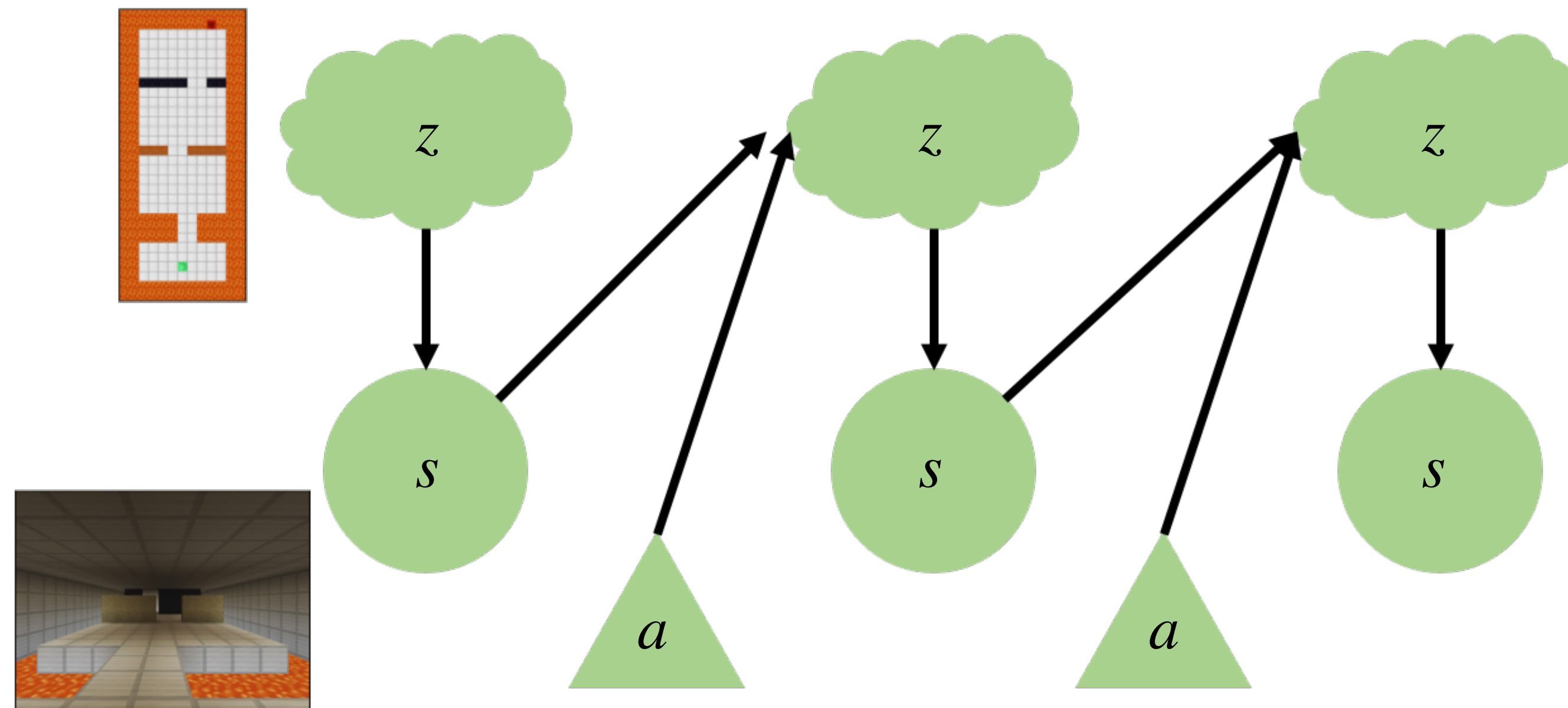
# Latent variable MDP

Latent variable MDP is captured by low-rank MDP, so it has small V-Bellman rank...



# Latent variable MDP

Latent variable MDP is captured by low-rank MDP, so it has small V-Bellman rank...

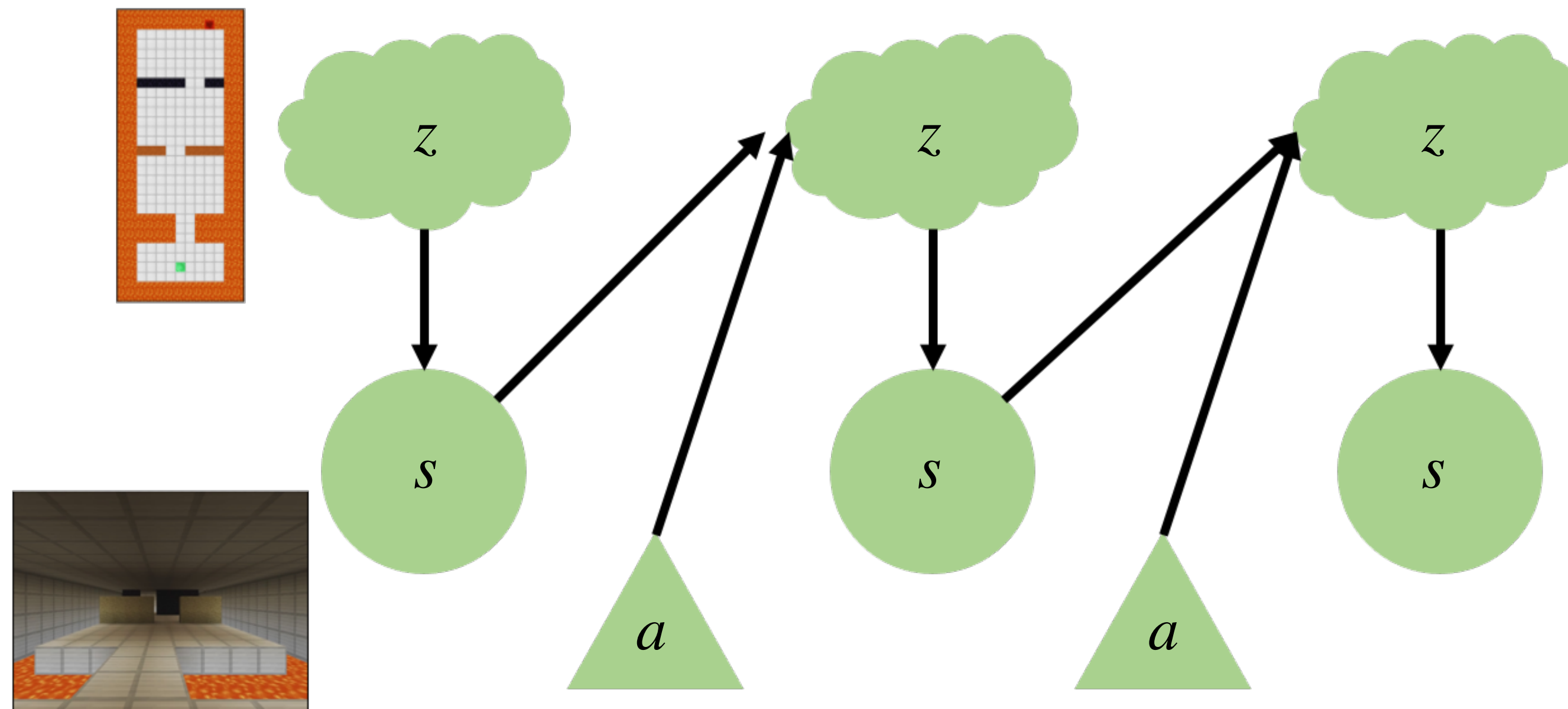


Given  $s, a: z \sim \phi^*(s, a), s' \sim \nu^*(z)$



# Latent variable MDP

Latent variable MDP is captured by low-rank MDP, so it has small V-Bellman rank...



Given  $s, a: z \sim \phi^*(s, a), s' \sim \nu^*(z)$

V-Bellman rank = Number of latent states



# Summary

1. Q-Bellman rank: related to the Bellman error of a Q function estimate  $g$ :

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

# Summary

1. Q-Bellman rank: related to the Bellman error of a Q function estimate  $g$ :

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

2. V-Bellman rank: related to the Bellman error of a V function estimate

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

# Summary

1. Q-Bellman rank: related to the Bellman error of a Q function estimate  $g$ :

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

2. V-Bellman rank: related to the Bellman error of a V function estimate

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

3. Small Bellman rank means that:

$$\forall f, g \in \mathcal{F} : \mathcal{E}(g; f, h) = \langle W_h(g), X_h(f) \rangle$$

where  $X_h(f)$ ,  $W_h(f)$  are  
low-dim vectors

# Summary

1. Q-Bellman rank: related to the Bellman error of a Q function estimate  $g$ :

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right] \right]$$

2. V-Bellman rank: related to the Bellman error of a V function estimate

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

3. Small Bellman rank means that:

$$\forall f, g \in \mathcal{F} : \mathcal{E}(g; f, h) = \langle W_h(g), X_h(f) \rangle$$

where  $X_h(f), W_h(f)$  are  
low-dim vectors

4. Many models (more in the book chapter) indeed have low-Q or V Bellman rank

## Next week:

A general algorithm that can learn an  $\epsilon$  near optimal policy w/ # of samples

$\text{poly}(H, 1/\epsilon, \ln(|\mathcal{H}|), \text{b-rank})$