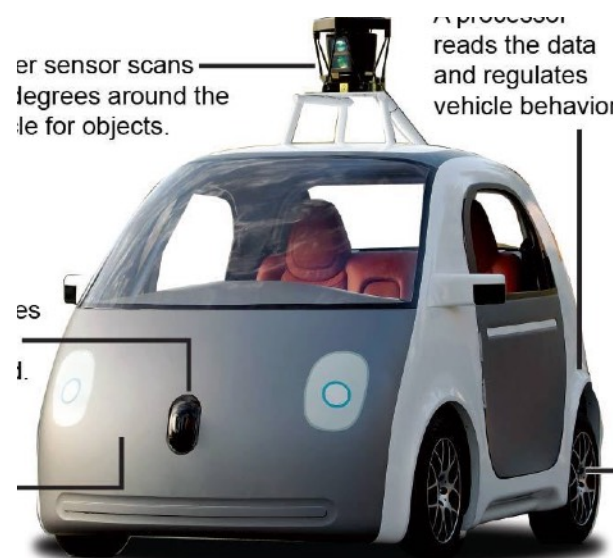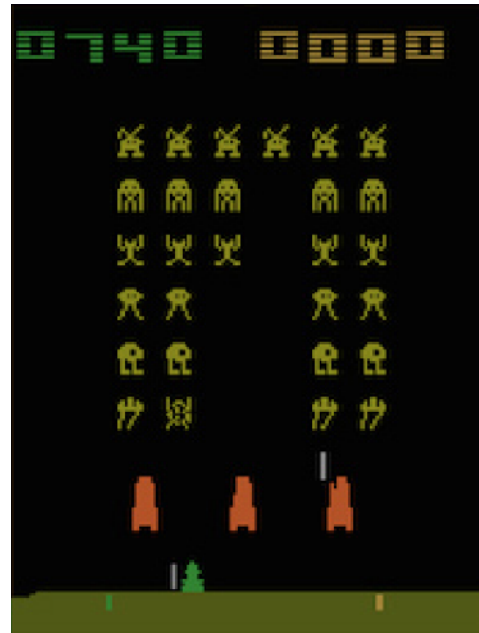# PAC-Exploration in Reinforcement Learning with Function Approximation

**Nan Jiang** (University of Michigan),
Akshay Krishnamurthy (University of Massachusetts Amherst),
Alekh Agarwal, John Langford, Robert E. Schapire (Microsoft Research)

1

# RL applications

# Key aspects of RL

*Bellman Equations*

(Dynamic Programming)

Temporal credit assignment

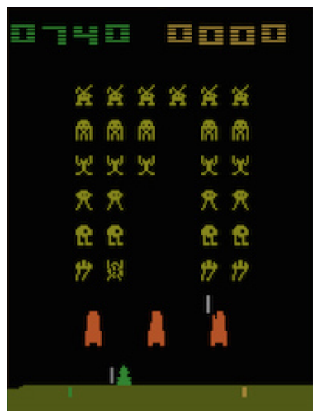Generalization

(Supervised Learning)

*Function Approx.*

Exploration
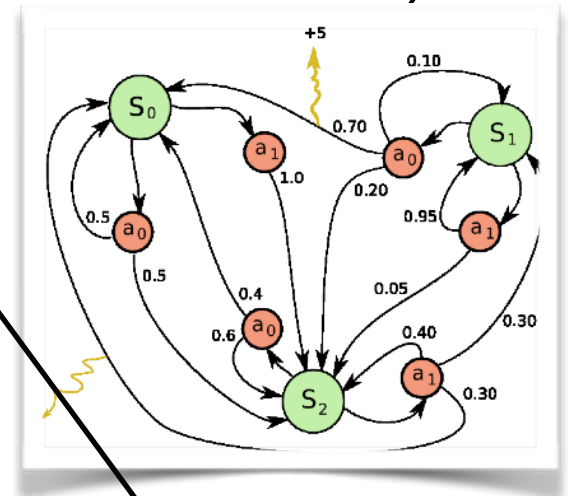
(Multi-Armed Bandit)

*Optimism*

*Bellman Equations*
(Dynamic Programming)
Temporal credit
assignment

(Approx. DP)

(PAC-MDP)

**?**

Generalization

(Contextual Bandit)
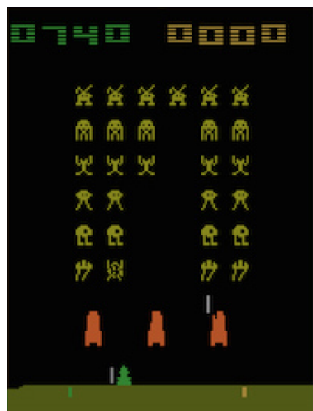
Exploration

(Supervised Learning)

(Multi-Armed Bandit)
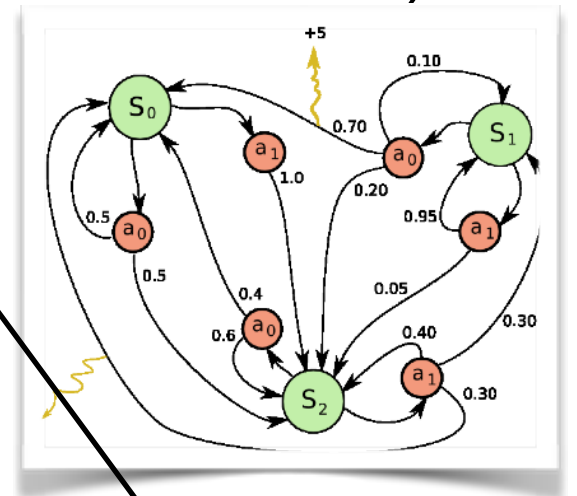
*Function Approx.*

*Optimism*

4

*Bellman Equations*
(Dynamic Programming)
Temporal credit
assignment

(Approx. DP)

(PAC-MDP)

Contextual Decision Processes
*Bellman Rank*

Generalization (Contextual Bandit) Exploration
(Supervised Learning) (Multi-Armed Bandit)
*Function Approx.* *Optimism*

4

# Contextual Decision Processes

Finite action space $\mathcal{A}$, context space $\mathcal{X}$, horizon $H$

For every episode (stochastic and stationary)

- $x_1$ is drawn, and the learner chooses $a_1$.

- $r_2$, $x_2$ are drawn, and the learner chooses $a_2$.

- $r_3$, $x_3$ are drawn, and the learner chooses $a_3$.

- …

- $r_H$ is drawn, and episode ends. (Next episode starts)

Policy $\pi : \mathcal{X} \to \mathcal{A}$

Goal: maximize $\quad V^\pi = \mathbb{E}\left[\sum_{h=1}^{H} r_h \mid a_{1:H} \sim \pi\right]$

# What are contexts?

- Similar to features (a design choice).

- The most detailed choice of context: full interaction history with the environment in this episode.

- Often can be simpler, e.g., when the problem is (short-order) Markov.

- For tabular MDPs: context = (state, time-step).

# Policy vs Value function

- A policy $\pi : \mathcal{X} \to \mathcal{A}$ tells you what to do

- <span style="color:red">Good policy achieves high value</span>

- A value function $f : \mathcal{X} \times \mathcal{A} \to \mathbb{R}$ tells you how much value you would get in the long term

  - induces a greedy policy $\pi_f = (x \mapsto \arg\max_{a \in \mathcal{A}} f(x, a))$

- <span style="color:red">Good value function…</span>

  - <span style="color:red">induces a good policy</span>

  - <span style="color:red">predicts its long-term value accurately</span>

# RL with value-function approximation

Given $\mathcal{F} \subset (\mathcal{X} \times \mathcal{A} \to \mathbb{R})$, learner identifies $f \in \mathcal{F}$

- $\log|\mathcal{F}|$ is small

- $\exists f$ that satisfies *Bellman Equations* ("valid")

$$\forall f' \in \mathcal{F}, h \in [H]$$

$$\mathbb{E}_{\substack{a_{1:h-1} \sim \pi_{f'} \\ a_{h:h+1} \sim \pi_f}} [f(x_h, a_h) - r_h - f(x_{h+1}, a_{h+1})] = 0$$

the optimal value we aim at is $V_{\mathcal{F}}^{\star} = \sup_{\text{valid } f} V^{\pi_f}$

# PAC Learning

- Ideally, we want to identify a near-optimal policy after acquiring

$$poly(|\mathcal{A}|, H, \log|\mathcal{F}|, 1/\delta, 1/\epsilon)$$

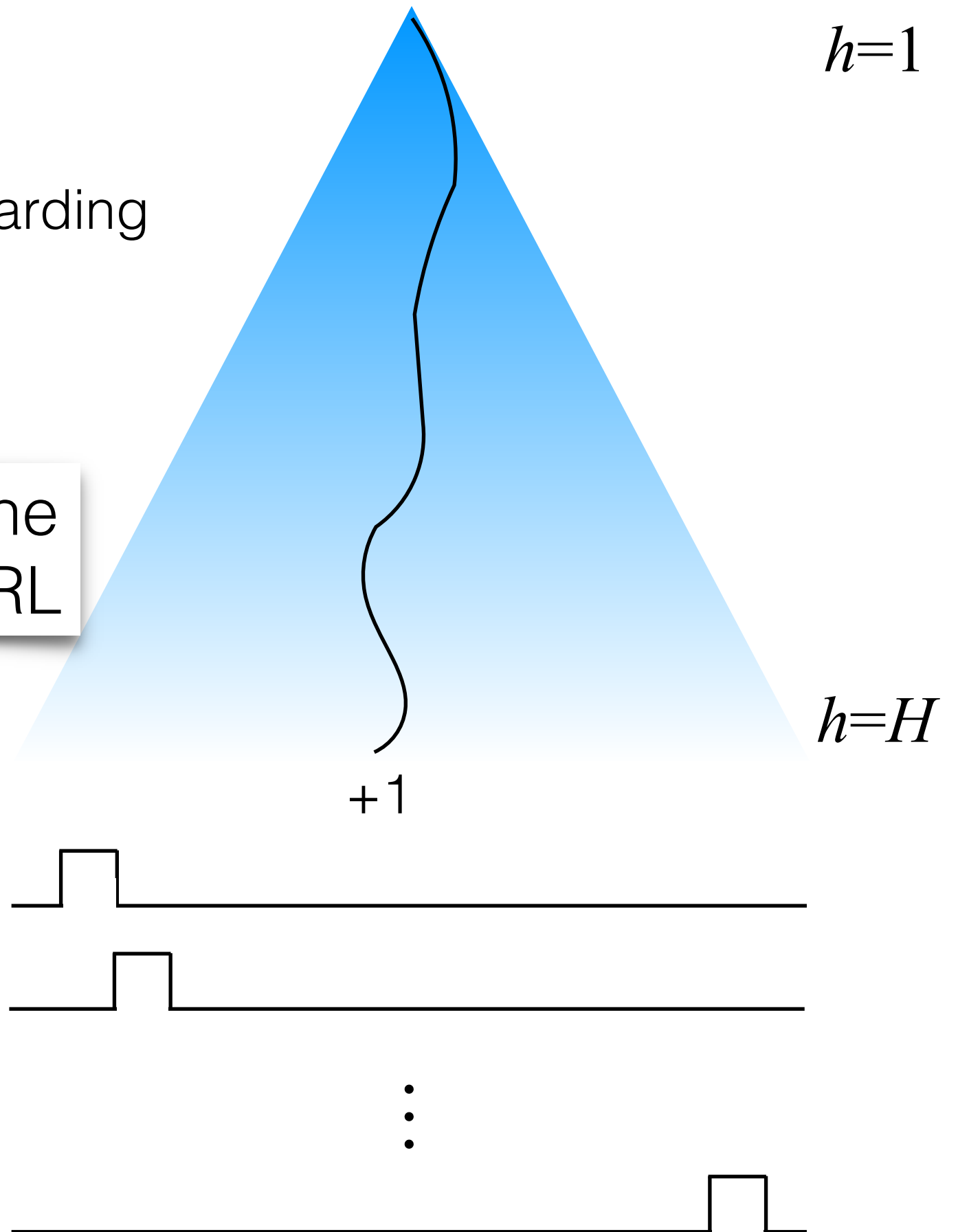episodes of data. (note: no $|\mathcal{X}|$)

- But, there is a lower bound exponential in $H...$

Deterministic, complete tree

One of the $|\mathcal{A}|^H$ leaves is rewarding

need a new measure for the
difficulty of exploration in RL

Size of function space
$\log|\mathcal{F}| = H\log|\mathcal{A}|$

$h$=1

$h$=H

+1

# *Bellman Rank* = rank of
# Bellman error matrix <span style="font-size:smaller">(largest over all level $h$)</span>
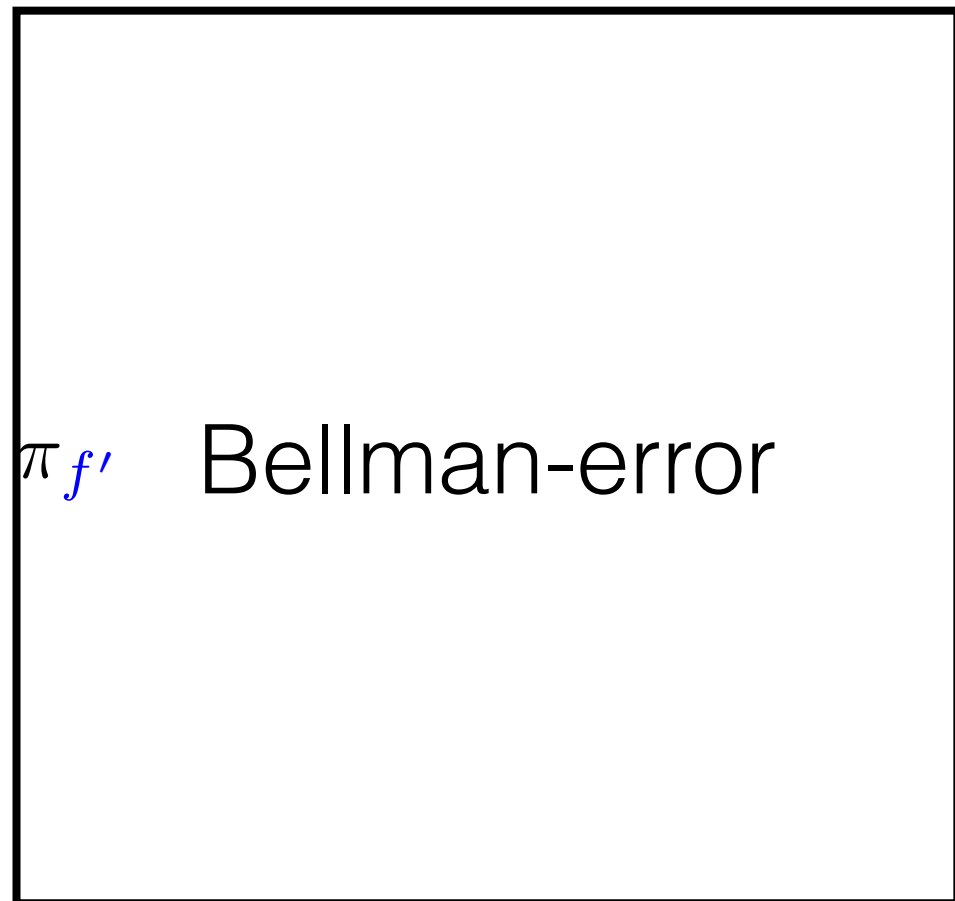
candidate value function
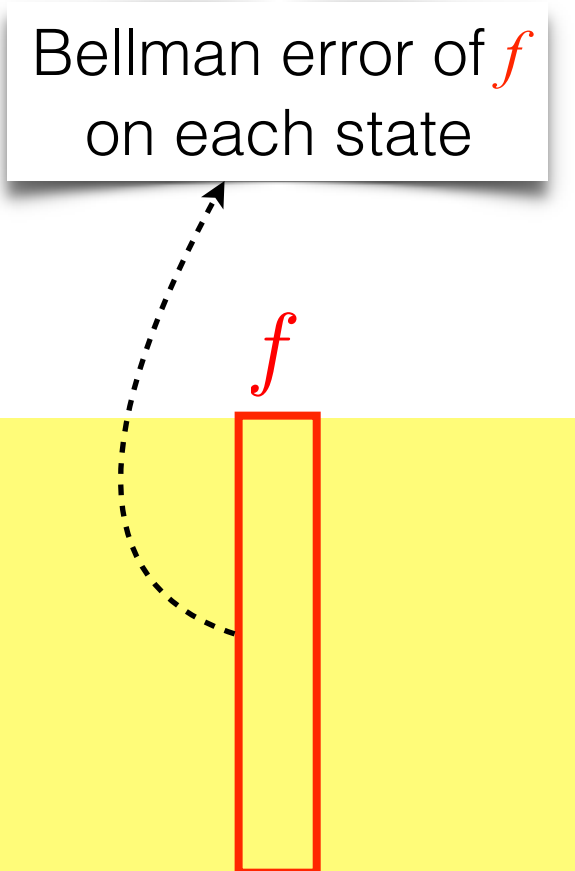
$$f \in \mathcal{F}$$

"roll-in" policy

$$\pi_{f'} : f' \in \mathcal{F} \quad \dots \mathbb{E}_{\substack{a_{1:h-1} \sim \pi_{f'} \\ a_{h:h+1} \sim \pi_f}} [f(x_h, a_h) - r_h - f(x_{h+1}, a_{h+1})]$$

11

# Tabular MDP has low Bellman Rank

$f$

distribution over states
induced by the roll-in policy

Bellman error of $f$
on each state

$f$

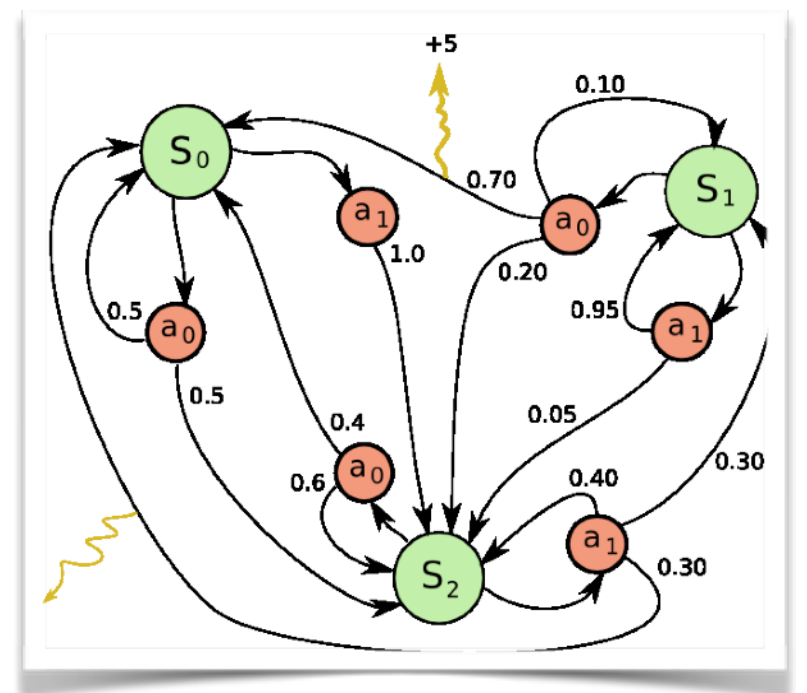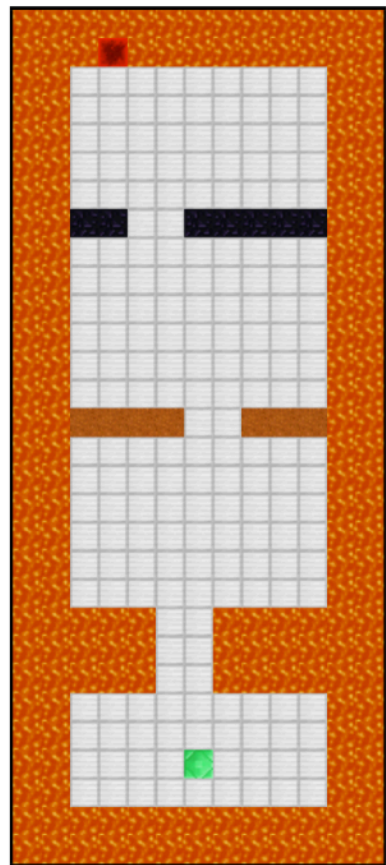$$\pi_{f'} \quad \text{Bellman-error} \quad = \pi_{f'} \quad \times$$

$$\mathbb{E}_{\substack{a_{1:h-1} \sim \pi_{f'} \\ a_{h:h+1} \sim \pi_f}} [f(x_h, a_h) - r_h - f(x_{h+1}, a_{h+1})]$$
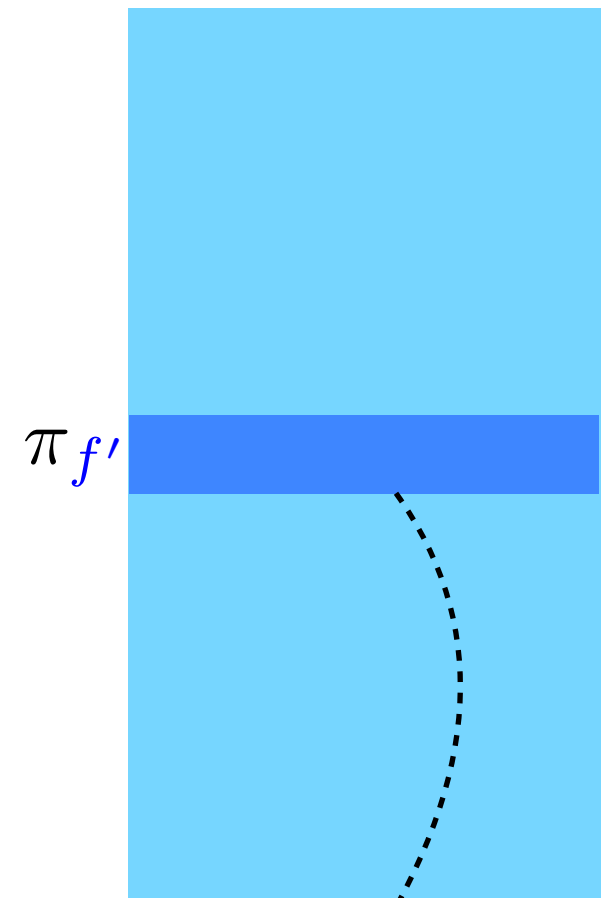
# "Visual grid-world" has low Bellman Rank



small hidden
state space

rich & complex observations
(and near-Markovian)

$\pi_{f'}$

distribution over hidden-states
induced by the roll-in policy

Figure from [Johnson et al. IJCAI'16]

# RL settings that yield low Bellman Rank ($M$)
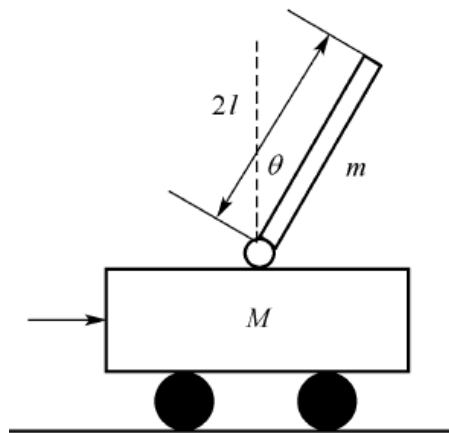


rich-obs POMDPs w/
reactive value function:
$M$ = #hidden states

$$P_{\mathcal{T}|h}$$

PSRs w/ similar set-up:
$M$ = poly(linear dim.)



large MDPs w/
low-rank dynamics:
$M$ = rank of transition matrix



LQR control*:
$M$ = poly(#state variables)

$$Q^{\star}$$

state abstraction
that preserves Q*:
$M$ = poly(#abstract states)

.. can measure any
(process, function space)

*our algorithm does not directly
apply to continuous actions

# Overview of the algorithm



$\pi_{f'}$

Bellman error

$$\mathbb{E}_{\substack{a_{1:h-1}\sim\pi_{f'} \\ a_{h:h+1}\sim\pi_f}}\left[f(x_h,a_h)-r_h-f(x_{h+1},a_{h+1})\right]$$

$>\phi \qquad\qquad >\phi \quad >\phi$

1. Pick an exploration policy

2. Estimate all entries in the corresponding row

3. Eliminate columns whose error > 0 (w/ statistical significance)

# Overview of the algorithm



$\pi_{f'}$

Bellman error

$$\mathbb{E}_{\substack{a_{1:h-1}\sim\pi_{f'} \\ a_{h:h+1}\sim\pi_f}} [f(x_h, a_h) - r_h - f(x_{h+1}, a_{h+1})]$$
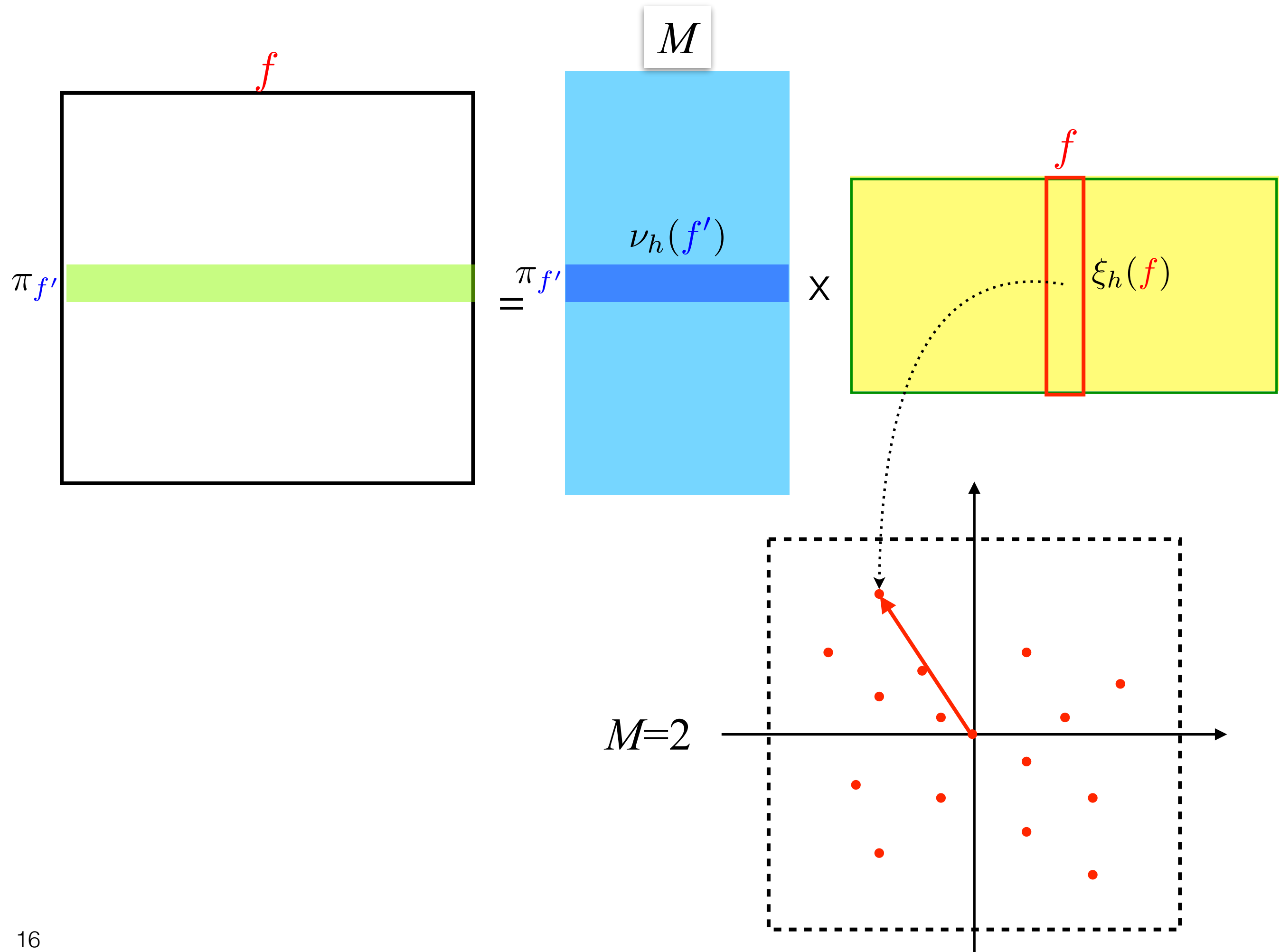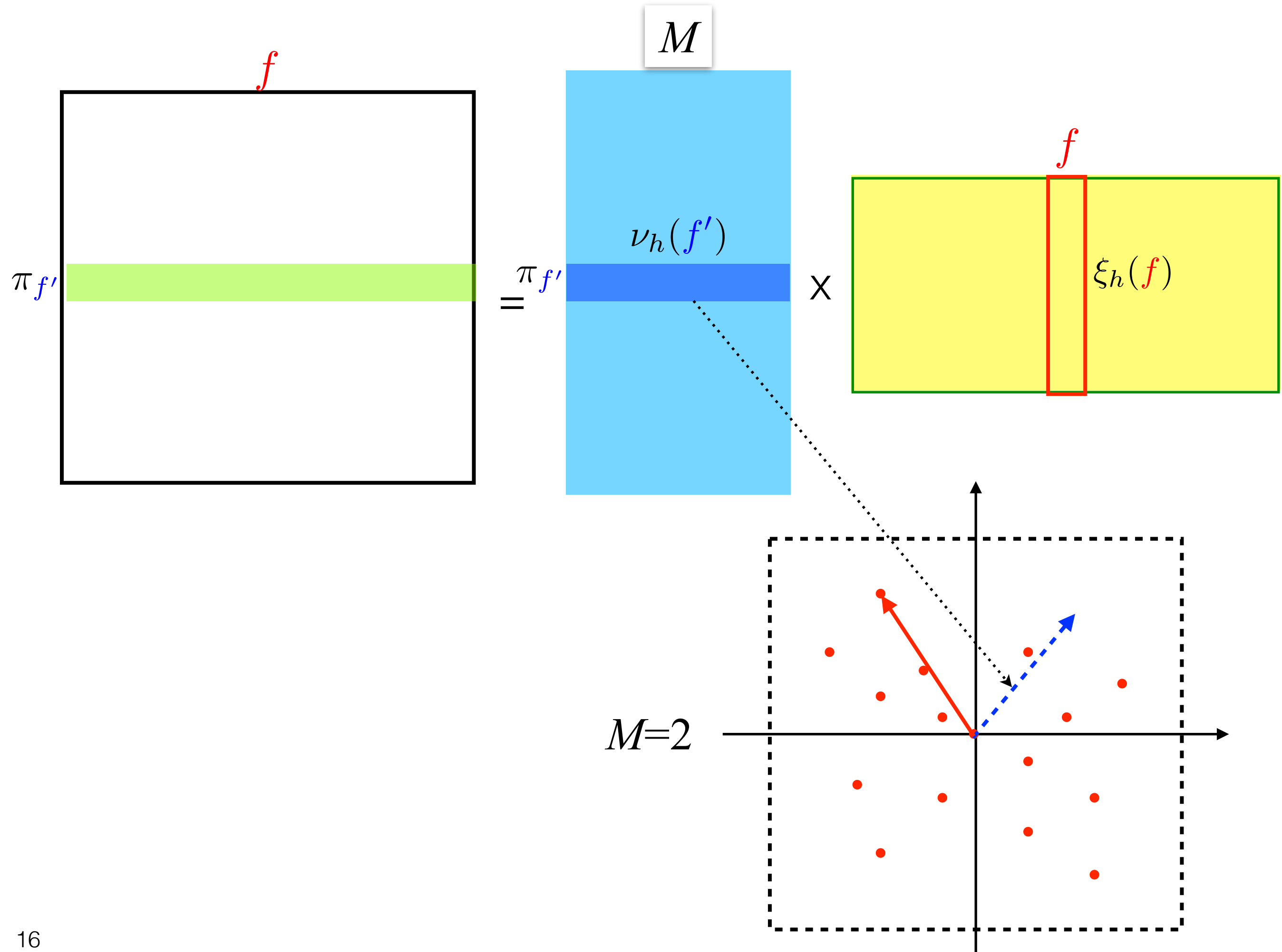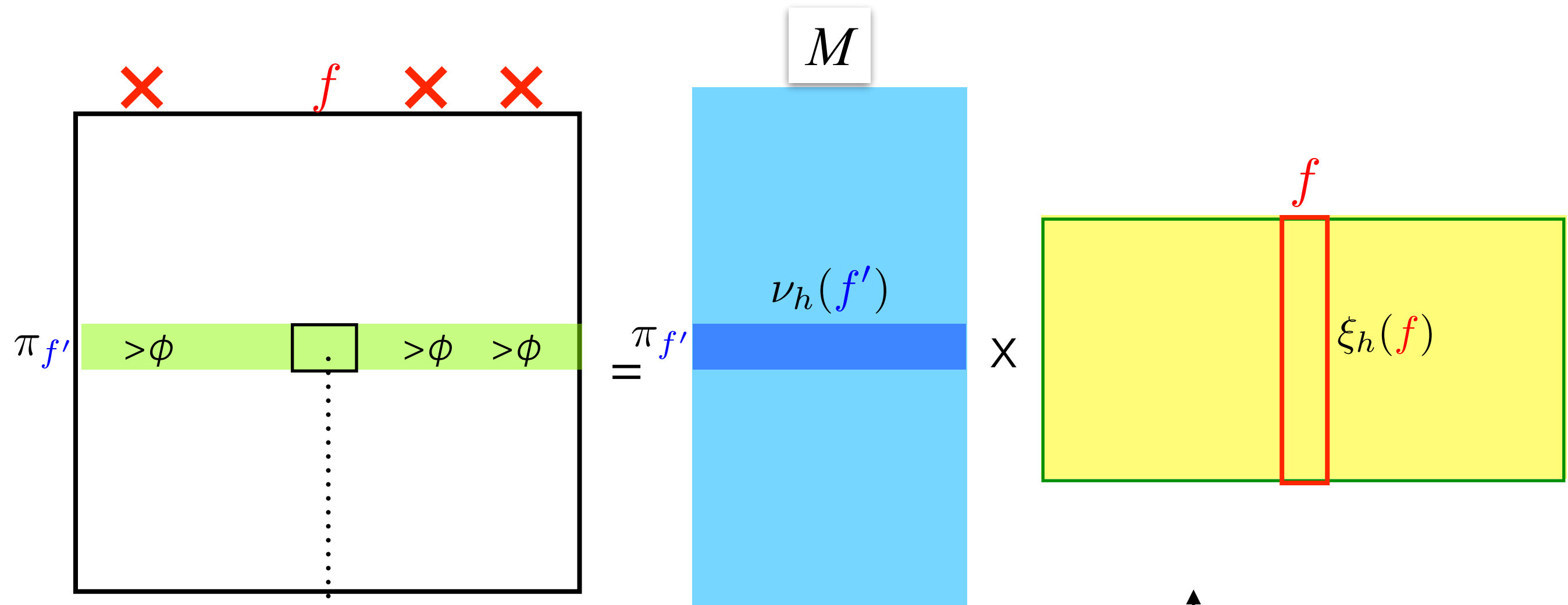
$> \phi$ $\qquad$ $> \phi$ $\quad$ $> \phi$

1. Pick an exploration policy

2. Estimate all entries in the corresponding row

3. Eliminate columns whose error > 0 (w/ statistical significance)
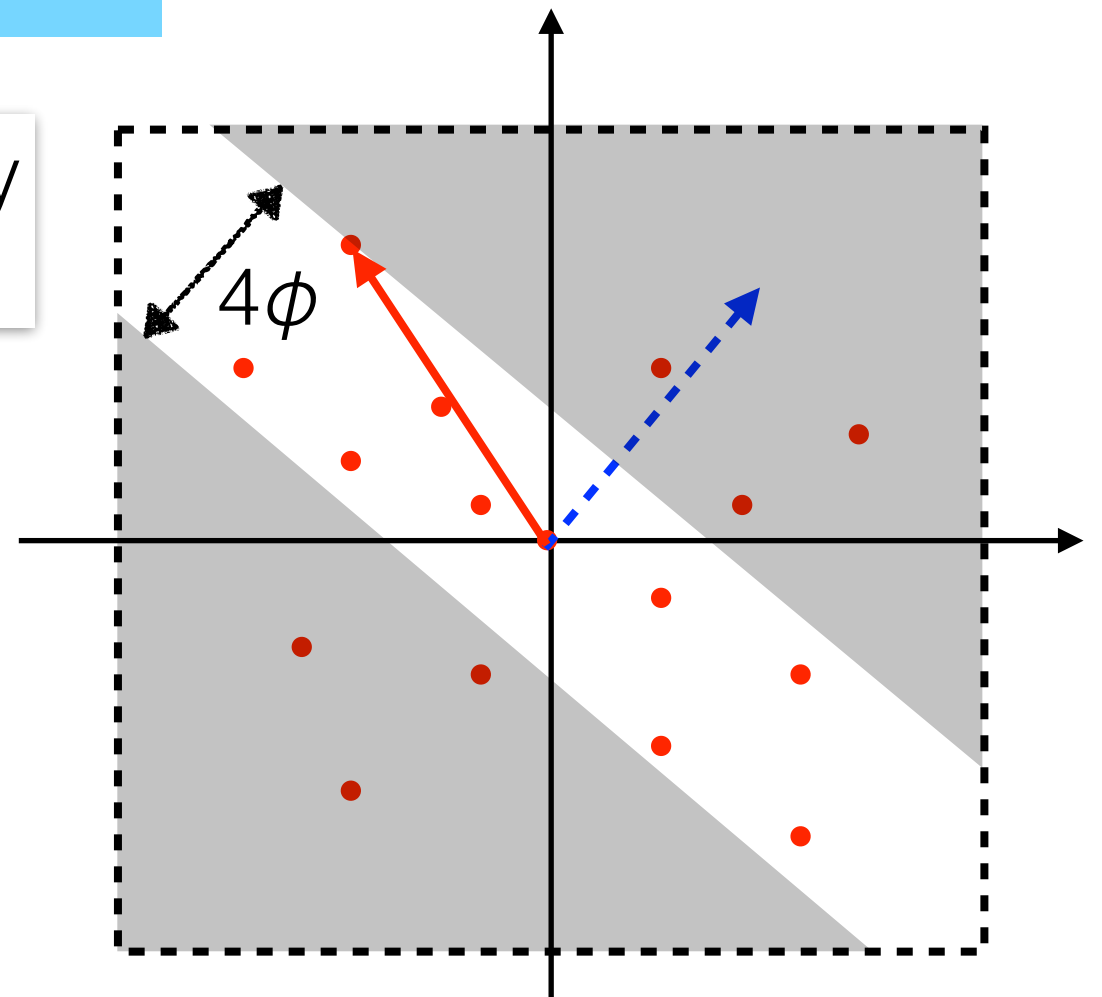
4. Repeat

*There are $H$ such matrices

15

$$\pi_{f'} = {}^{\pi_{f'}}\nu_h(f') \times \xi_h(f)$$

$M$

$f$

$M=2$

16

$\pi_f$ $f$

$= {}^\pi f'$ $M$ $\nu_h(f')$ $\times$ $\xi_h(f)$ $f$

$M=2$

16

$\times$    $f$    $\times$    $\times$

$\pi_{f'}$    $>\phi$    $>\phi$    $>\phi$

$\langle \rightarrow , \dashrightarrow \rangle$

$M$

$= {}^{\pi_{f'}} \quad \nu_h(f') \quad \times \quad \xi_h(f)$

$f$

$\phi$ controlled by sample size

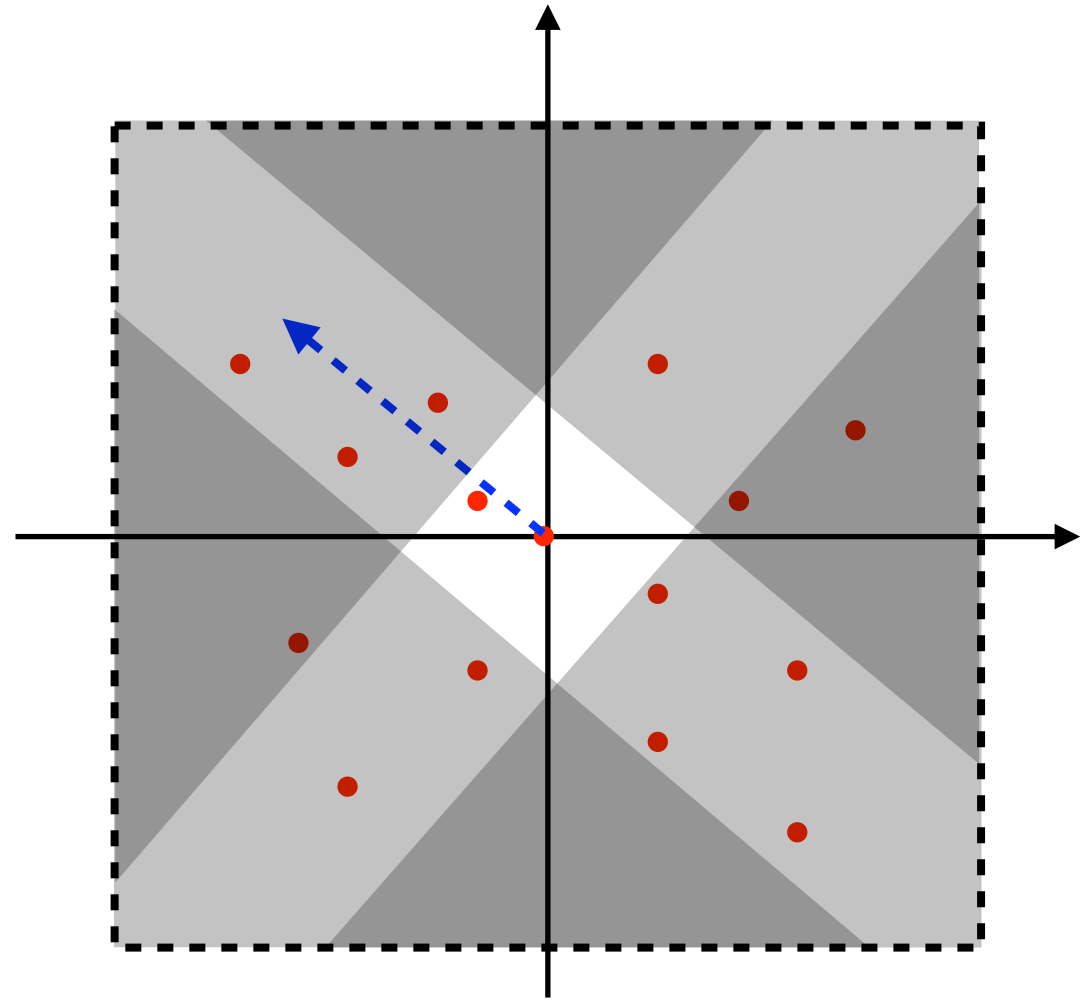$4\phi$

$M{=}2$

key observation:

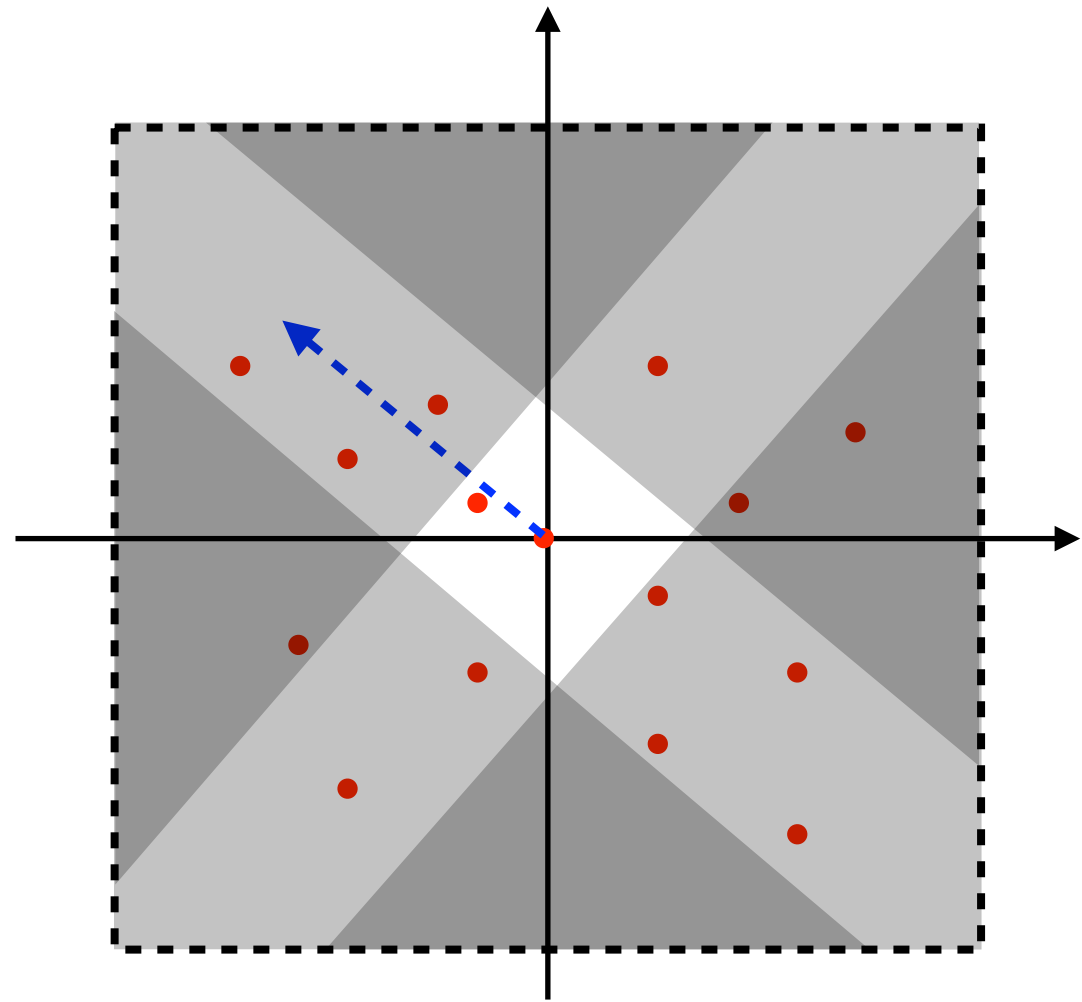$\rightarrow$ and $\dashrightarrow$ are roughly orthogonal

16

inefficient exploration

- new distribution is similar to previous ones
- area of while space shrinks slowly

efficient exploration

- new distribution is different from previous ones
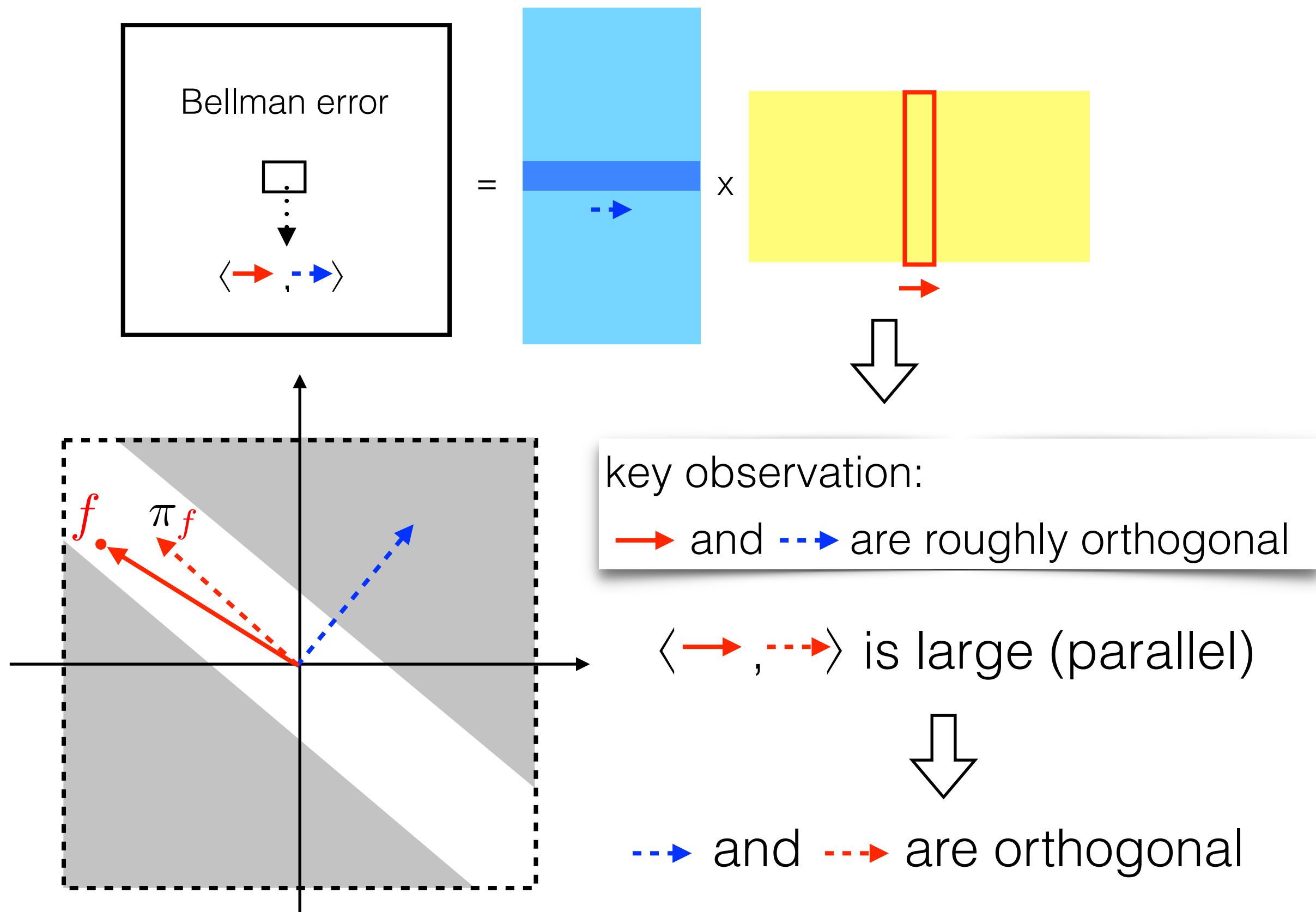- area of while space shrinks quickly

efficient exploration

algorithm
- new distribution is different from previous ones

analysis
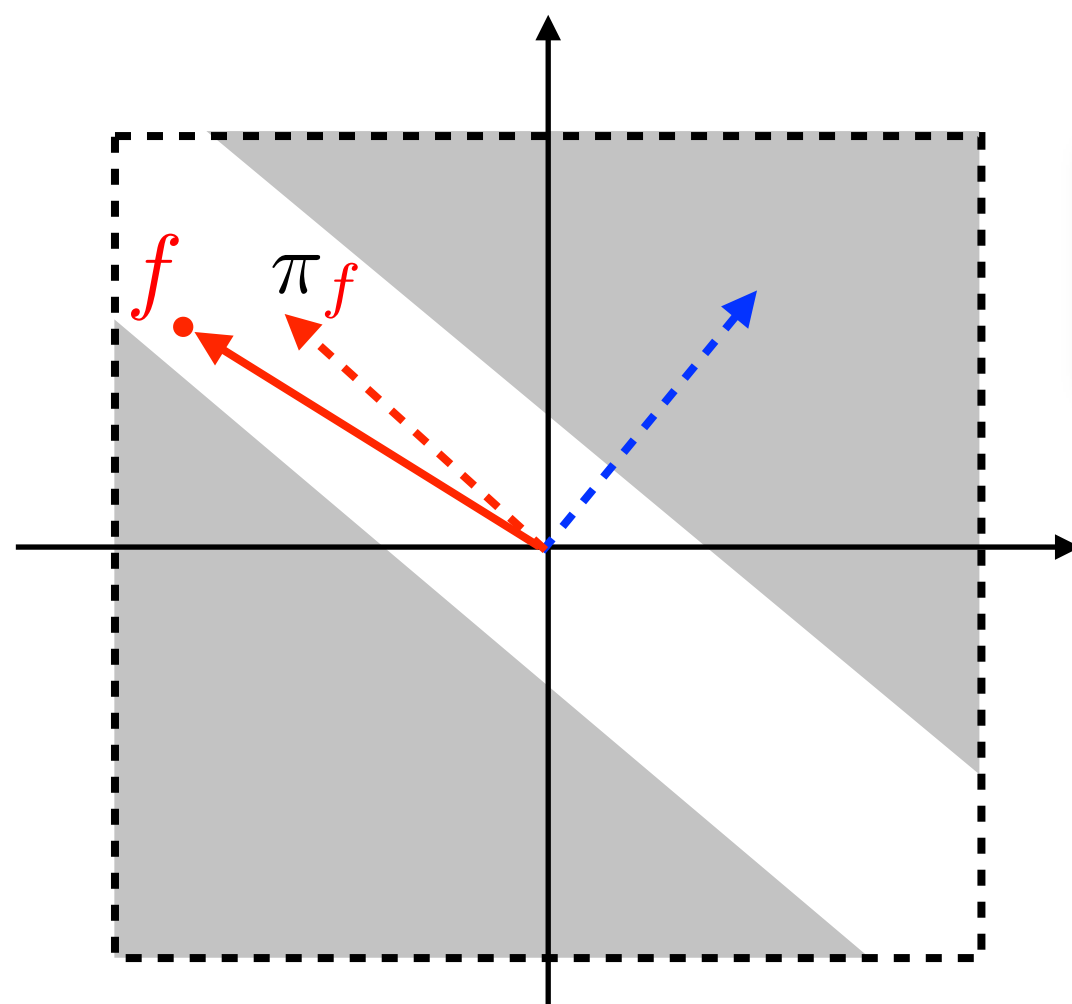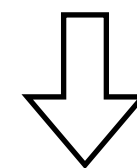- area of while space shrinks quickly

# Exploration strategy: pick $f$ optimistically and explore with $\pi_f$



Bellman error

$$\langle \textcolor{red}{\rightarrow} , \textcolor{blue}{\dashrightarrow} \rangle$$

$=$ $\times$

key observation:

$\textcolor{red}{\longrightarrow}$ and $\textcolor{blue}{\dashrightarrow}$ are roughly orthogonal

$\langle \textcolor{red}{\rightarrow} , \textcolor{red}{\dashrightarrow} \rangle$ is large (parallel)

$\textcolor{blue}{\dashrightarrow}$ and $\textcolor{red}{\dashrightarrow}$ are orthogonal

$f$ $\pi_f$

# Exploration strategy: pick $f$ optimistically and explore with $\pi_f$
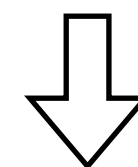
Exploration Lemma: for any $f$,

$$\mathbb{E}[\max_{a \in \mathcal{A}} f(x_1, a)] - V^{\pi_f} = \sum_{h=1}^{H} \mathbb{E}_{\substack{a_{1:h-1} \sim \pi_f \\ a_{h:h+1} \sim \pi_f}} \left[ f(x_h, a_h) - r_h - f(x_{h+1}, a_{h+1}) \right]$$
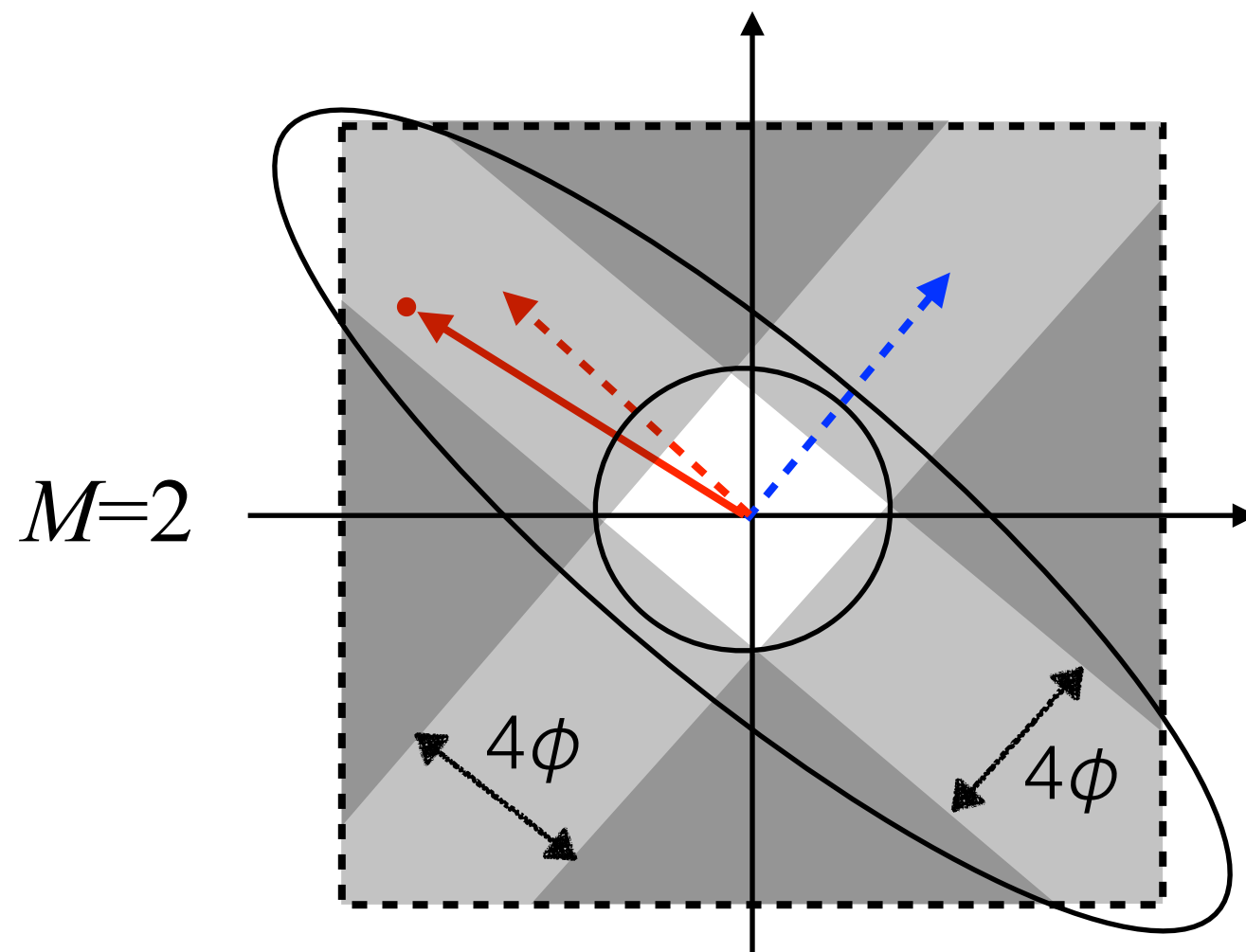
$\langle \longrightarrow , \dashrightarrow \rangle$

key observation:

$\longrightarrow$ and $\dashrightarrow$ are roughly orthogonal

$\langle \longrightarrow , \dashrightarrow \rangle$ is large (parallel)

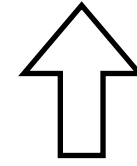$\dashrightarrow$ and $\dashrightarrow$ are orthogonal

$M=2$

$4\phi$

$4\phi$

Adaptation of [Todd,1982]:
Ellipsoid volume shrinks exponentially if

$$|\langle \textcolor{red}{\longrightarrow} , \textcolor{red}{\dashrightarrow} \rangle| \ \geq \ 3\sqrt{M} \ \times \ 2\phi$$

controlled by sub-optimality

controlled by sample size

# Sample complexity

We can identify a policy $\varepsilon$-suboptimal compared to $V_{\mathcal{F}}^{\star}$ with probability at least 1-$\delta$, after acquiring this many episodes of data:

$$\tilde{O}\left(\frac{M^2 H^3 |\mathcal{A}|}{\epsilon^2} \log(|\mathcal{F}|/\delta)\right)$$