# On the Curses of Future and History in Partially Observed Off-policy Evaluation

Nan Jiang

University of Illinois at Urbana-Champaign

March 2024
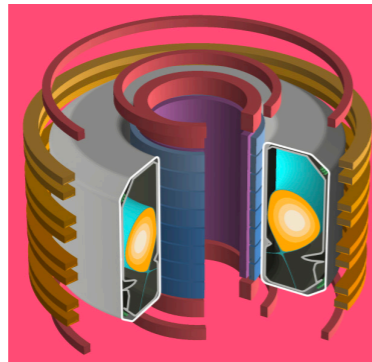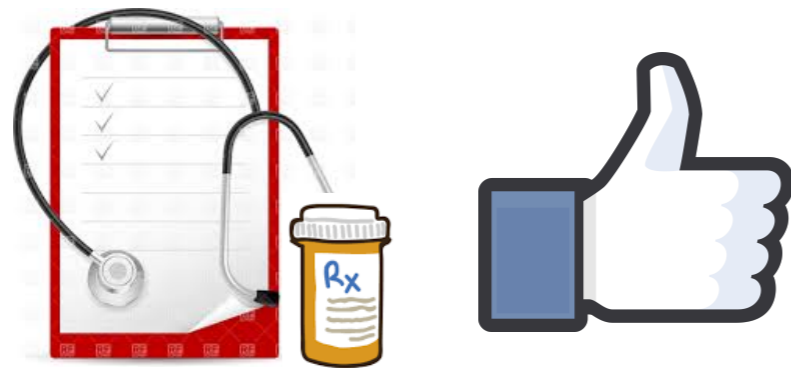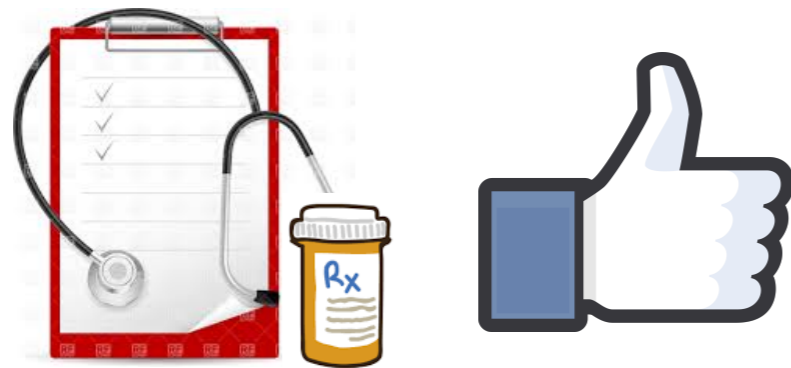
Masatoshi
Uehara

Yuheng
Zhang

Key ingredient: **simulator**

- Unlimited data

- Decision w/o real consequences

- Can easily evaluate new strategy

Key ingredient: **simulator**

- Unlimited data  **X**

- Decision w/o real consequences **X**

- Can easily evaluate new strategy **X**

Offline Reinforcement Learning

Key ingredient: **simulator**

- Unlimited data  **X**

- Decision w/o real consequences **X**

- Can easily evaluate new strategy **X**

# Supervised learning pipeline


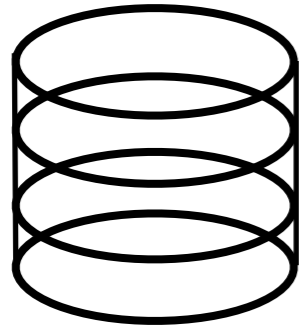
data

# Supervised learning pipeline

# Supervised learning pipeline

# Offline RL pipeline



data
$\{ (s, a, r, s') \}$

training    validation    test

# Offline RL pipeline



training     validation     test

data

$\{\,(s,\,a,\,r,\,s')\,\}$

# Offline RL pipeline



data
$\{ (s, \boxed{a,} r, s') \}$

training    validation    test

# Offline RL pipeline



data
$\{ (s, a, r, s') \}$

training        validation        test

Off-policy Evaluation
(OPE)

# Unbiased OPE

## Importance sampling (IS) [Precup'00]

Behavior

Target

Precup. 2000. Eligibility traces for off-policy policy evaluation.

# Unbiased OPE

## Importance sampling (IS) [Precup'00]

Behavior

Target

Precup. 2000. Eligibility traces for off-policy policy evaluation.

# Unbiased OPE

## Importance sampling (IS) [Precup'00]



Behavior          Target

Precup. 2000. Eligibility traces for off-policy policy evaluation.

# Unbiased OPE

**Importance sampling (IS)** [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓

Behavior        Target

Precup. 2000. Eligibility traces for off-policy policy evaluation.

# Unbiased OPE

**Importance sampling (IS)** [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)

- No Markovianity required  ✓

- Exponential-in-horizon variance!

Data

Candidate

Behavior

Target

Precup. 2000. Eligibility traces for off-policy policy evaluation.

# Unbiased OPE

## Importance sampling (IS) [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)

- No Markovianity required ✓

- Exponential-in-horizon variance!

Behavior          Target

Data

Candidate

$$o_1, a_1, r_1, \ldots, o_H, a_H, r_H \quad \Rightarrow \quad \left( \prod_{h=1}^{H} \frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)} \right) \left( \sum_{h=1}^{H} r_h \right)$$

Precup. 2000. Eligibility traces for off-policy policy evaluation.

# Unbiased OPE

**Importance sampling (IS)** [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)

- No Markovianity required ✓

- Exponential-in-horizon variance!

Behavior        Target

Data

Candidate

$$o_1, a_1, r_1, \ldots, o_H, a_H, r_H \quad \Rightarrow \quad \left( \prod_{h=1}^{H} \frac{\pi(a_h | o_h)}{\pi_b(a_h | o_h)} \right) \left( \sum_{h=1}^{H} r_h \right)$$

Precup. 2000. Eligibility traces for off-policy policy evaluation.

# Unbiased OPE

**Importance sampling (IS)** [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)

- No Markovianity required ✓

- **Exponential-in-horizon** variance!

Behavior    Target

Data

Candidate

$$o_1, a_1, r_1, \ldots, o_H, a_H, r_H \;\;\Rightarrow\;\; \left( \prod_{h=1}^{H} \frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)} \right) \left( \sum_{h=1}^{H} r_h \right)$$

- Or, can only evaluate $\pi$ when $\prod_{h=1}^{H} \frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)}$ small

Precup. 2000. Eligibility traces for off-policy policy evaluation.

4

# Unbiased OPE

## Importance sampling (IS) [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)

- No Markovianity required ✓

- **Exponential-in-horizon** variance!

Behavior        Target

Data

Candidate

$$o_1, a_1, r_1, \ldots, o_H, a_H, r_H \quad \Rightarrow \quad \left( \prod_{h=1}^{H} \frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)} \right) \left( \sum_{h=1}^{H} r_h \right)$$

- Or, can only evaluate $\pi$ when $\prod_{h=1}^{H} \frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)}$ small

IS' measure of *coverage*

Precup. 2000. Eligibility traces for off-policy policy evaluation.

4

# Better coverage?

# Better coverage?

# Better coverage?

# Better coverage?

**FQE** [Munos, Szepesvari… CJ'19, …] **/ MIS** [Liu et al'18, Nachum et al'19, UHJ'20, …]

- Assume MDPs: $o_1(= s_1), a_1, r_1, \ldots, o_H(= s_H), a_H, r_H$

# Better coverage?

**FQE** [Munos, Szepesvari… CJ'19, …] **/ MIS** [Liu et al'18, Nachum et al'19, UHJ'20, …]

- Assume MDPs: $o_1(= s_1), a_1, r_1, \ldots, o_H(= s_H), a_H, r_H$

- Learn value functions: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^{H} r_{h'} | s_h]$

  - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.

# Better coverage?

**FQE** [Munos, Szepesvari… CJ'19, …] / **MIS** [Liu et al'18, Nachum et al'19, UHJ'20, …]

- Assume MDPs: $o_1(= s_1), a_1, r_1, \ldots, o_H(= s_H), a_H, r_H$

- Learn value functions: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^{H} r_{h'} | s_h]$

  - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.

  - Estimate from class $\mathcal{V}$ by $\arg\min_{V \in \mathcal{V}} \mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]$

\* Also needs Bellman completeness

# Better coverage?

**FQE** [Munos, Szepesvari… CJ'19, …] / **MIS** [Liu et al'18, Nachum et al'19, UHJ'20, …]

- Assume MDPs: $o_1(= s_1), a_1, r_1, \ldots, o_H(= s_H), a_H, r_H$

- Learn value functions: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^{H} r_{h'}|s_h]$

  - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.

  - Estimate from class $\mathcal{V}$ by $\arg\min_{V \in \mathcal{V}} \mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]$

- Coverage: *marginal* state distribution

  - require $\frac{d^\pi(s_h)}{d^{\pi_b}(s_h)}$ and $\frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)}$ small



5

* Also needs Bellman completeness

# Better coverage?

**FQE** [Munos, Szepesvari… CJ'19, …] **/ MIS** [Liu et al'18, Nachum et al'19, UHJ'20, …]

- Assume MDPs: $o_1(= s_1), a_1, r_1, \ldots, o_H(= s_H), a_H, r_H$

- Learn value functions: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^{H} r_{h'} | s_h]$

  - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.

  - Estimate from class $\mathcal{V}$ by $\arg\min_{V \in \mathcal{V}} \mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]$

- Coverage: *marginal* state distribution

  - require $\frac{d^\pi(s_h)}{d^{\pi_b}(s_h)}$ and $\frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)}$ small

  - Can further improve to e.g.,
    $$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}}$$

* Also needs Bellman completeness

# Better coverage?

**FQE** [Munos, Szepesvari… CJ'19, …] **/ MIS** [Liu et al'18, Nachum et al'19, UHJ'20, …]

- Assume MDPs: $o_1 (= s_1), a_1, r_1, \ldots, o_H (= s_H), a_H, r_H$

- Learn value functions: $V^\pi(s_h) := \mathbb{E}_\pi [\sum_{h'=h}^{H} r_{h'} | s_h]$

  - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.

  - Estimate from class $\mathcal{V}$ by $\arg\min_{V \in \mathcal{V}} \mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]$

- Coverage: *marginal* state distribution

  - require $\frac{d^\pi(s_h)}{d^{\pi_b}(s_h)}$ and $\frac{\pi(a_h | o_h)}{\pi_b(a_h | o_h)}$ small

  - Can further improve to e.g.,

    $\sup_{V \in \mathcal{V}} \dfrac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}}$

- Fundamental to offline training

  & online exploration

\* Also needs Bellman completeness

# Partially Observed (non-Markov) Problems

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, \underbrace{o_h, a_h, r_h}, \ldots o_H, a_H, r_H$$

- Define new state $(\tau_h, o_h)$. Problem solved?

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

- Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^{\pi}(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)}$

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

- Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^\pi(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)} = \displaystyle\prod_{h'=1}^{h-1} \dfrac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})}$ **!!**

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$\underbrace{o_1, a_1, r_1, \ldots,}_{} o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

- Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^\pi(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)} = \displaystyle\prod_{h'=1}^{h-1} \dfrac{\color{red}\pi(a_{h'}|o_{h'})}{\color{blue}\pi_b(a_{h'}|o_{h'})}$ **!!**

- Structured can help: e.g., if $\mathcal{V}$ is linear in feature $\phi : \mathcal{S} \to \mathbb{R}^d$

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

$$\underbrace{\phantom{o_1, a_1, r_1, \ldots,}}$$

- Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^\pi(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)} = \displaystyle\prod_{h'=1}^{h-1} \dfrac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})}$ **!!**

- Structured can help: e.g., if $\mathcal{V}$ is linear in feature $\phi : \mathcal{S} \to \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

\* Also needs Bellman completeness

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

$$\underbrace{\phantom{o_1, a_1, r_1, \ldots,}}$$

- Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^\pi(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)} = \displaystyle\prod_{h'=1}^{h-1} \dfrac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})}$ **‼**

- Structured can help: e.g., if $\mathcal{V}$ is linear in feature $\phi : \mathcal{S} \to \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

- this assumes given low-dim linear feature to encode history…

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

$\underbrace{\hspace{3cm}}$

- Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^\pi(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)} = \displaystyle\prod_{h'=1}^{h-1} \dfrac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})}$ **!!**

- Structured can help: e.g., if $\mathcal{V}$ is linear in feature $\phi : \mathcal{S} \to \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

  - this assumes given low-dim linear feature to encode history…
  - side q: what structure in $\mathcal{V}$ balances expressivity and coverage

* Also needs Bellman completeness

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

$$\underbrace{\phantom{o_1, a_1, r_1, \ldots}}$$

  - Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^\pi(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)} = \displaystyle\prod_{h'=1}^{h-1} \dfrac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})}$ **‼**

- Structured can help: e.g., if $\mathcal{V}$ is linear in feature $\phi : \mathcal{S} \to \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$
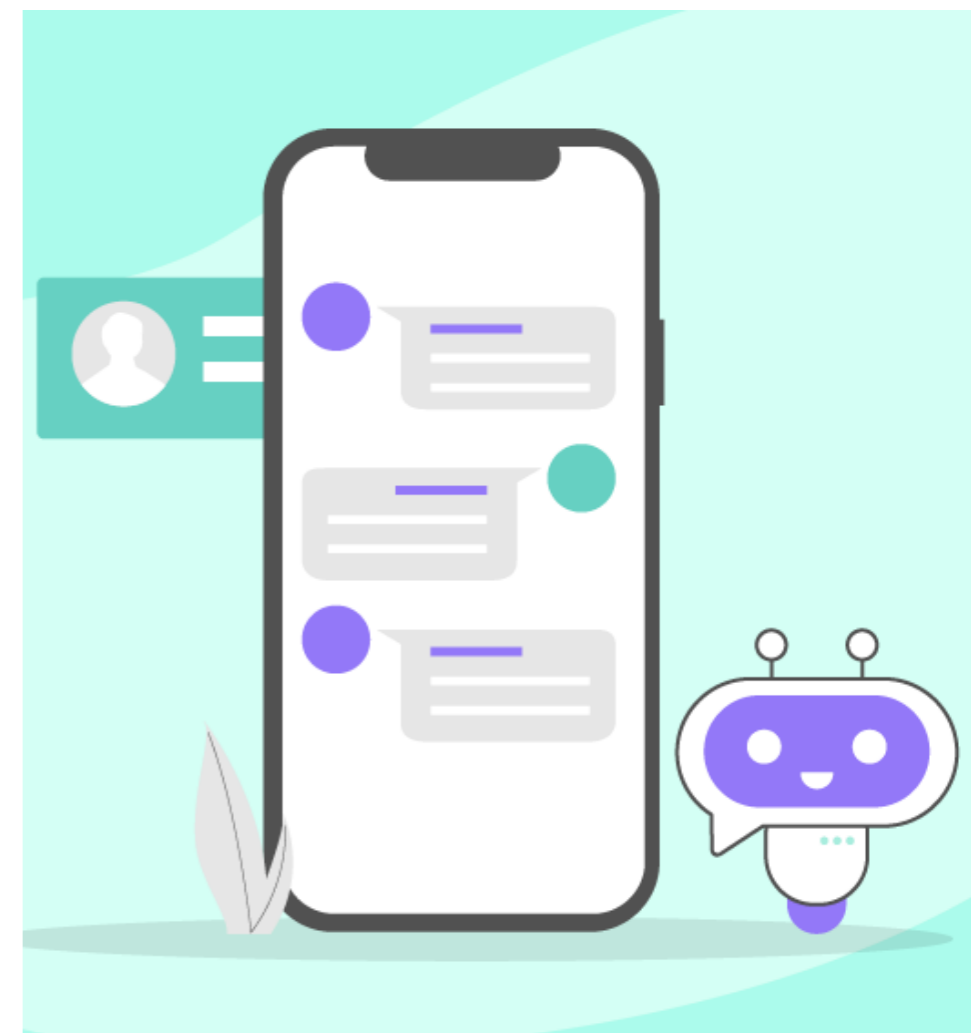
  - this assumes given low-dim linear feature to encode history…
  - side q: what structure in $\mathcal{V}$ balances expressivity and coverage
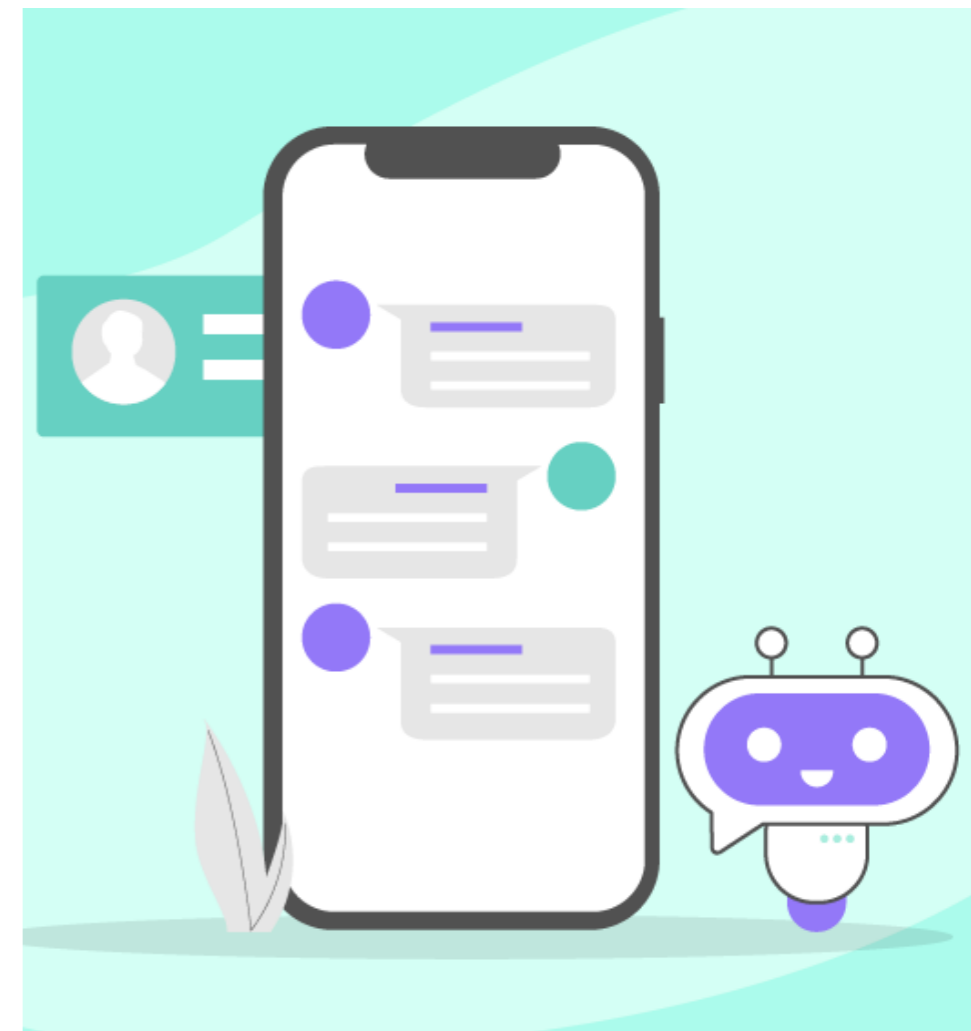  - connection to known empirical evidence (LLMs, RLHF, etc.)

* Also needs Bellman completeness

# Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \ldots, o_h, a_h, r_h, \ldots o_H, a_H, r_H$$

$$\underbrace{\qquad\qquad\qquad}$$

- Define new state $(\tau_h, o_h)$. Problem solved?

- State density ratio: $\dfrac{d^\pi(o_1, a_1, \ldots, o_h)}{d^{\pi_b}(o_1, a_1, \ldots, o_h)} = \displaystyle\prod_{h'=1}^{h-1} \dfrac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})}$ **‼**

- Structured can help: e.g., if $\mathcal{V}$ is linear in feature $\phi : \mathcal{S} \to \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \le \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

- this assumes given low-dim linear feature to encode history…

Can we avoid the exponentials in OPE in PO settings, without relying on structured function classes?

* Also needs Bellman completeness

# Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \ldots, H,$

    - nature generates *latent* state $s_h \in S_h$ (small?)

$s_1$

$h{=}1$

# Partially Observable MDPs (POMDPs)

$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,

  - nature generates *latent* state $s_h \in S_h$ (small?)

  - agent observes $o_h \in O_h$ (large), $o_h \sim \mathsf{O}(\,\cdot\mid s_h)$

  emission process
  $$\mathsf{O}: S \to \Delta(O)$$

# Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \ldots, H,$

  - nature generates *latent* state $s_h \in S_h$ (small?)

  - agent observes $o_h \in O_h$ (large), $o_h \sim \mathsf{O}(\,\cdot\,|\,s_h)$

  - chooses action $a_h \in A$   (small)

emission process
$$\mathsf{O}: S \rightarrow \Delta(O)$$



$h=1$

7

# Partially Observable MDPs (POMDPs)

$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,

  - nature generates *latent* state $s_h \in S_h$ (small?)

  - agent observes $o_h \in O_h$ (large), $o_h \sim \mathsf{O}(\cdot \mid s_h)$

  - chooses action $a_h \in A$ (small)

  - receives reward $r_h = R(o_h, a_h)$

emission process
$$\mathsf{O}: S \to \Delta(O)$$

reward function
$$R: S \times A \to [0,1]$$



$o_1$

$\mathsf{O}$

$a_1$

$s_1$

$h=1$

7

# Partially Observable MDPs (POMDPs)

$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,
  - nature generates *latent* state $s_h \in S_h$ (small?)
    - $s_h \sim P(\cdot \mid s_{h-1}, a_{h-1})$ for $h \geq 2$
  - agent observes $o_h \in O_h$ (large), $o_h \sim O(\cdot \mid s_h)$
  - chooses action $a_h \in A$   (small)
  - receives reward $r_h = R(o_h, a_h)$

transition dynamics
$$P: S \times A \to \Delta(S)$$

emission process
$$O: S \to \Delta(O)$$

reward function
$$R: S \times A \to [0,1]$$



7

# Partially Observable MDPs (POMDPs)

$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,
  - nature generates *latent* state $s_h \in S_h$ (small?)
    - $s_h \sim P(\,\cdot\mid s_{h-1}, a_{h-1})$ for $h \geq 2$
  - agent observes $o_h \in O_h$ (large), $o_h \sim \mathsf{O}(\,\cdot\mid s_h)$
  - chooses action $a_h \in A$   (small)
  - receives reward $r_h = R(o_h, a_h)$

transition dynamics
$$P: S \times A \to \Delta(S)$$

emission process
$$\mathsf{O}: S \to \Delta(O)$$

reward function
$$R: S \times A \to [0,1]$$

7

$h=1$   $2$

# Partially Observable MDPs (POMDPs)

$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,

  - nature generates *latent* state $s_h \in S_h$ (small?)

    - $s_h \sim P(\,\cdot \mid s_{h-1}, a_{h-1})$ for $h \geq 2$

  - agent observes $o_h \in O_h$ (large), $o_h \sim O(\,\cdot \mid s_h)$

  - chooses action $a_h \in A$  (small)

  - receives reward $r_h = R(o_h, a_h)$

transition dynamics
$$P: S \times A \to \Delta(S)$$

emission process
$$O: S \to \Delta(O)$$

reward function
$$R: S \times A \to [0,1]$$



7

# Partially Observable MDPs (POMDPs)

$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,
  - nature generates *latent* state $s_h \in S_h$ (small?)
    - $s_h \sim P(\,\cdot\mid s_{h-1}, a_{h-1})$ for $h \geq 2$
  - agent observes $o_h \in O_h$ (large), $o_h \sim \mathsf{O}(\,\cdot\mid s_h)$
  - chooses action $a_h \in A$   (small)
  - receives reward $r_h = R(o_h, a_h)$

- *Memoryless* policies $\pi : \mathcal{O} \to \Delta(\mathcal{A})$

transition dynamics
$$P : S \times A \to \Delta(S)$$

emission process
$$\mathsf{O} : S \to \Delta(O)$$

reward function
$$R : S \times A \to [0,1]$$



7

# Partially Observable MDPs (POMDPs)

$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,

  - nature generates *latent* state $s_h \in S_h$ (small?)

    - $s_h \sim P(\cdot \mid s_{h-1}, a_{h-1})$ for $h \geq 2$

  - agent observes $o_h \in O_h$ (large), $o_h \sim \mathsf{O}(\cdot \mid s_h)$

  - chooses action $a_h \in A$ (small)

  - receives reward $r_h = R(o_h, a_h)$

- *Memoryless* policies $\pi : \mathcal{O} \to \Delta(\mathcal{A})$

- Goal: estimate $J(\pi) := \mathbb{E}_\pi \left[ \sum_{h=1}^H r_h \right]$ using episodes collected by $\pi_b$

transition dynamics
$$P: S \times A \to \Delta(S)$$

emission process
$$\mathsf{O}: S \to \Delta(O)$$

reward function
$$R: S \times A \to [0,1]$$



7

# Partially Observable MDPs (POMDPs)
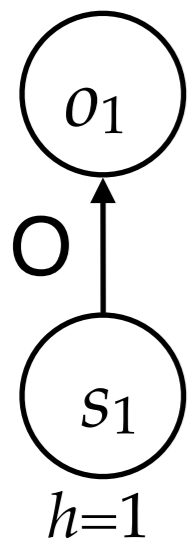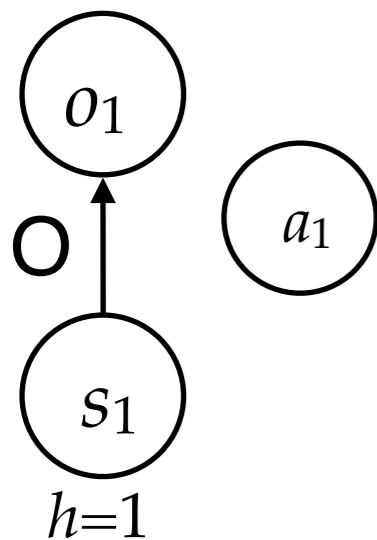
$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$
$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

- For $h = 1, 2, \ldots, H$,

  - nature generates *latent* state $s_h \in S_h$ (small?)

    - $s_h \sim P(\,\cdot\,|\,s_{h-1}, a_{h-1})$ for $h \geq 2$

  - agent observes $o_h \in O_h$ (large), $o_h \sim \mathsf{O}(\,\cdot\,|\,s_h)$

  - chooses action $a_h \in A$   (small)

  - receives reward $r_h = R(o_h, a_h)$

- *Memoryless* policies $\pi : \mathcal{O} \rightarrow \Delta(\mathcal{A})$

- Goal: estimate $J(\pi) := \mathbb{E}_\pi[\sum_{h=1}^{H} r_h]$ using episodes collected by $\pi_b$

transition dynamics
$$P: S \times A \rightarrow \Delta(S)$$

emission process
$$\mathsf{O}: S \rightarrow \Delta(O)$$

reward function
$$R: S \times A \rightarrow [0,1]$$



Coverage over *latent* state?

7

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}\left[\sum_{h'=h}^{H} r_{h'} | s_h\right] \in [0, H]$$

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'} | s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}|s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}|s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

9

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\textstyle\sum_{h'=h}^{H} r_{h'} | s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

  - $\mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

$$\underbrace{o_1, a_1, \ldots, \overbrace{o_h, a_h, \underbrace{\ldots o_H, a_H}_{f_{h+1}}}^{f_h}}_{\tau_h}$$

$s_h$ here

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\textstyle\sum_{h'=h}^{H} r_{h'}|s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

  - $\mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}|s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

  - $\mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}|s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

  - $\mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

1. Does (well-behaved) $V_{\mathcal{F}}^{\pi}$ even exist?

2. Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

$f_h$

$s_h \cdots\cdots\cdots \Pr_{\pi_b}[f_h|s_h]$

Outcome Matrix $M_{\mathcal{F}}$

$$o_1, a_1, \ldots, \underbrace{o_h, a_h, \ldots o_H, a_H}_{}$$

$s_h$ here     $f_h$

$\underbrace{\phantom{o_1, a_1, \ldots}}_{\tau_h}$    $\underbrace{\phantom{o_h, a_h, \ldots}}_{f_{h+1}}$

$\times \begin{vmatrix} V_{\mathcal{F}}^{\pi} \end{vmatrix} = \begin{vmatrix} V_{\mathcal{S}}^{\pi} \end{vmatrix}$

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'} | s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

  - $\mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

1. Does (well-behaved) $V_{\mathcal{F}}^{\pi}$ even exist?
2. Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?



$s_h$ here

$$\underbrace{o_1, a_1, \ldots,}_{\tau_h} \underbrace{o_h, a_h, \ldots o_H, a_H}_{}$$

$f_h$

$f_{h+1}$

$f_h$

$s_h \text{-----------} \mathrm{Pr}_{\pi_b}[f_h | s_h]$

Outcome Matrix $M_{\mathcal{F}}$

$\times \begin{vmatrix} V_{\mathcal{F}}^{\pi} \end{vmatrix} = \begin{vmatrix} V_{\mathcal{S}}^{\pi} \end{vmatrix}$

9

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

$V_{\mathcal{F}}^{\pi}$

$V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

| | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
|---|---|---|
| Candidate | $V_{\mathcal{F}}$ | |
| | | |

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

| | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
|---|---|---|
| Candidate | $V_{\mathcal{F}}$ | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| | | |

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

| | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
|---|---|---|
| Candidate | $V_{\mathcal{F}}$ | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

| | | | |
|---|---|---|---|
| | $V_{\mathcal{F}}^{\pi}$ | | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
| Candidate | $V_{\mathcal{F}}$ | | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ | $=$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
|---|---|---|
| Candidate | $V_{\mathcal{F}}$ | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$  **=** | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h} \sim \pi \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}}]$ | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |

$$V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$$

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
|---|---|---|
| Candidate | $V_{\mathcal{F}}$ | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ **=** | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h}\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}]$ **=** | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |

$$V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$$

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | | | |
|---|---|---|---|
| | $V_{\mathcal{F}}^{\pi}$ | | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
| Candidate | $V_{\mathcal{F}}$ | | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ | **=** | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h}\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}]$ | **=** | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |
| | $V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$ | | |
| Learning objective | | | $\mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi}V_{\mathcal{S}})(s_h))^2]$ |

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | | |
|---|---|---|
| | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
| Candidate | $V_{\mathcal{F}}$ | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ $\quad=\quad$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h}\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}] \quad=\quad \sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ <br> $V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$ | |
| Learning objective | $\sum_{h=1}^{H} \mathbb{E}_{s_h\sim\pi_b}\left[\mathbb{E}_{\substack{a_h\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|s_h]^2\right] \;=\; \mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi}V_{\mathcal{S}})(s_h))^2]$ | |

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | | |
|---|---|---|
| | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
| Candidate | $V_{\mathcal{F}}$ | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ $=$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h} \sim \pi \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}}]$ $=$ | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |
| | $V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$ | |
| Learning objective | $\sum_{h=1}^{H} \mathbb{E}_{s_h \sim \pi_b}\left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}}|s_h]^2\right]$ **✗** $=$ | $\mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi}V_{\mathcal{S}})(s_h))^2]$ |

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | | |
|---|---|---|
| | $V_{\mathcal{F}}^{\pi}$ | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
| Candidate | $V_{\mathcal{F}}$ | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ $\quad=\quad$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h}\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}] \;=\; V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$ | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |
| Learning objective | $\sum_{h=1}^{H} \mathbb{E}_{s_h\sim\pi_b}\left[\mathbb{E}_{\substack{a_h\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|s_h]^2\right]$ **✗** $=$ | $\mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi}V_{\mathcal{S}})(s_h))^2]$ |

$$\sum_{h=1}^{H} \mathbb{E}_{\tau_h\sim\pi_b}\left[\mathbb{E}_{\substack{a_h\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|\tau_h]^2\right]$$

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | $V_{\mathcal{F}}^{\pi}$ | | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
|---|---|---|---|
| Candidate | $V_{\mathcal{F}}$ | | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ | $=$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h}\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}]$  $V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$ | $=$ | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |
| Learning objective | $\sum_{h=1}^{H} \mathbb{E}_{s_h\sim\pi_b}\left[\mathbb{E}_{\substack{a_h\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|s_h]^2\right]$  **X** | $=$ | $\mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi}V_{\mathcal{S}})(s_h))^2]$ |

$$\sum_{h=1}^{H} \mathbb{E}_{\tau_h\sim\pi_b}\left[\mathbb{E}_{\substack{a_h\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|\tau_h]^2\right]$$

$$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h\sim\pi_b}\left[\left(\sum_{s_h}\mathbb{E}_{\substack{a_h\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|s_h]\cdot\Pr[s_h|\tau_h]\right)^2\right]$$

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | | | |
|---|---|---|---|
| | $V_{\mathcal{F}}^{\pi}$ | | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
| Candidate | $V_{\mathcal{F}}$ | | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ | $=$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h} \sim \pi \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}}]$  $V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$ | $=$ | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |
| Learning objective | $\sum_{h=1}^{H} \mathbb{E}_{s_h \sim \pi_b}\left[\mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}}|s_h]^2\right]$ X | $=$ | $\mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi} V_{\mathcal{S}})(s_h))^2]$ |

$$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b}\left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}}|\tau_h]^2\right]$$

$$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} \mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}}|s_h] \cdot \Pr[s_h|\tau_h]\right)^2\right]$$

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

|  | $V_{\mathcal{F}}^{\pi}$ | | $V_{\mathcal{S}}^{\pi}(s_h) = \mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h]$ |
|---|---|---|---|
| Candidate | $V_{\mathcal{F}}$ | | $V_{\mathcal{S}}(s_h) := \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_h)|s_h]$ |
| Prediction | $J(\pi) \approx \mathbb{E}_{\pi_b}[V_{\mathcal{F}}(f_1)]$ | $=$ | $\mathbb{E}[V_{\mathcal{S}}(s_1)]$ |
| Error | $\sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h}\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}]$ $\quad\downarrow$ $V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$ | $=$ | $\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$ |
| Learning objective | $\sum_{h=1}^{H} \mathbb{E}_{s_h\sim\pi_b}\left[\mathbb{E}_{\substack{a_h\sim g(s_h) \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|s_h]^2\right]$  ✗ | $=$ | $\mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi}V_{\mathcal{S}})(s_h))^2]$ |

$$\sum_{h=1}^{H} \mathbb{E}_{\tau_h\sim\pi_b}\left[\mathbb{E}_{\substack{a_h\sim\pi \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|\tau_h]^2\right]$$

$$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h\sim\pi_b}\left[\left(\sum_{s_h} \mathbb{E}_{\substack{a_h\sim g(s_h) \\ a_{h+1:H}\sim\pi_b}}[\Delta_h V_{\mathcal{F}}|s_h] \cdot \Pr[s_h|\tau_h]\right)^2\right]$$

belief state

linear measure

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?



$S = 2$

$\sum_{h=1}^{H} \mathbb{E}_{\pi} [V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$

Learning objective $\sum_{h=1}^{H} \mathbb{E}_{s_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | s_h]^2 \right] = \mathbb{E}_{\pi_b} [(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi} V_{\mathcal{S}})(s_h))^2]$

**X**

$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | \tau_h]^2 \right]$

belief state

$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \left( \sum_{s_h} \mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | s_h] \cdot \Pr[s_h | \tau_h] \right)^2 \right]$

linear measure

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?
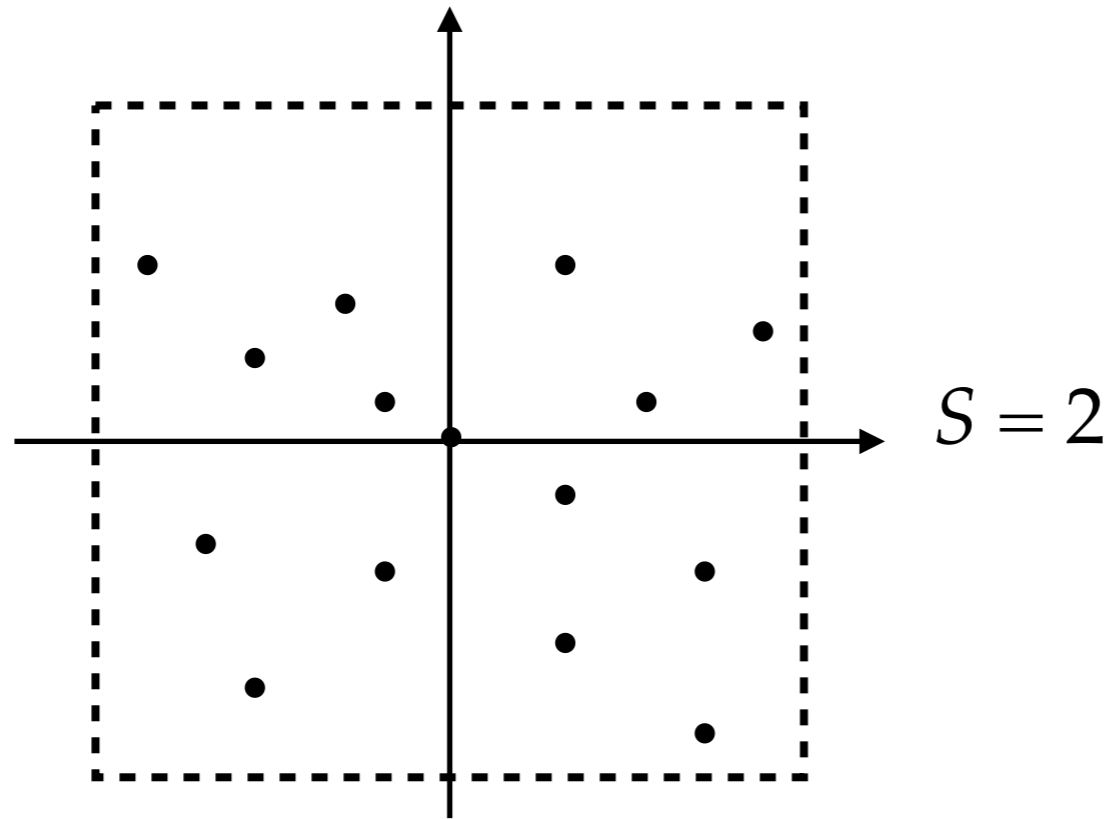


$\Pr[s_h|\tau_h]$

$S = 2$

$\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$

Learning objective $\sum_{h=1}^{H} \mathbb{E}_{s_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}}|s_h]^2 \right]$ ✗ $= \mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi} V_{\mathcal{S}})(s_h))^2]$
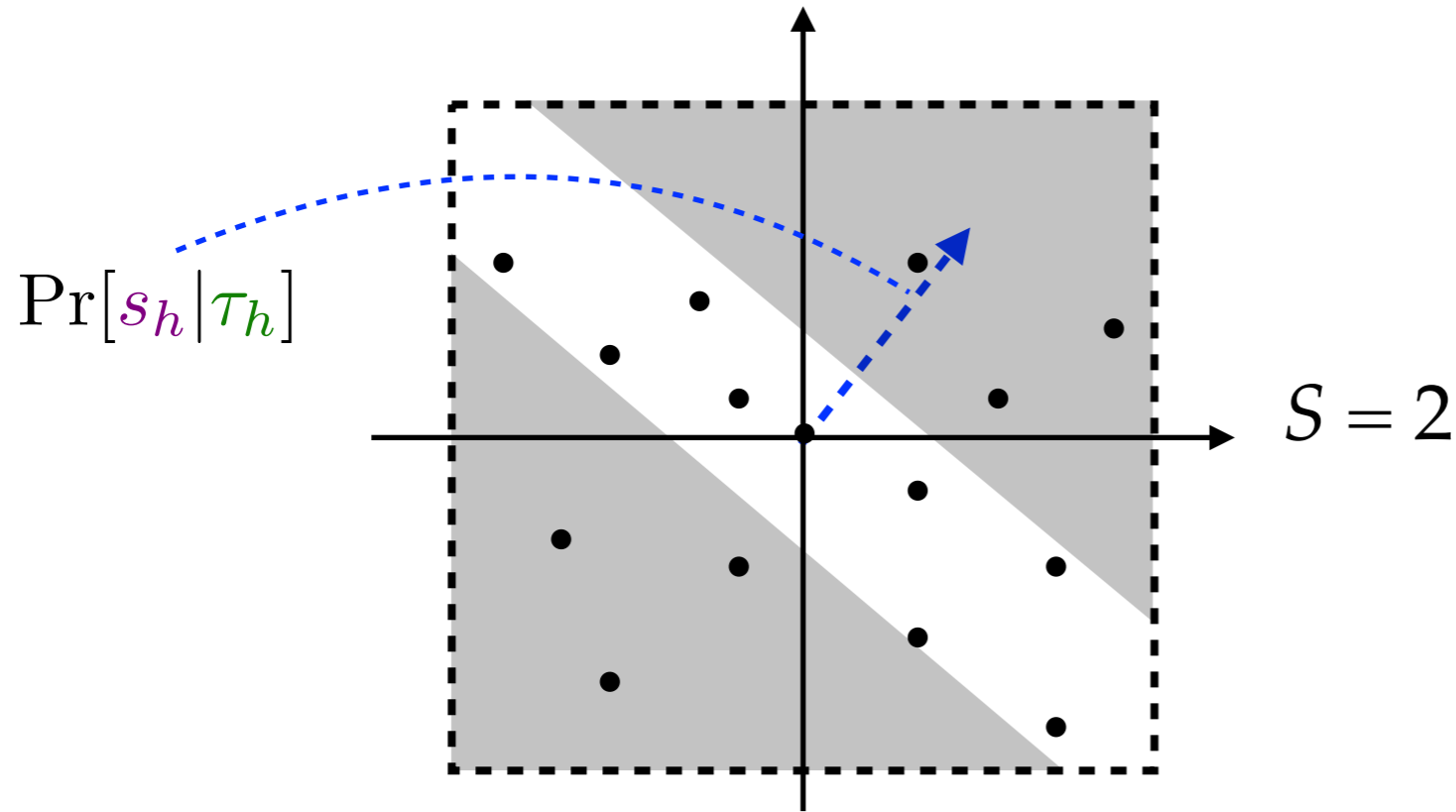
$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}}|\tau_h]^2 \right]$

belief state

$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \left( \sum_{s_h} \mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}}|s_h] \cdot \Pr[s_h|\tau_h] \right)^2 \right]$

linear measure

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

$\Pr[s_h | \tau_h]$

$S = 2$

$\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$

Learning objective $\sum_{h=1}^{H} \mathbb{E}_{s_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | s_h]^2 \right] = \mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi} V_{\mathcal{S}})(s_h))^2]$
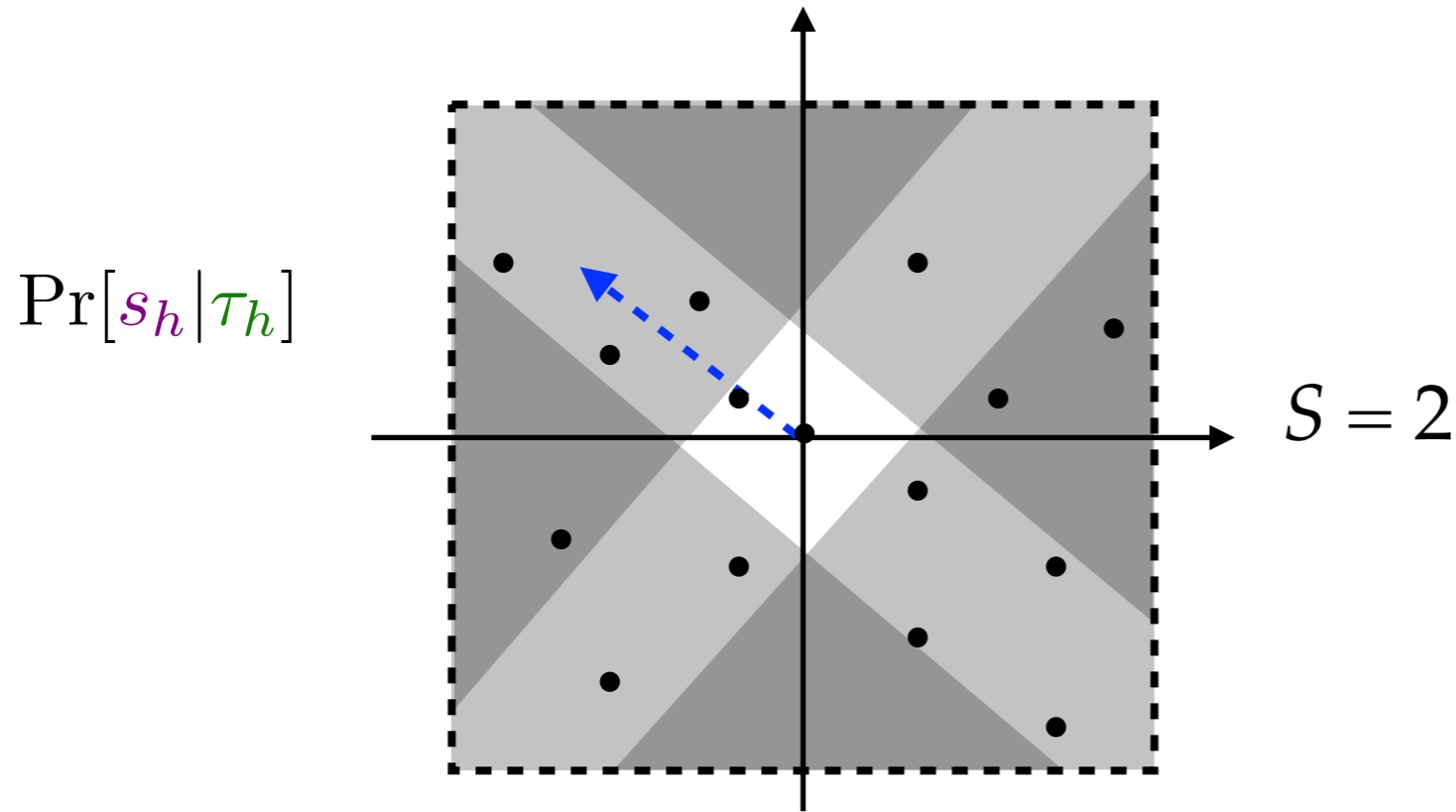
**X**

$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | \tau_h]^2 \right]$

belief state

$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \left( \sum_{s_h} \mathbb{E}_{\substack{a_h \sim g(s_h) \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | s_h] \cdot \Pr[s_h | \tau_h] \right)^2 \right]$

linear measure

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?

$\Pr[s_h | \tau_h]$

$S = 2$

$$\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$$

$$= \sum_{h=1}^{H} \sum_{s_h} g(s_h) d^{\pi}(s_h)$$

Learning objective $\sum_{h=1}^{H} \mathbb{E}_{s_h \sim \pi_b}\left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} g(s_h) [\Delta_h V_{\mathcal{F}} | s_h]^2\right] = \mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi} V_{\mathcal{S}})(s_h))^2]$
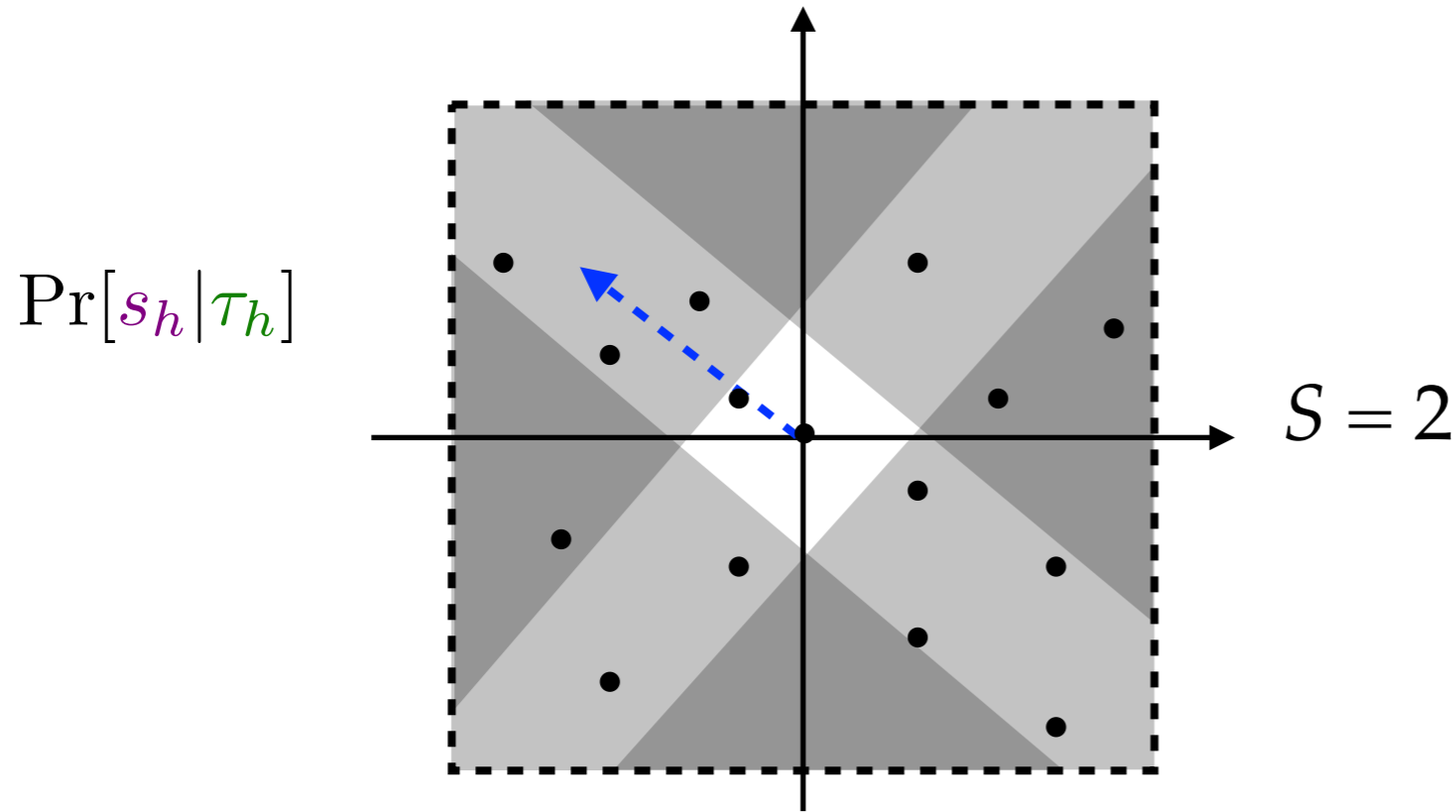
✗

$$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b}\left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}}[\Delta_h V_{\mathcal{F}} | \tau_h]^2\right]$$

belief state

$$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} g(s_h)[\Delta_h V_{\mathcal{F}} | s_h] \cdot \Pr[s_h | \tau_h]\right)^2\right]$$

linear measure

# Does it work in the same way as $V_{\mathcal{S}}^{\pi}$?
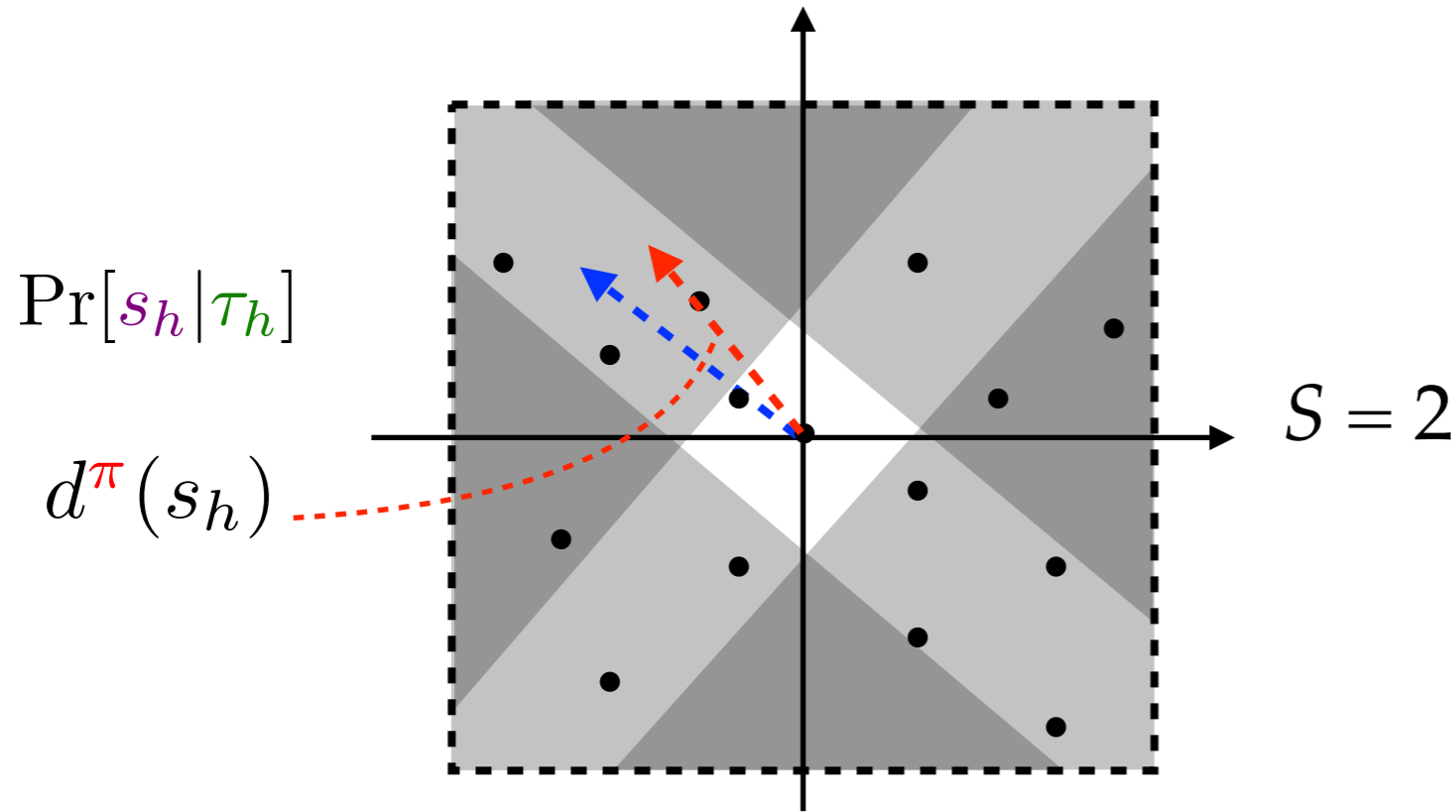
$\Pr[s_h|\tau_h]$

$d^{\pi}(s_h)$

$S = 2$

$\sum_{h=1}^{H} \mathbb{E}_{\pi}[V_{\mathcal{S}}(s_h) - r_h - V_{\mathcal{S}}(s_{h+1})]$

$$= \sum_{h=1}^{H} \sum_{s_h} g(s_h) d^{\pi}(s_h)$$

Learning objective $\sum_{h=1}^{H} \mathbb{E}_{s_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} g(s_h) \Delta_h V_{\mathcal{F}}|s_h]^2 \right] = \mathbb{E}_{\pi_b}[(V_{\mathcal{S}}(s_h) - (\mathcal{T}^{\pi} V_{\mathcal{S}})(s_h))^2]$ **X**

$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}}|\tau_h]^2 \right]$

belief state

$$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \left( \sum_{s_h} \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} g(s_h) \Delta_h V_{\mathcal{F}}|s_h] \cdot \Pr[s_h|\tau_h] \right)^2 \right]$$

linear measure

$$= \sum_{h=1}^{H} \sum_{s_h} g(s_h) d^{\pi}(s_h)$$

cover

$$= \sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \left( \sum_{s_h} \underset{\substack{a_h \sim \\ a_{h+1:H} \sim \pi_b}}{\mathbb{E}} g(s_h) [\Delta_h V_{\mathcal{F}} | s_h] \cdot \Pr[s_h | \tau_h] \right)^2 \right]$$

linear measure

$$\sum_{h=1}^{H} \sum_{s_h} g(s_h) d^{\pi}(s_h)$$

$$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b} \left[ \left( \sum_{s_h} \mathbb{E}_{\substack{a_h \sim \\ a_{h+1:H} \sim \pi_b}} g(s_h)_h [\Delta_h V_{\mathcal{F}} | s_h] \cdot \Pr[s_h | \tau_h] \right)^2 \right]$$

cover

linear measure

11

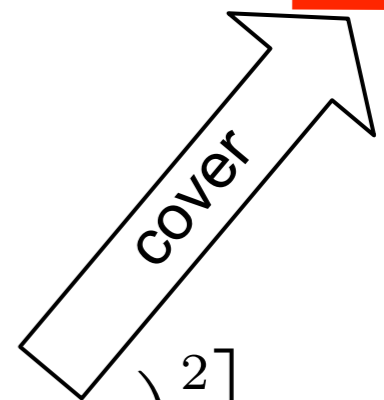**Theorem** (Informal): Assume

$$(d_h^{\pi})^{\top} \mathbb{E}_{\tau_h \sim \pi_b}[\mathbf{b}(\tau_h)\mathbf{b}(\tau_h)^{\top}]^{-1} d_h^{\pi} \leq C_{\mathcal{H}}$$

$$\sum_{h=1}^{H} \sum_{s_h} g(s_h) d^{\pi}(s_h)$$

cover

$$\sum_{h=1}^{H} \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} g(s_h) \cdot \Pr[s_h|\tau_h]\right)^2\right]$$

linear measure

belief state $\mathbf{b}(\tau_h)$

11

**Theorem** (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b}[\mathbf{b}(\tau_h)\mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \le C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- Coverage parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$

- Ranges & complexities of function classes (e.g., that of $\mathcal{V}$)

$$\sum_{h=1}^H \sum_{s_h} g(s_h)\underline{d^\pi(s_h)}$$

cover

$$\sum_{h=1}^H \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} \underset{\substack{a_h \sim \\ a_{h+1:H} \sim \pi_b}}{\mathbb{E}} g(s_h) [\Delta]_h V_{\mathcal{F}}|s_h] \cdot \Pr[s_h|\tau_h]\right)^2\right]$$

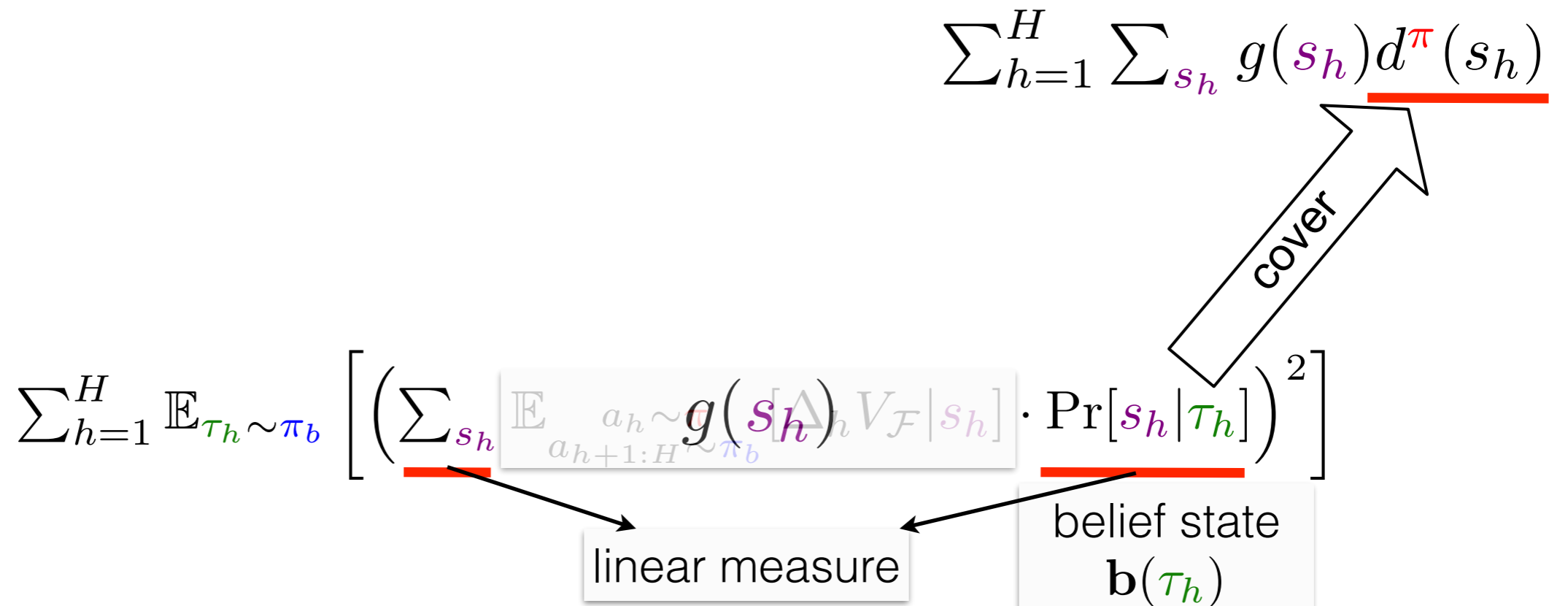linear measure     belief state $\mathbf{b}(\tau_h)$

11

**Theorem** (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b}[\mathbf{b}(\tau_h)\mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \le C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- Coverage parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \dfrac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$

- Ranges & complexities of function classes (e.g., that of $\mathcal{V}$)

- Similar to $\mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1}\mathbb{E}_\pi[\phi]$

$$\mathbb{E}_{\tau_h \sim \pi}[\Pr[s_h|\tau_h]]$$

$$\sum_{h=1}^H \sum_{s_h} g(s_h)d^\pi(s_h)$$

cover

$$\sum_{h=1}^H \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} \mathbb{E}_{\substack{a_h \sim \\ a_{h+1:H} \sim \pi_b}} g(s_h) V_{\mathcal{F}}|s_h] \cdot \Pr[s_h|\tau_h]\right)^2\right]$$

linear measure

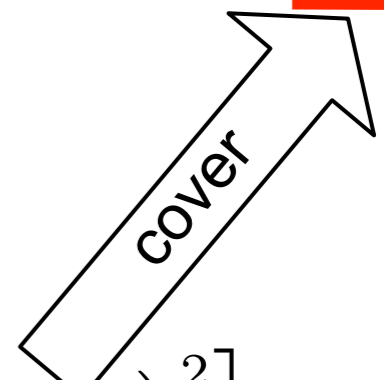belief state $\mathbf{b}(\tau_h)$

**Theorem** (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b}[\mathbf{b}(\tau_h)\mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_\mathcal{H}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- Coverage parameters: $C_\mathcal{H}$ and $C_\mathcal{A} := \max\limits_{s_h, a_h} \dfrac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$

- Ranges & complexities of function classes (e.g., that of $\mathcal{V}$)

- Similar to $\mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1}\mathbb{E}_\pi[\phi]$
- $\pi = \pi_b$ : $C_\mathcal{H}$ =1

$$\mathbb{E}_{\tau_h \sim \pi}[\Pr[s_h|\tau_h]]$$

$$\sum_{h=1}^H \sum_{s_h} g(s_h)\underline{d^\pi(s_h)}$$

cover

$$\sum_{h=1}^H \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} \underset{\substack{a_h \sim \\ a_{h+1:H} \sim \pi_b}}{\mathbb{E}}[g(s_h)\phi_h V_\mathcal{F}|s_h] \cdot \Pr[s_h|\tau_h]\right)^2\right]$$

linear measure

belief state
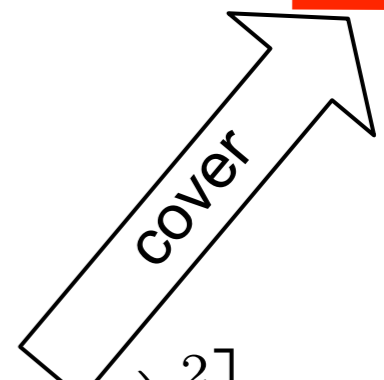$\mathbf{b}(\tau_h)$

11

**Theorem** (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b}[\mathbf{b}(\tau_h)\mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_\mathcal{H}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- Coverage parameters: $C_\mathcal{H}$ and $C_\mathcal{A} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$

- Ranges & complexities of function classes (e.g., that of $\mathcal{V}$)

- Similar to $\mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$
- $\pi = \pi_b : C_\mathcal{H} = 1$
- 1-hot $\mathbf{b}(\tau_h)$: $\mathbb{E}_{\pi_b}\left[(d_h^\pi/d_h^{\pi_b})^2\right]$

$$\mathbb{E}_{\tau_h \sim \pi}[\Pr[s_h|\tau_h]]$$

$$\sum_{h=1}^H \sum_{s_h} g(s_h) d^\pi(s_h)$$

cover

$$\sum_{h=1}^H \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} \mathbb{E}_{\substack{a_h \sim \\ a_{h+1:H} \sim \pi_b}} g(s_h)_h V_\mathcal{F}|s_h] \cdot \Pr[s_h|\tau_h]\right)^2\right]$$

linear measure

belief state
$\mathbf{b}(\tau_h)$

**Theorem** (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b}[\mathbf{b}(\tau_h)\mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in
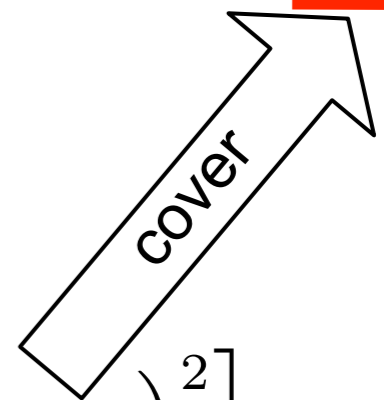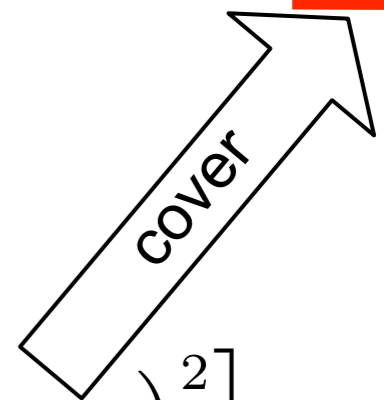
- Coverage parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$

- Ranges & complexities of function classes (e.g., that of $\mathcal{V}$)

- Similar to $\mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$
- $\pi = \pi_b : C_{\mathcal{H}} = 1$
- 1-hot $\mathbf{b}(\tau_h)$: $\mathbb{E}_{\pi_b}[(d_h^\pi/d_h^{\pi_b})^2]$
- $\|d_h^\pi/d_h^{\pi_b}\|_\infty \Rightarrow \|\mathbb{E}_{\pi_b}[\mathbf{b}(\tau_h)\mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi\|_\infty$

$$\mathbb{E}_{\tau_h \sim \pi}[\Pr[s_h|\tau_h]]$$

$$\downarrow$$

$$\sum_{h=1}^H \sum_{s_h} g(s_h) d^\pi(s_h)$$

cover

$$\sum_{h=1}^H \mathbb{E}_{\tau_h \sim \pi_b}\left[\left(\sum_{s_h} \mathbb{E}_{\substack{a_h \sim \\ a_{h+1:H} \sim \pi_b}} g(s_h)[\cdot]_h V_{\mathcal{F}}|s_h] \cdot \Pr[s_h|\tau_h]\right)^2\right]$$
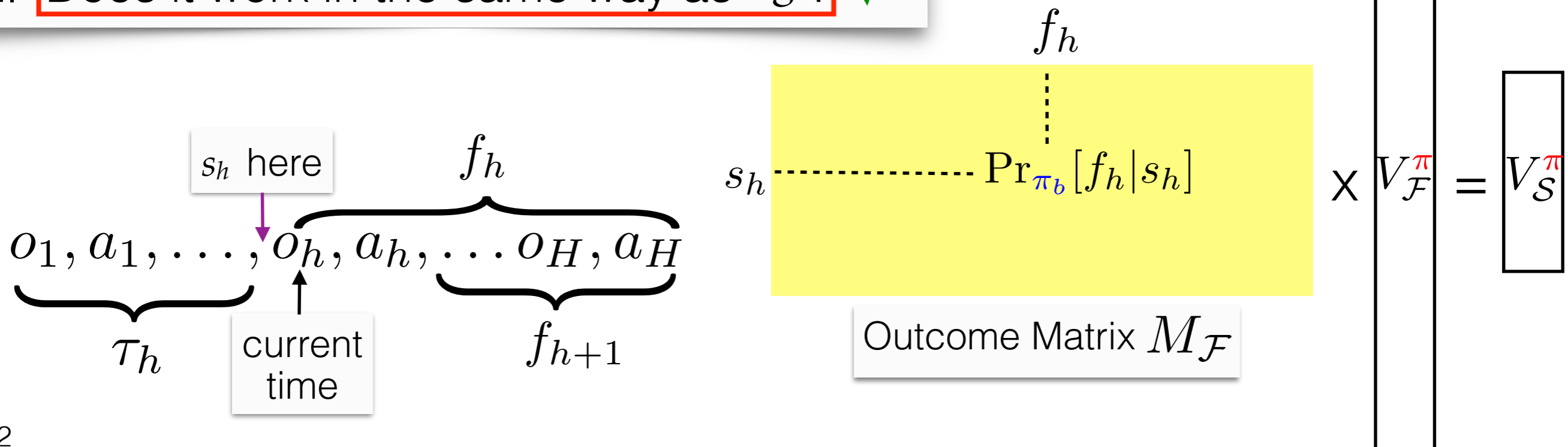
linear measure

belief state
$\mathbf{b}(\tau_h)$

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}|s_h] \in [0, H]$$

  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

  - $\mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

1. Does (well-behaved) $V_{\mathcal{F}}^{\pi}$ even exist?
2. Does it work in the same way as $V_{\mathcal{S}}^{\pi}$? ✓

$$f_h$$

$$o_1, a_1, \ldots, o_h, a_h, \ldots o_H, a_H$$

$s_h$ here

$f_h$

$\tau_h$

current time

$f_{h+1}$

$s_h \text{-------------} \Pr_{\pi_b}[f_h|s_h]$

Outcome Matrix $M_{\mathcal{F}}$

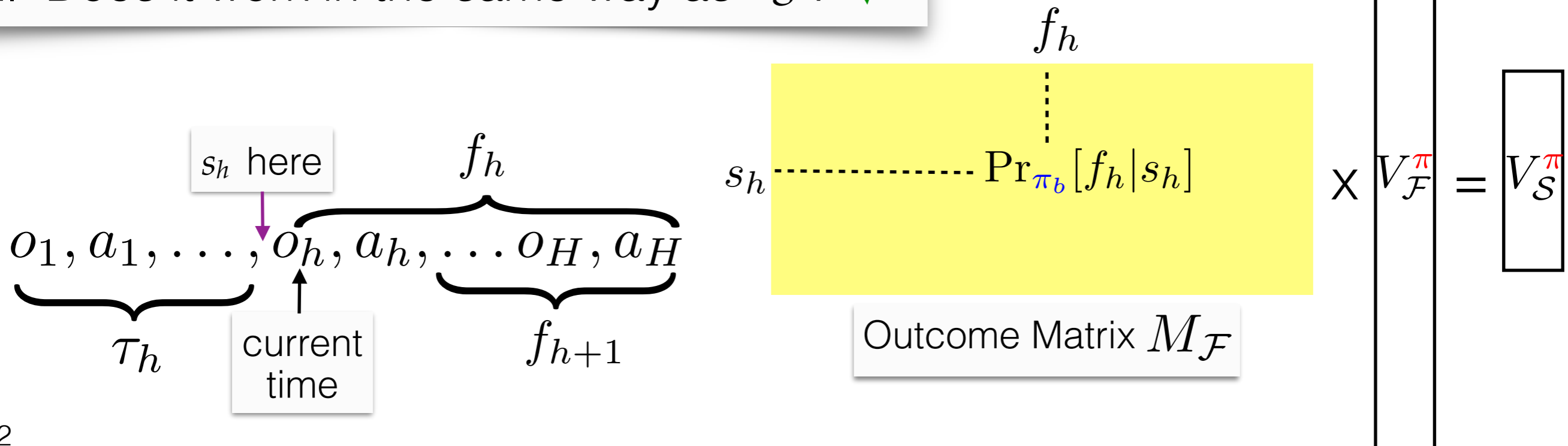$\times \, |V_{\mathcal{F}}^{\pi}| = |V_{\mathcal{S}}^{\pi}|$

# Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi}[\sum_{h'=h}^{H} r_{h'}|s_h] \in [0, H]$$
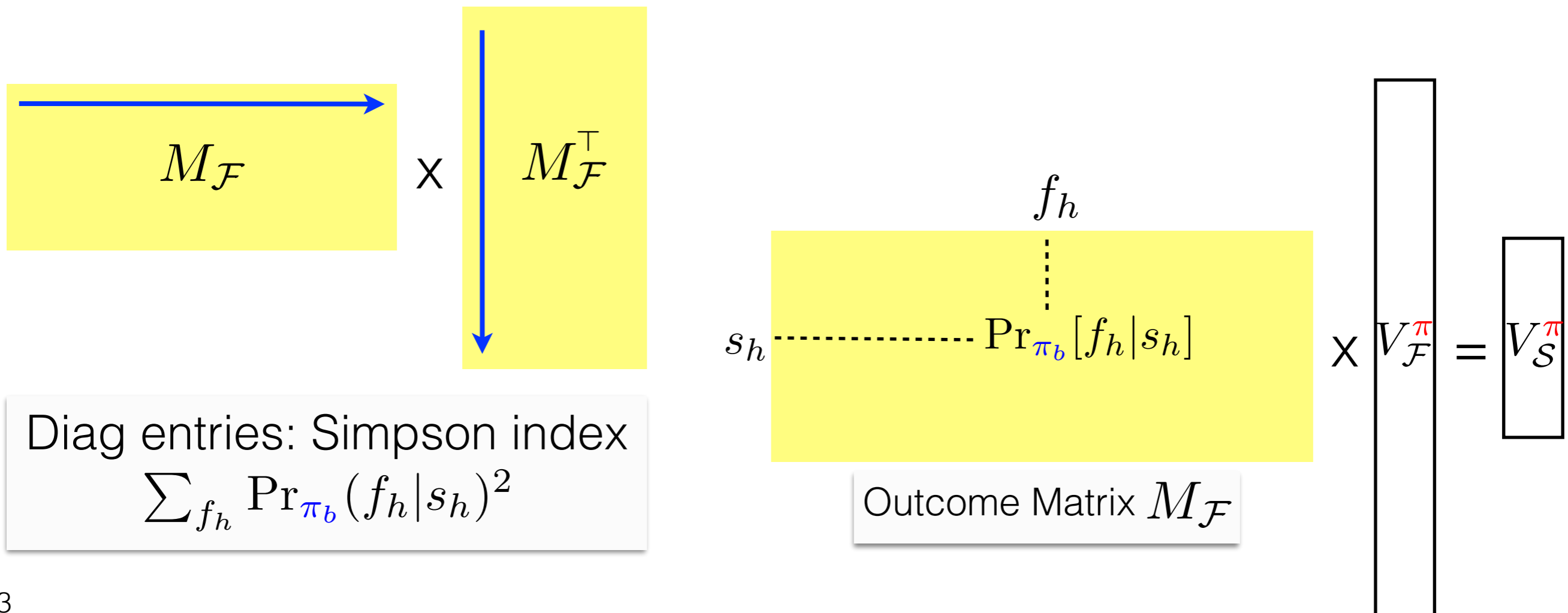
  - Problem: $s_h$ is latent — can't even use this function!

- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using *future* as input!

  - $\mathbb{E}_{\pi_b}[V_{\mathcal{F}}^{\pi}(f_h)|s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

1. Does (well-behaved) $V_{\mathcal{F}}^{\pi}$ even exist?
2. Does it work in the same way as $V_{\mathcal{S}}^{\pi}$? ✓

$$f_h$$

$$s_h \text{ here} \qquad f_h$$

$$o_1, a_1, \ldots, o_h, a_h, \ldots o_H, a_H$$

$$\tau_h \qquad \begin{array}{c}\text{current}\\\text{time}\end{array} \qquad f_{h+1}$$

$$s_h \text{----------} \Pr_{\pi_b}[f_h|s_h]$$

Outcome Matrix $M_{\mathcal{F}}$

$$\times \begin{vmatrix} V_{\mathcal{F}}^{\pi} \end{vmatrix} = \begin{vmatrix} V_{\mathcal{S}}^{\pi} \end{vmatrix}$$

12
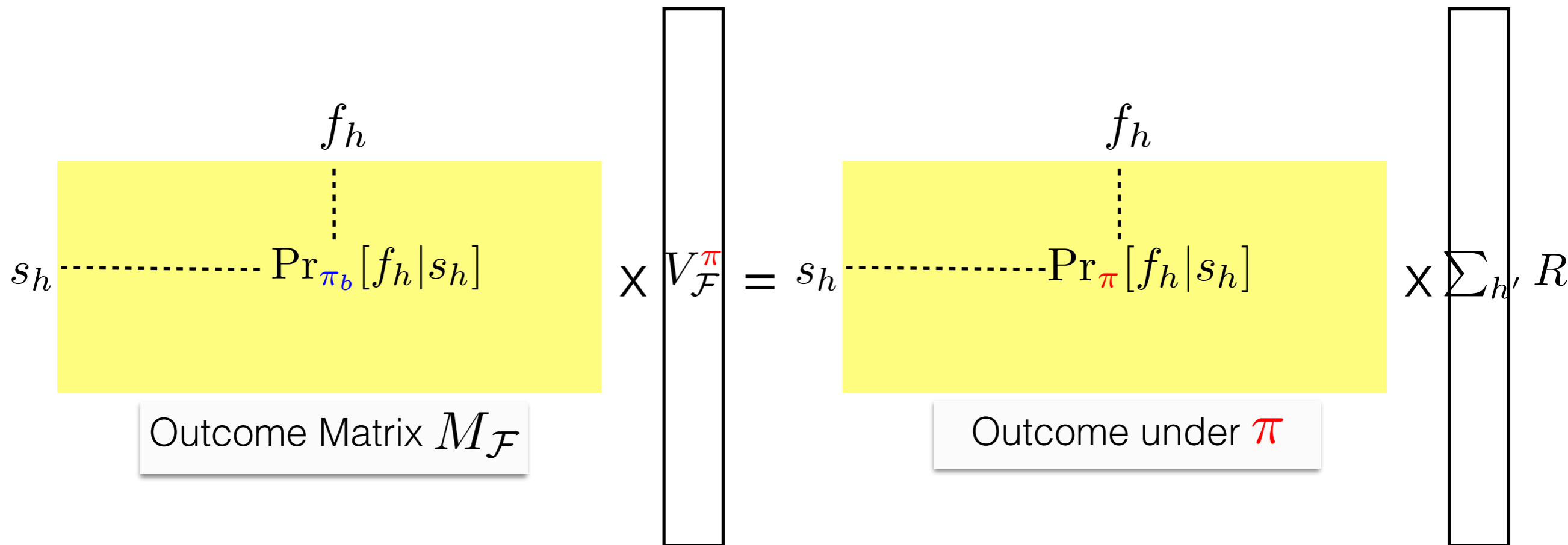
# Constructing $V_{\mathcal{F}}^{\pi}$: pseudo-inverse

- Pay $1/\sigma_{\min}$ of outcome matrix, which is 1/min-eigenvalue of

  - Exponentially small when system is stochastic!
  - Problem: "linear regression" with exp. small covariates
  - Similar quantity appears in online RL (Liu et al'22)



$M_{\mathcal{F}}$  x  $M_{\mathcal{F}}^{\top}$

Diag entries: Simpson index
$\sum_{f_h} \mathrm{Pr}_{\pi_b}(f_h|s_h)^2$

$f_h$

$s_h \cdots\cdots\cdots\cdots \mathrm{Pr}_{\pi_b}[f_h|s_h]$

Outcome Matrix $M_{\mathcal{F}}$

x $V_{\mathcal{F}}^{\pi}$ = $V_{\mathcal{S}}^{\pi}$

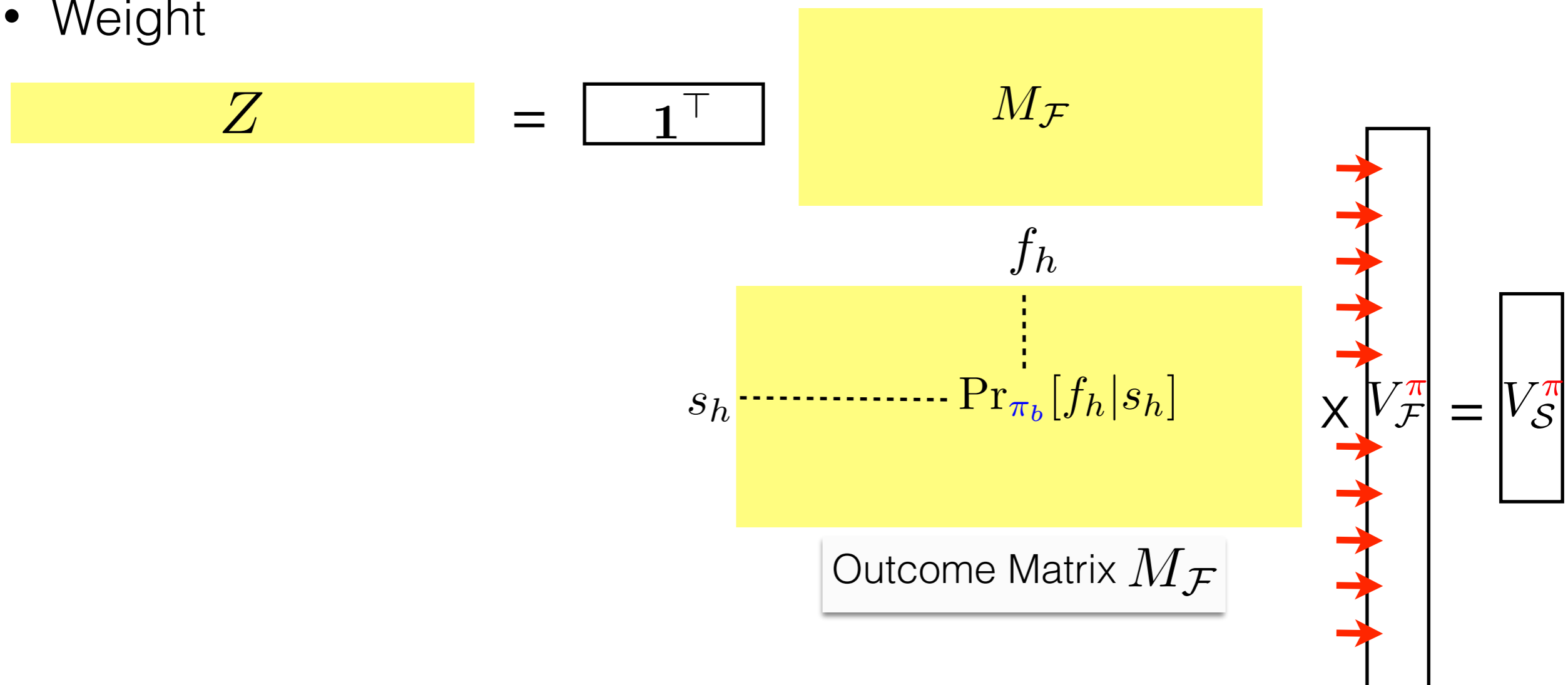# Constructing $V_{\mathcal{F}}^{\pi}$: reweighting

- General case

  **!!**

  - $V_{\mathcal{F}}^{\pi}(f_h) = \displaystyle\prod_{h'=h}^{H} \frac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})} \left( \sum_{h'=h}^{H} R(o_{h'}, a_{h'}) \right)$



$f_h$

$s_h \text{-----} \Pr_{\pi_b}[f_h|s_h]$

Outcome Matrix $M_{\mathcal{F}}$

$\times \boxed{V_{\mathcal{F}}^{\pi}} = s_h \text{-----} \Pr_{\pi}[f_h|s_h]$

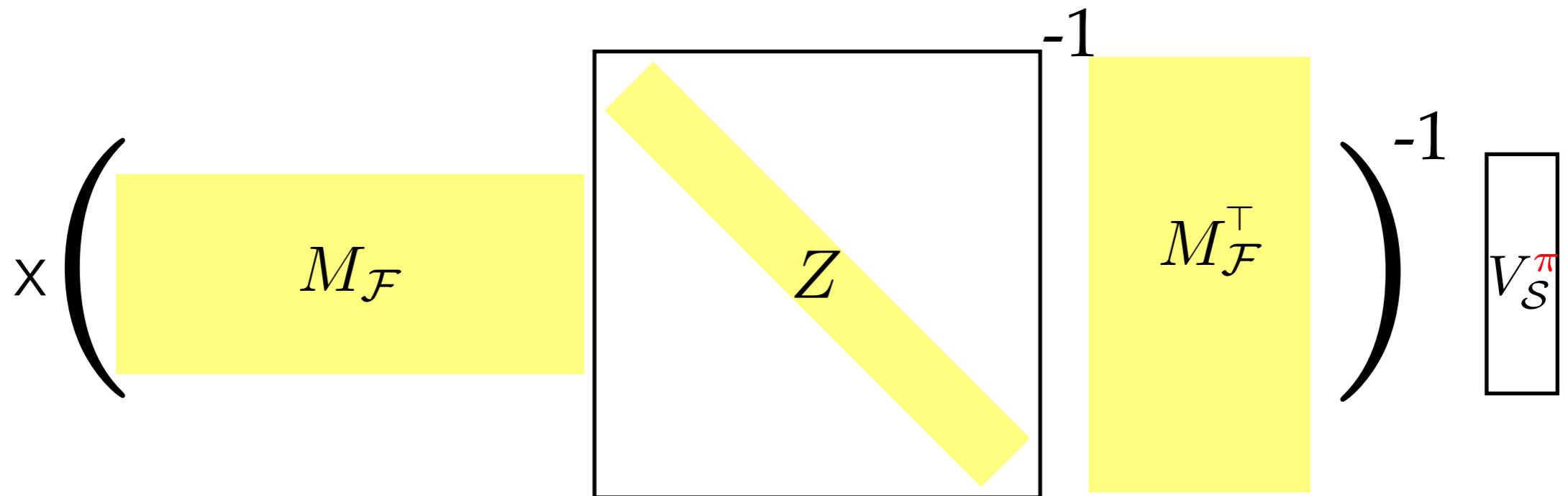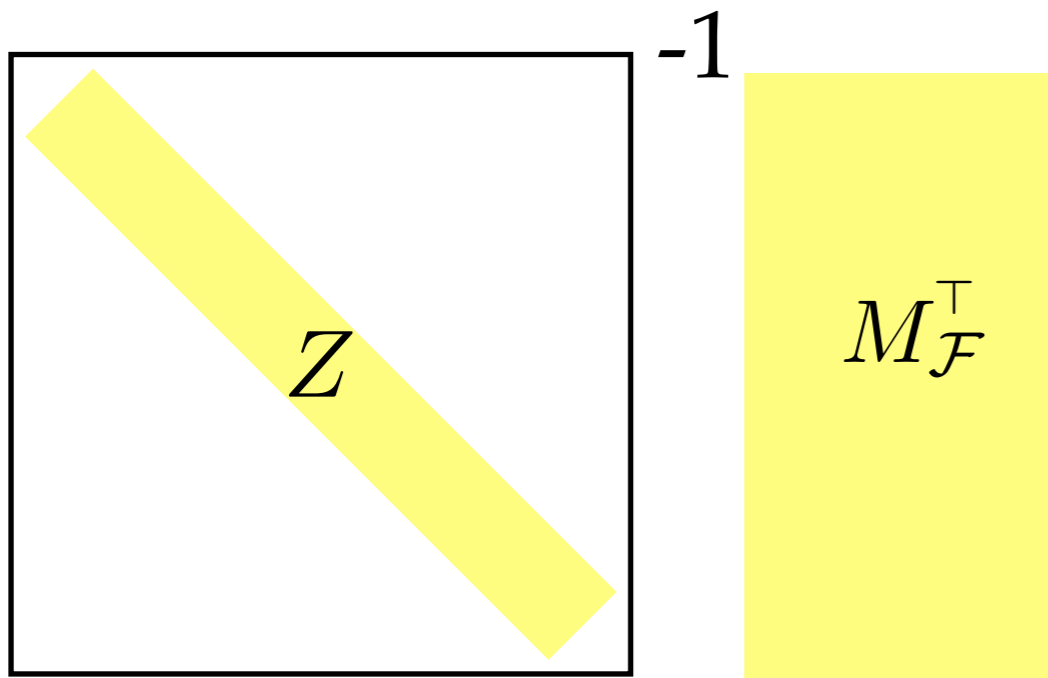Outcome under $\pi$

$\times \sum_{h'} R$

# Constructing $V_{\mathcal{F}}^{\pi}$: *weighted* pseudo-inverse

- Pseudo-inverse minimizes $L_2$ norm (we want $L_\infty$)

- $L_2$ norm treats all exponentially many coordinates equally — not informative

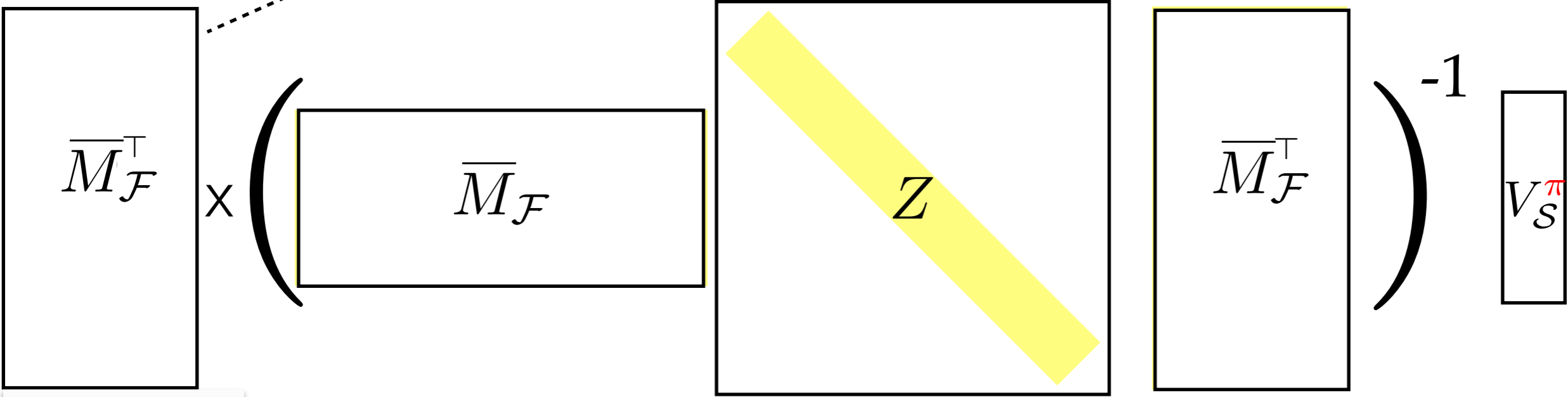- Find solution that minimizes *weighted $L_2$* norm

- Weight

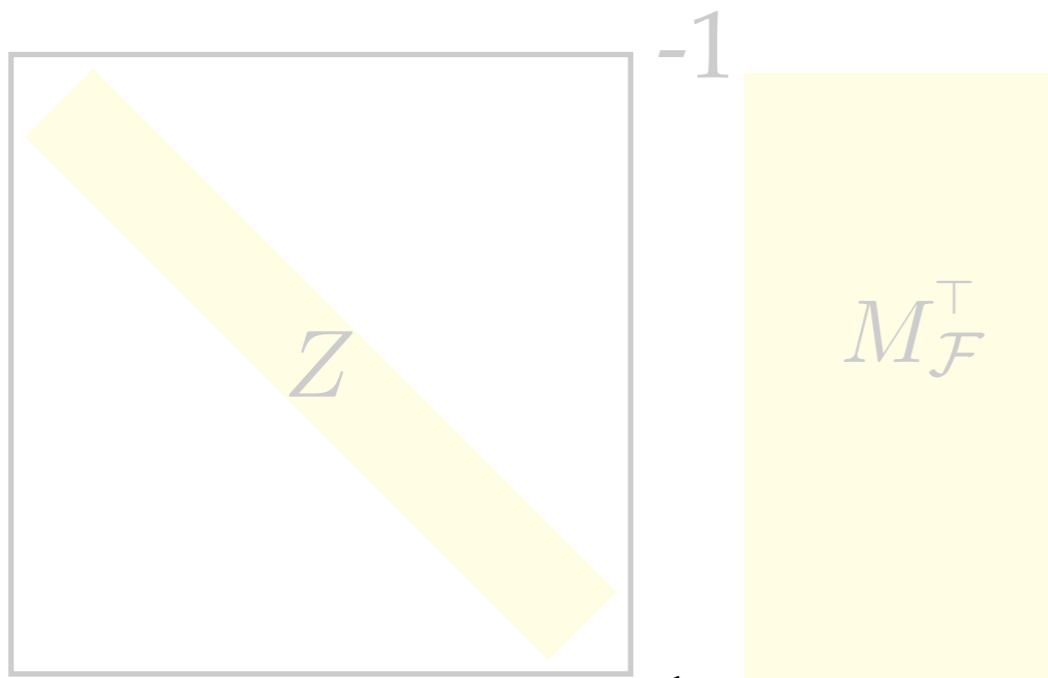$$\boxed{Z} = \boxed{\mathbf{1}^\top} \boxed{M_{\mathcal{F}}}$$

$$f_h$$
$$\vdots$$
$$s_h \text{-------------} \Pr_{\pi_b}[f_h | s_h]$$

Outcome Matrix $M_{\mathcal{F}}$

$$\times \; V_{\mathcal{F}}^{\pi} = V_{\mathcal{S}}^{\pi}$$

# Constructing $V_{\mathcal{F}}^{\pi}$: *weighted* pseudo-inverse



Col of $\overline{M}_{\mathcal{F}}$: posterior of $s_h$ given $f_h$ under uniform prior

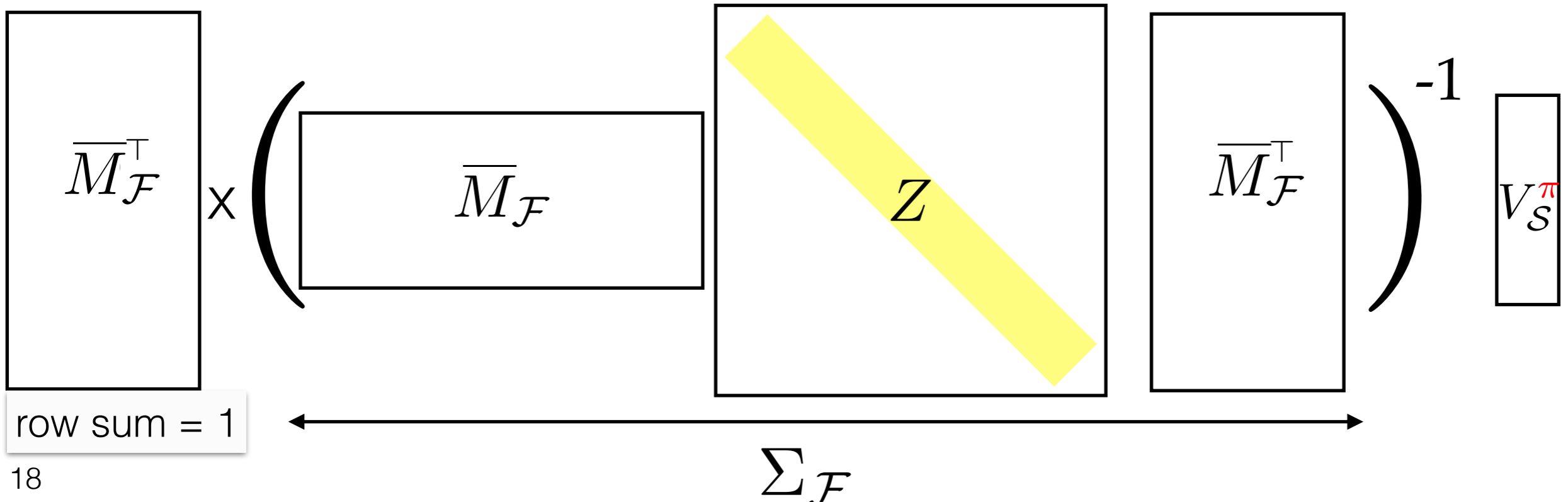# Constructing $V_{\mathcal{F}}^{\pi}$: *weighted* pseudo-inverse



Col of $\overline{M}_{\mathcal{F}}$: posterior of $s_h$ given $f_h$ under uniform prior
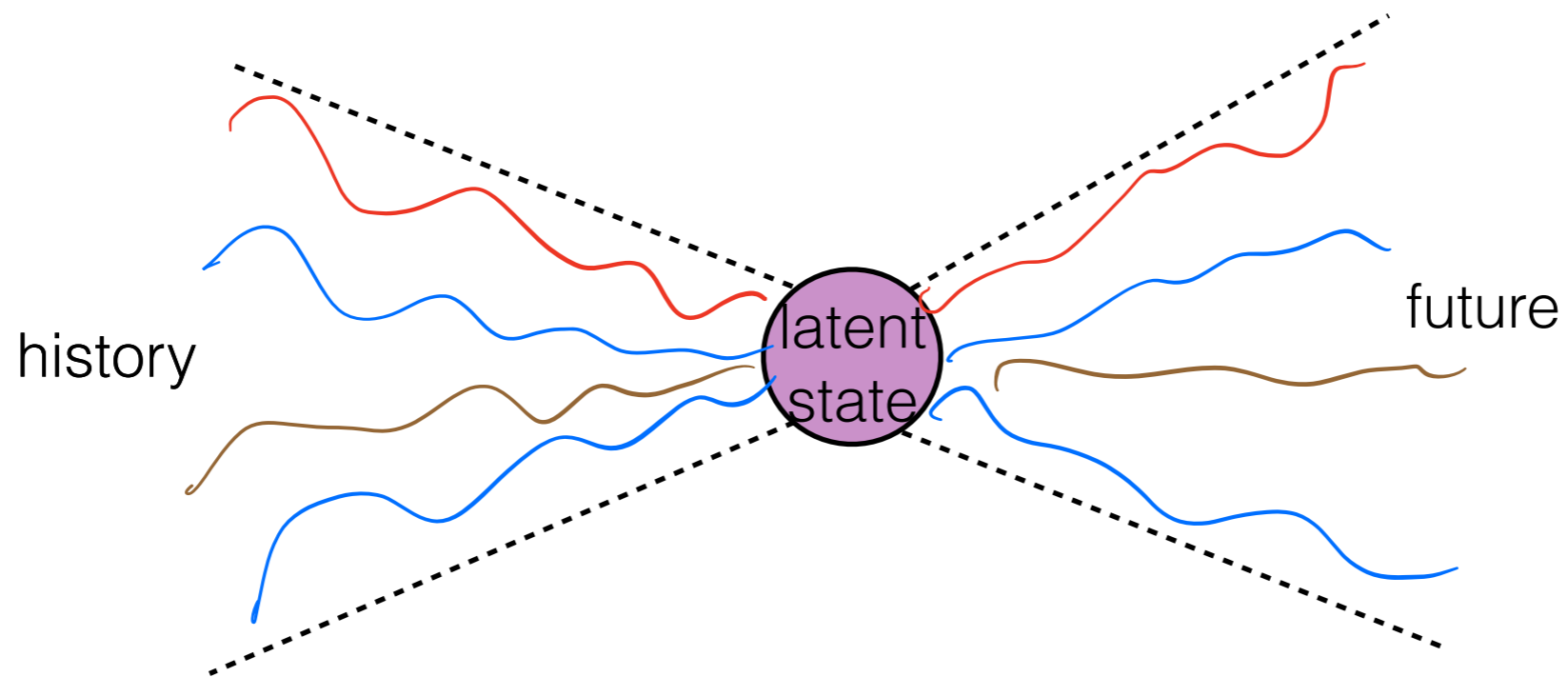
# Properties of $\Sigma_{\mathcal{F}}$

- $\Sigma_{\mathcal{F}}$ is doubly stochastic

- When $f_h$ reveals $s_h$, $\Sigma_{\mathcal{F}} = \mathbf{I}$

- More generally, confusion matrix of predicting $s_h$ from $f_h$

- **Outcome Coverage**: $\|\Sigma_{\mathcal{F}}^{-1} V_{\mathcal{S}}^{\pi}\|_{\infty} \leqslant C_{\mathcal{F}}$ (not $\|\Sigma_{\mathcal{F}}^{-1/2} V_{\mathcal{S}}^{\pi}\|_{2}$)

  - $\Rightarrow$ $\|V_{\mathcal{F}}^{\pi}\|_{\infty} \leqslant C_{\mathcal{F}}$



$$\overline{M}_{\mathcal{F}}^{\top} \times \left( \overline{M}_{\mathcal{F}} \quad Z \quad \overline{M}_{\mathcal{F}}^{\top} \right)^{-1} V_{\mathcal{S}}^{\pi}$$
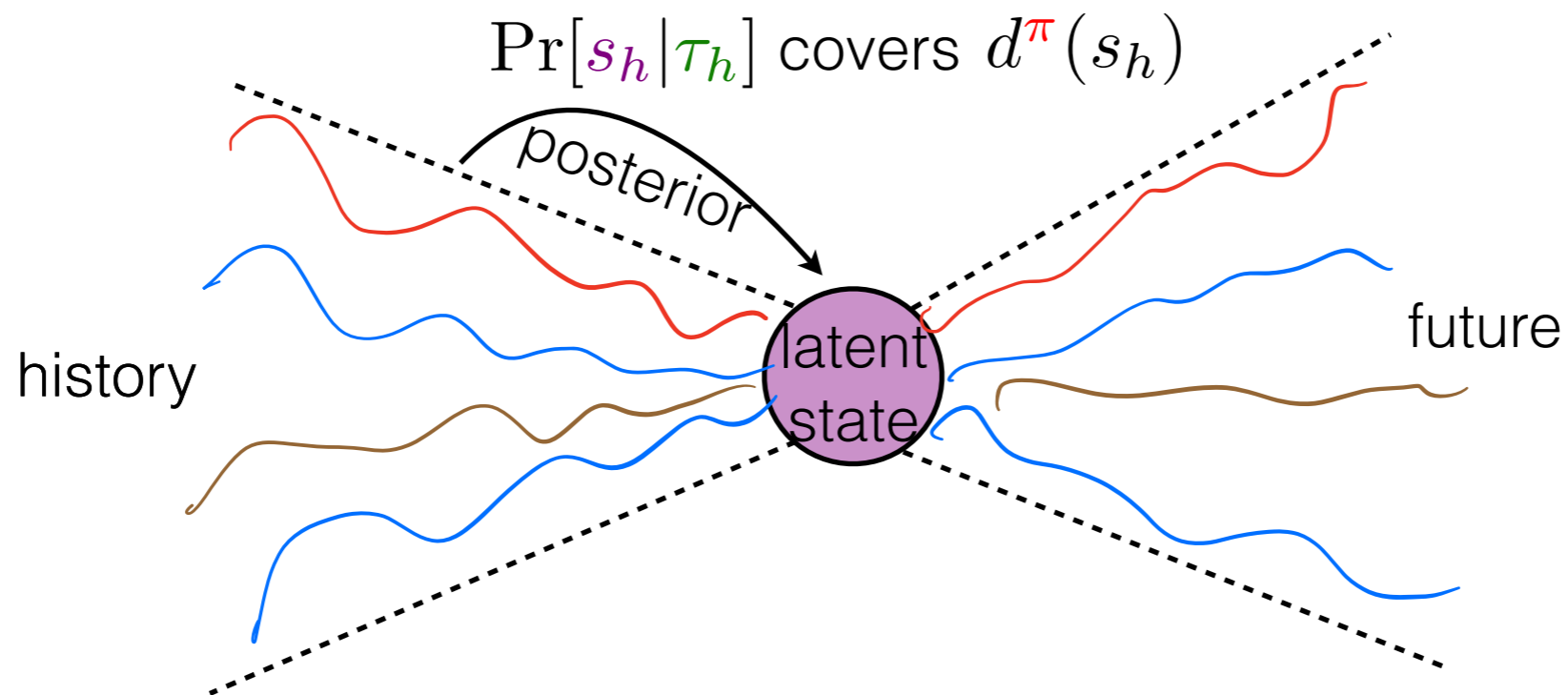
row sum = 1

$\Sigma_{\mathcal{F}}$

# Conclusion

- Problem: OPE in POMDPs

- New framework: future-dependent value function $V_{\mathcal{F}}^{\pi}$

# Conclusion

- **Problem**: OPE in POMDPs

- New **framework**: **future-dependent** value function $V_{\mathcal{F}}^{\pi}$

- New **assumptions**:

  - **Belief coverage =>** error transfer

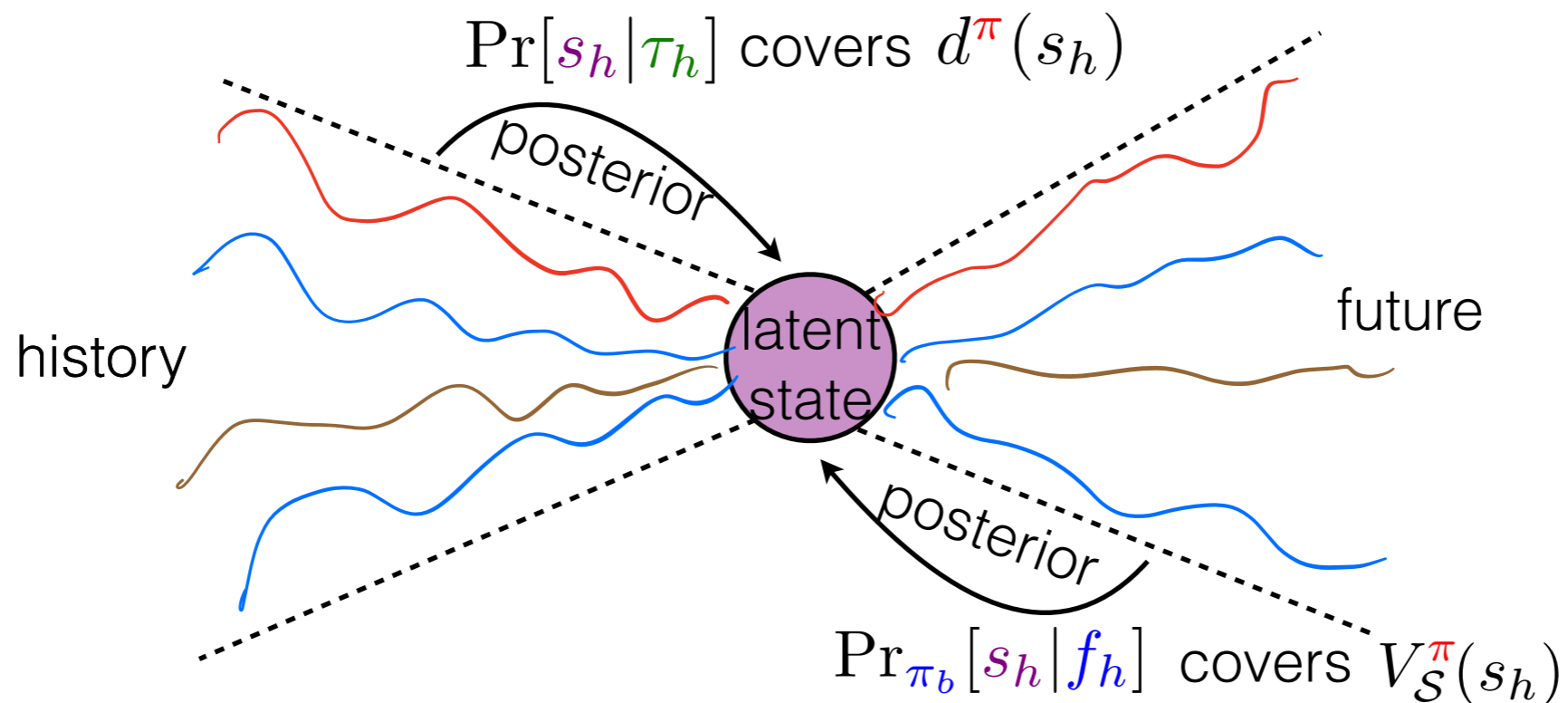$$\Pr[s_h | \tau_h] \text{ covers } d^{\pi}(s_h)$$

# Conclusion

- **Problem**: OPE in POMDPs

- New **framework**: **future-dependent** value function $V_{\mathcal{F}}^{\pi}$

- New **assumptions**:

  - **Belief coverage** => error transfer

  - **Outcome coverage** => bounded $V_{\mathcal{F}}^{\pi}$

$\Pr[s_h | \tau_h]$ covers $d^{\pi}(s_h)$

posterior

history

latent state

future

posterior
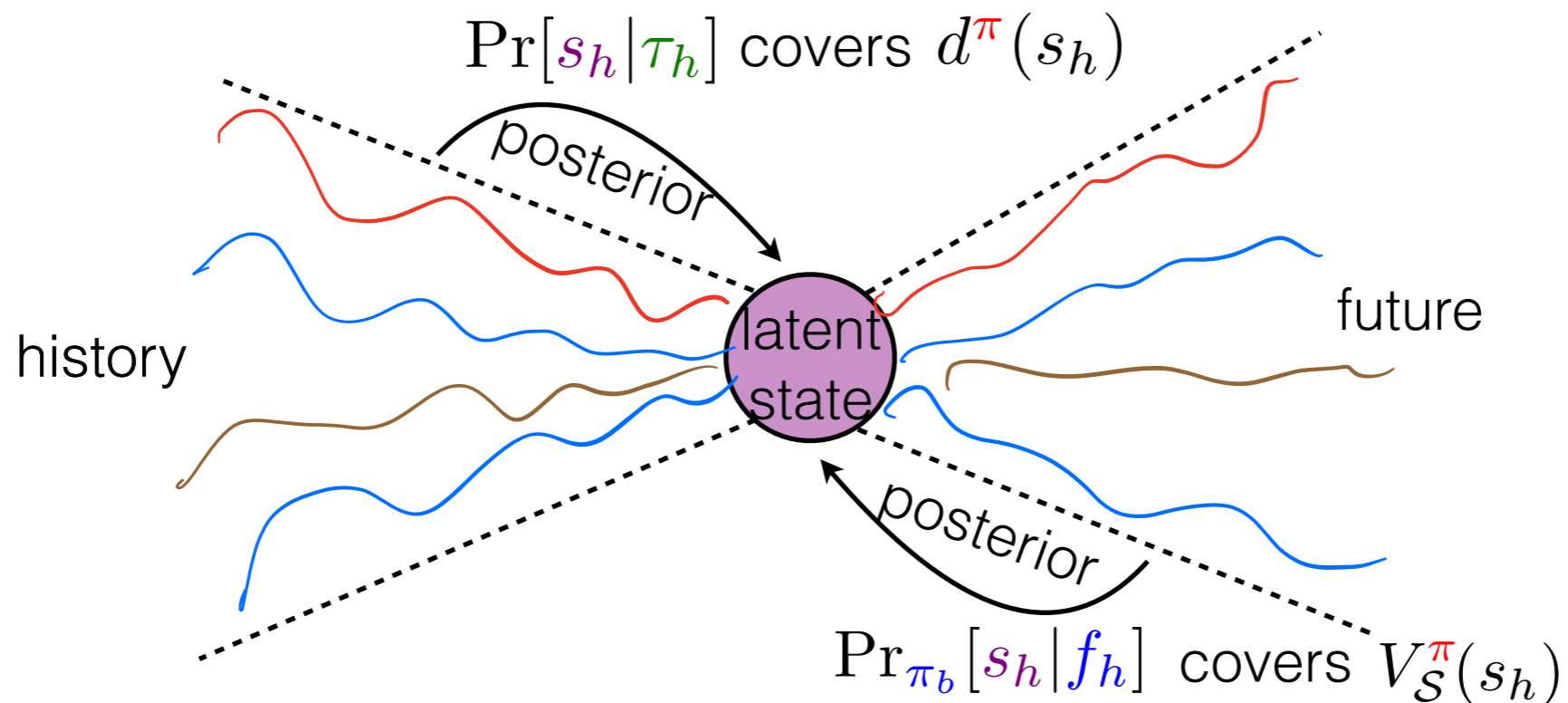
$\Pr_{\pi_b}[s_h | f_h]$ covers $V_{\mathcal{S}}^{\pi}(s_h)$

# Conclusion

- Problem: OPE in POMDPs

- New framework: future-dependent value function $V_{\mathcal{F}}^{\pi}$

- New assumptions:

  - **Belief coverage =>** error transfer

  - **Outcome coverage =>** bounded $V_{\mathcal{F}}^{\pi}$

- Open question: beyond memoryless & FSM policies

# Conclusion

Masatoshi Uehara

Yuheng Zhang

- **Problem**: OPE in POMDPs

- New framework: future-dependent value function $V_{\mathcal{F}}^{\pi}$

- New assumptions:

  - **Belief coverage => error transfer**

  - **Outcome coverage => bounded $V_{\mathcal{F}}^{\pi}$**

- Open question: beyond memoryless & FSM policies

**Thank you! Questions?**



$\Pr[s_h|\tau_h]$ covers $d^{\pi}(s_h)$

posterior

history

latent state

future

posterior

$\Pr_{\pi_b}[s_h|f_h]$ covers $V_{\mathcal{S}}^{\pi}(s_h)$