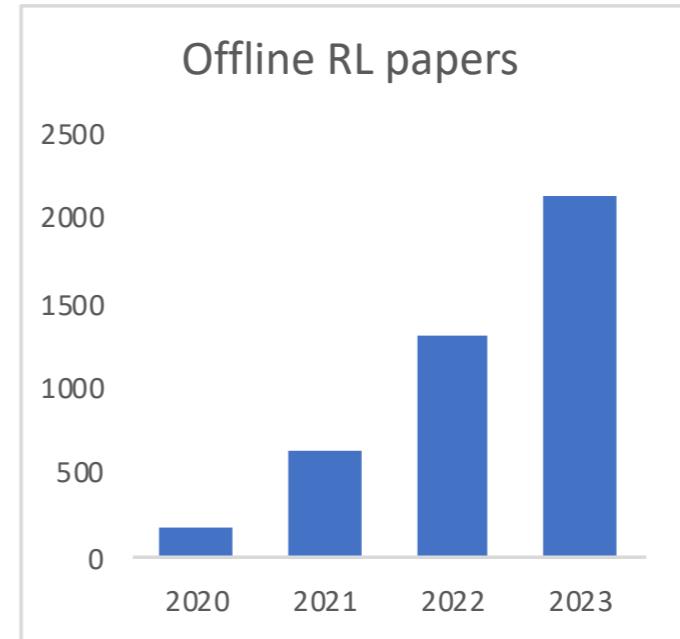
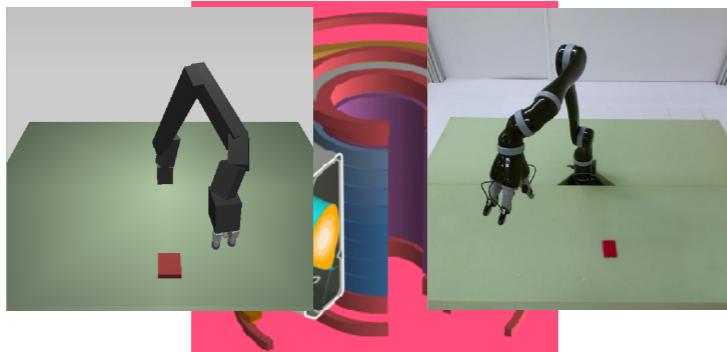


Rethinking the theoretical foundation of reinforcement learning

Nan Jiang
University of Illinois at Urbana-Champaign
April 24, 2025
@Shariff University of Technology

- (Offline) RL in **real life**

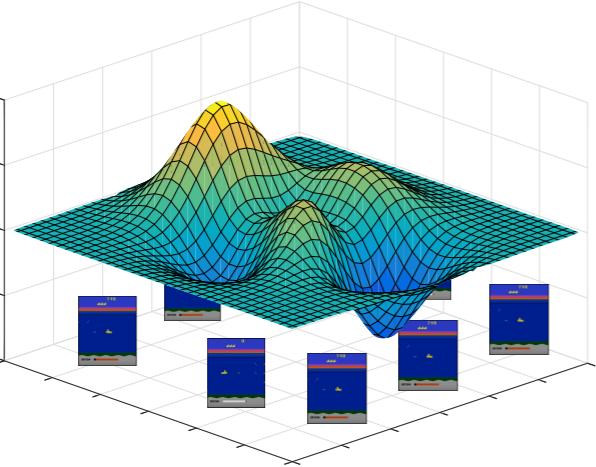
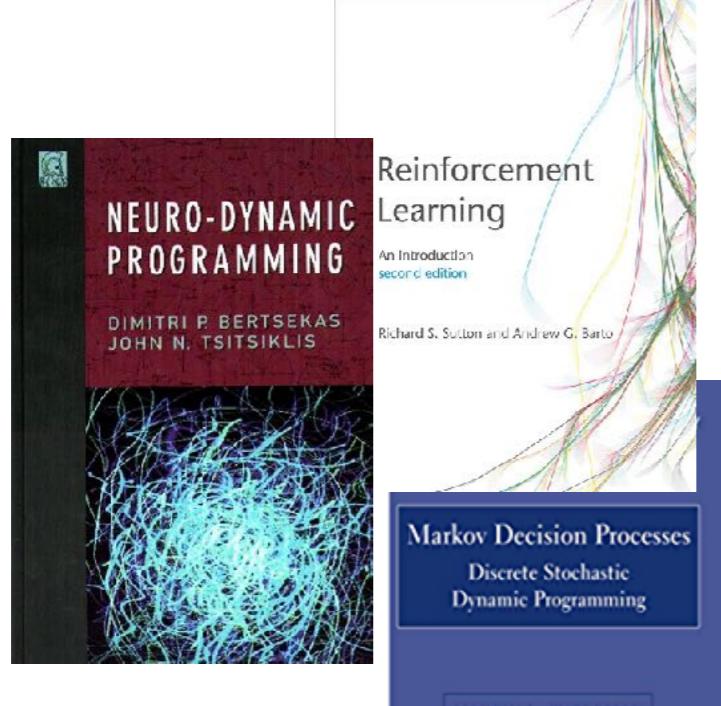


Key ingredient: simulator

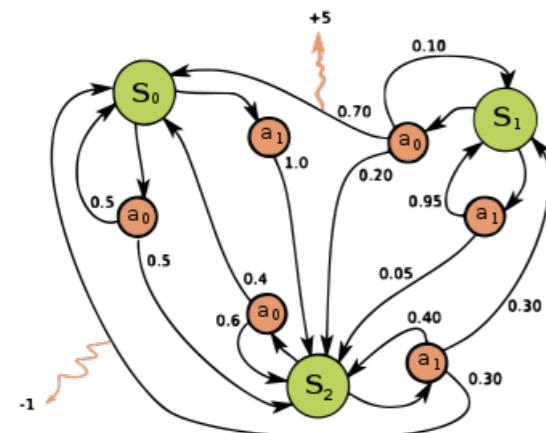
- Unlimited data **X**
- Decision w/o real consequences **X**
- Can easily evaluate new strategy **X**

Why are we **not** seeing (offline) RL deployed everywhere already?

~2000

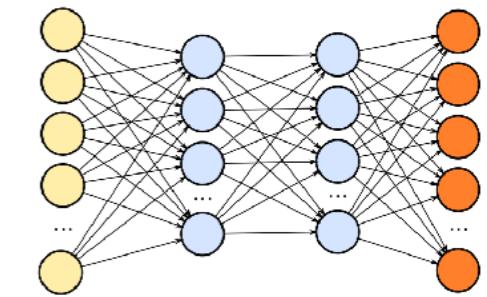


*Bellman rank,
Eluder dimension,
Concentrability, ...*

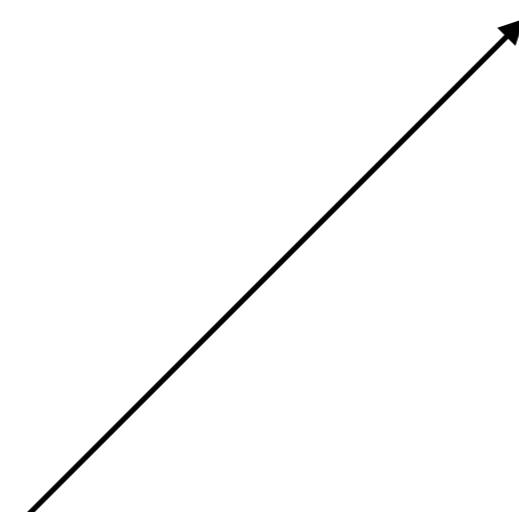


\sqrt{HSAT} regret,
 SAH^2/ϵ^2 sample
complexity, ...

- (Offline) RL in **real life**
- Role of theory in **modern RL**

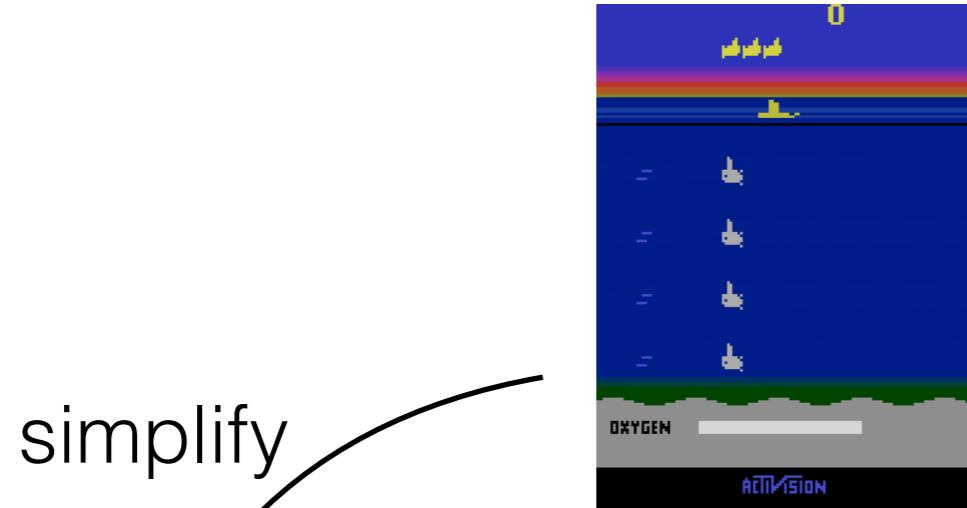


*Empirical: Atari, Mujoco,
OpenAI Gym, target
network, architecture, ...*

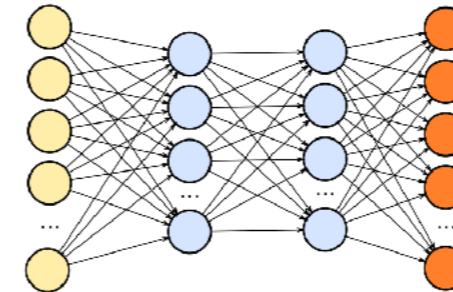


* finite-sample analysis of ADP & MCTS 00~10

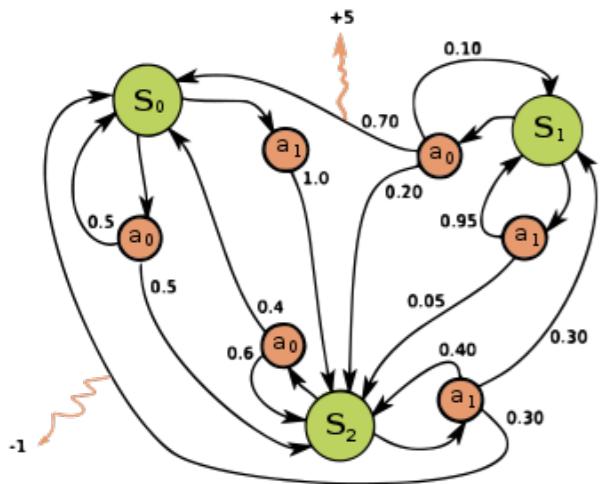
- (Offline) RL in **real life**
- Role of theory in **modern RL**
- Theoretical foundation



simplify

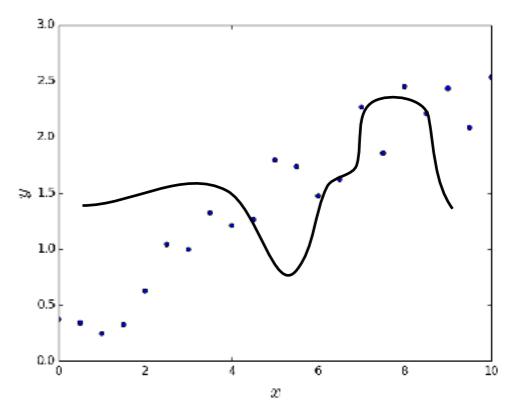
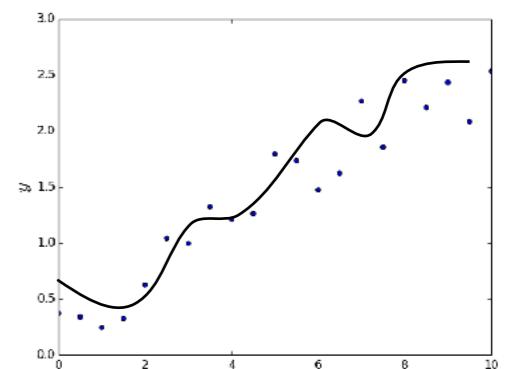


extend

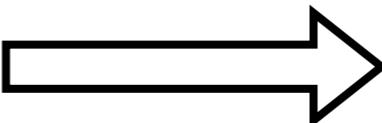


	$Q^*(s, a)$
(S_1, a_1)	...
(S_1, a_2)	...
(S_2, a_1)	...

“tabular” RL

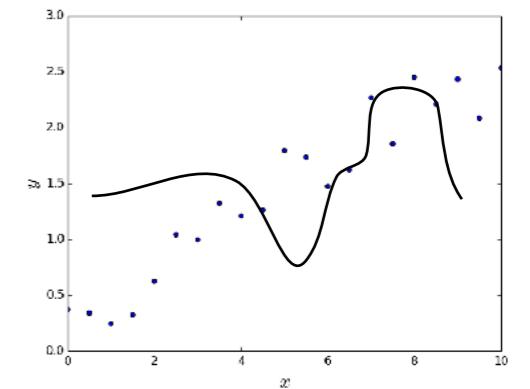
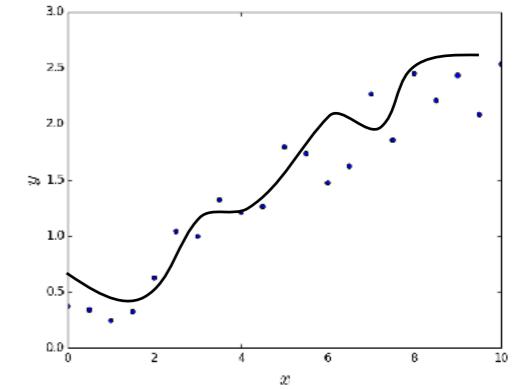
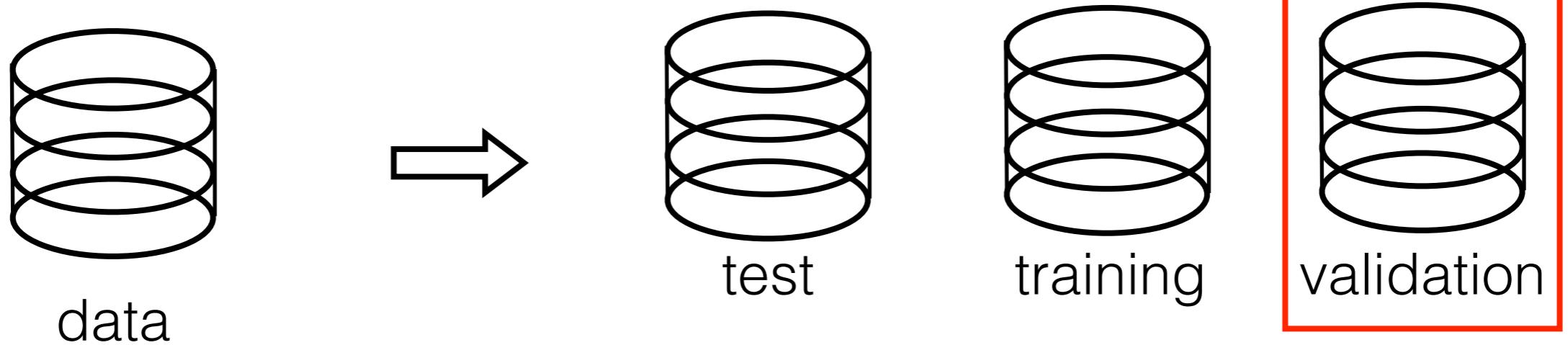


- (Offline) RL in real life
- Role of theory in modern RL
- Theoretical foundation

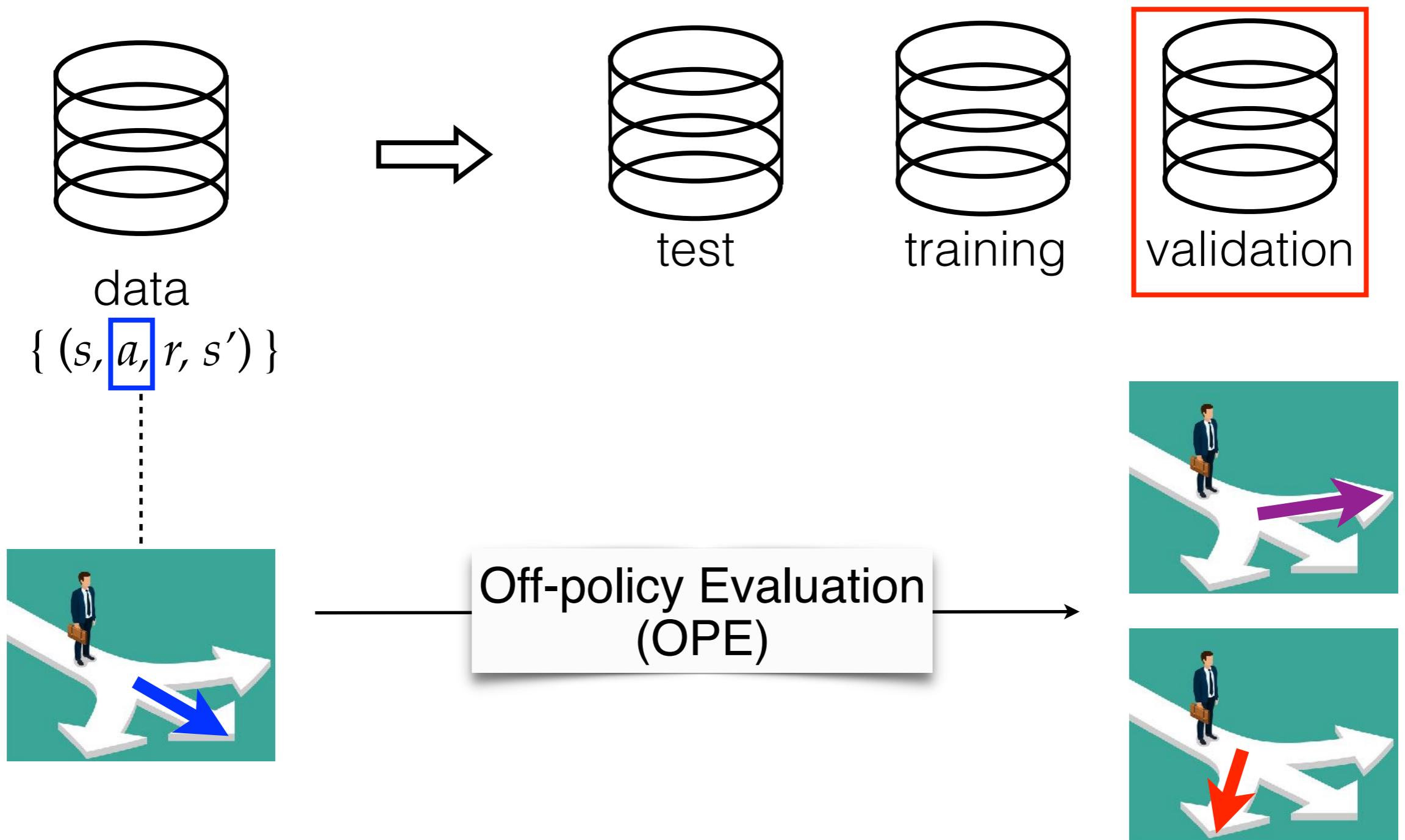
modern RL  real life

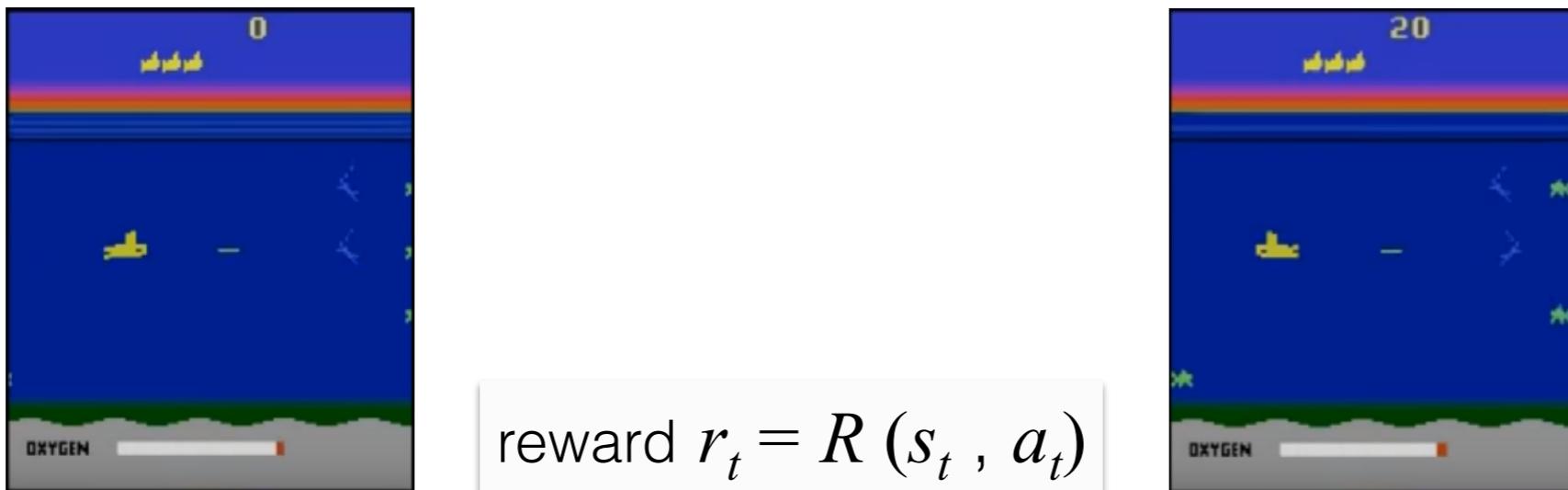
Rethinking
theoretical foundation

Supervised learning pipeline

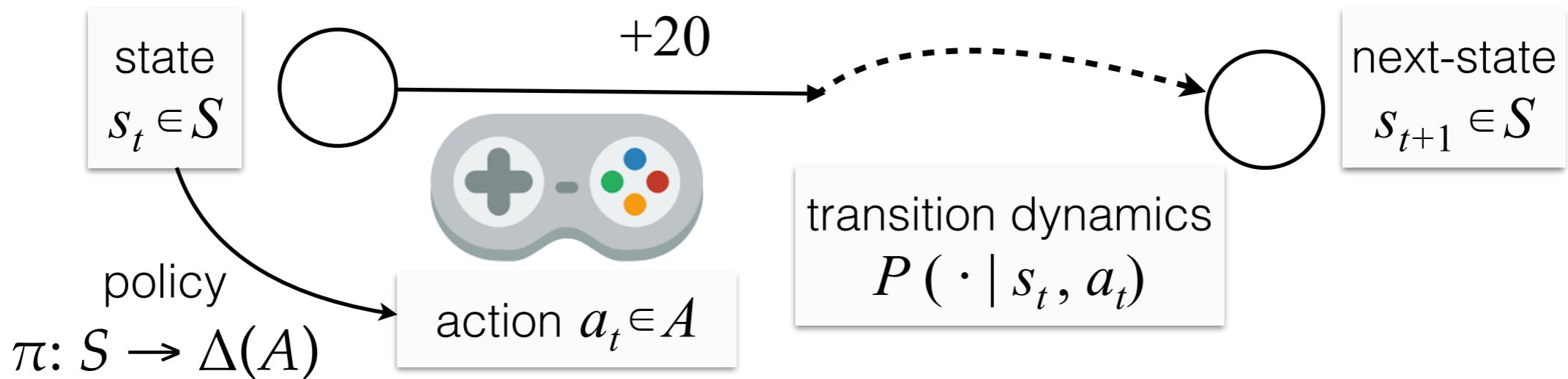


Offline RL pipeline





$$\text{reward } r_t = R(s_t, a_t)$$

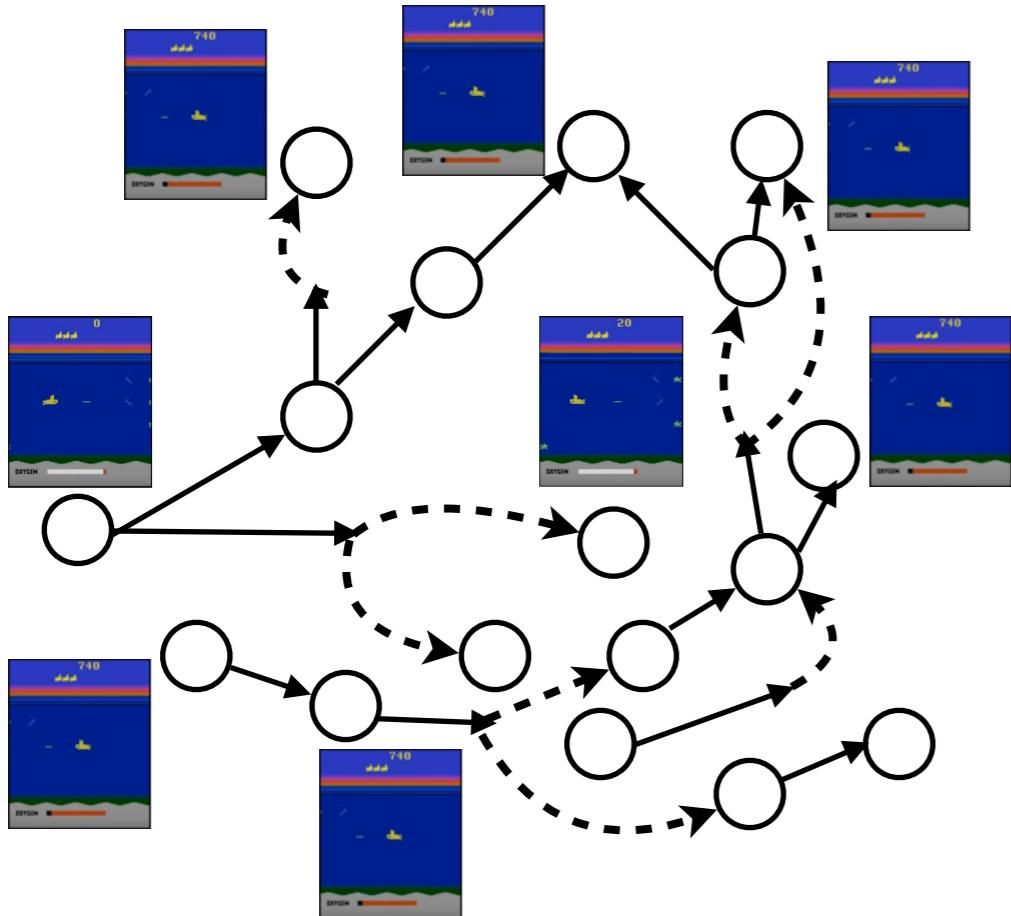


Policy evaluation: estimate $J(\pi) := \mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t r_t | s_0]$ given π

Policy optimization: $\max_\pi J(\pi) = Q^\pi(s_0, \pi)$

How to find Q^π ?

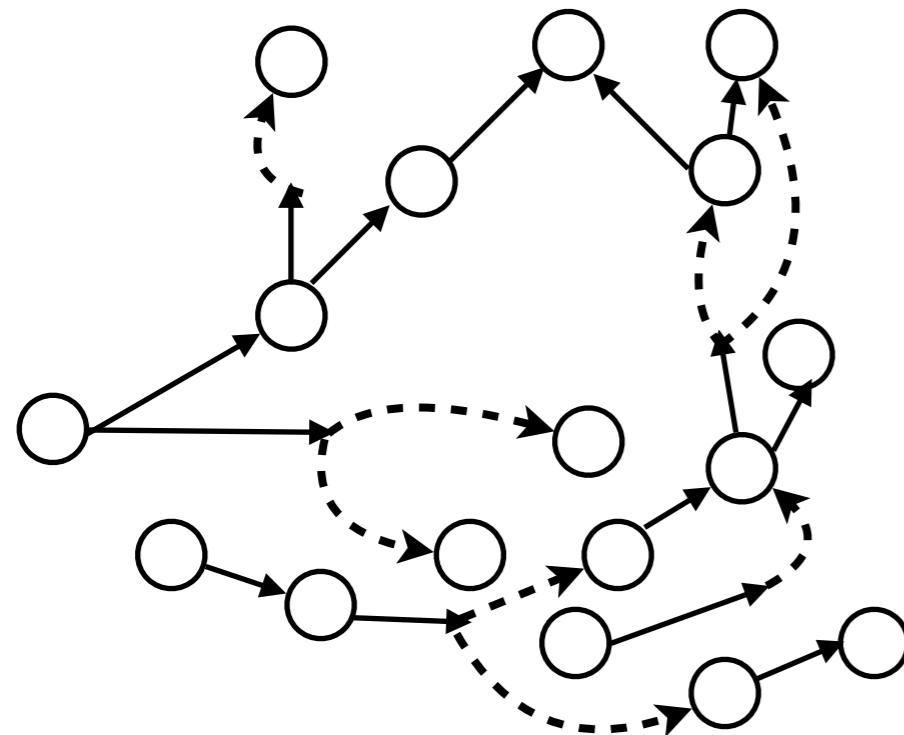
$Q^\pi = \mathcal{T}^\pi Q^\pi \rightarrow |S \times A| \text{ equations}$



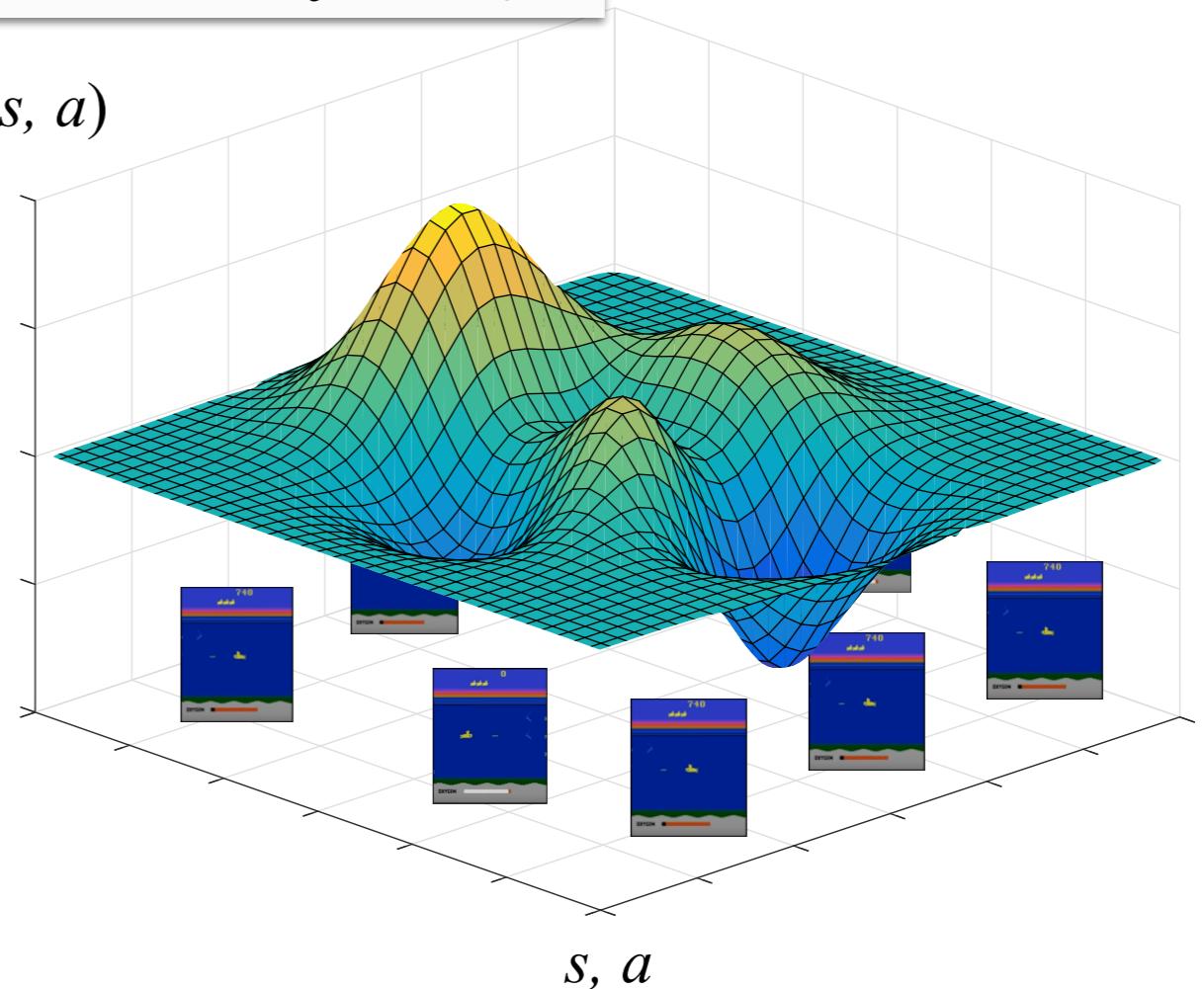
How to find Q^π ?

$$Q^\pi = \mathcal{T}^\pi Q^\pi \rightarrow |SxA| \text{ equations } \mathbf{X}$$

Find θ s.t. $f_\theta \approx Q^\pi$



$f_\theta(s, a)$



$s_0, a_0, r_0, s_1, a_1, r_1, s_2, a_2, r_2, s_3, \dots$



$(s, a, r, s') \sim D$

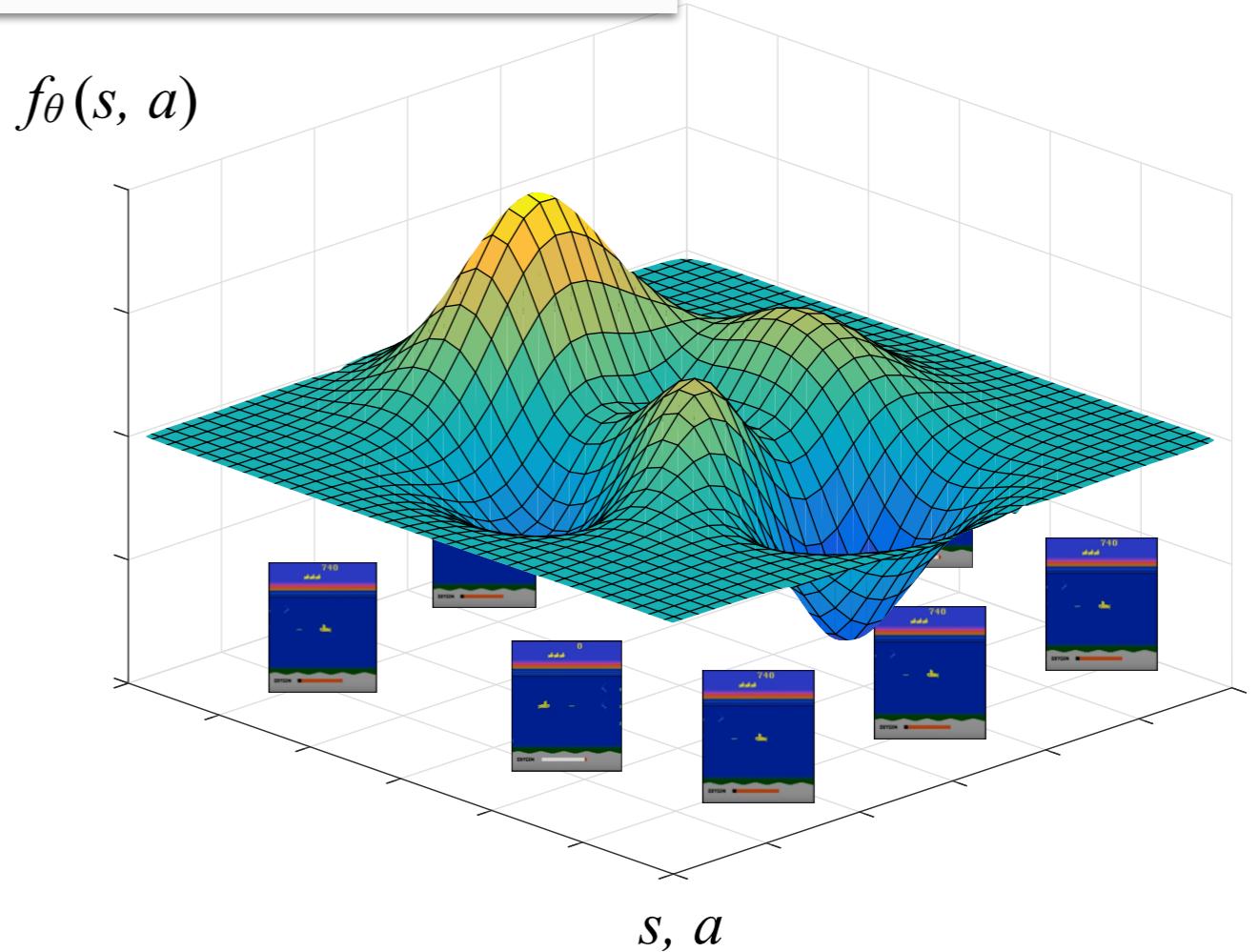
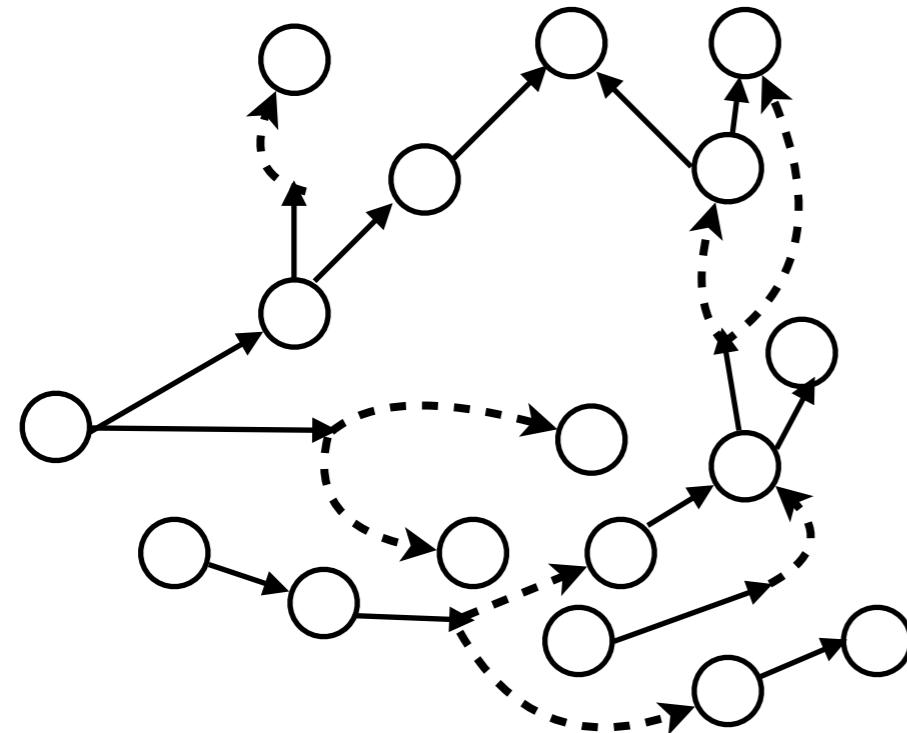
Validation:

(FQE: learn Q^π)

$$f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma f_{k-1}(s', \pi))^2]$$

iterative

Find θ s.t. $f_\theta \approx Q^\pi$



Validation:

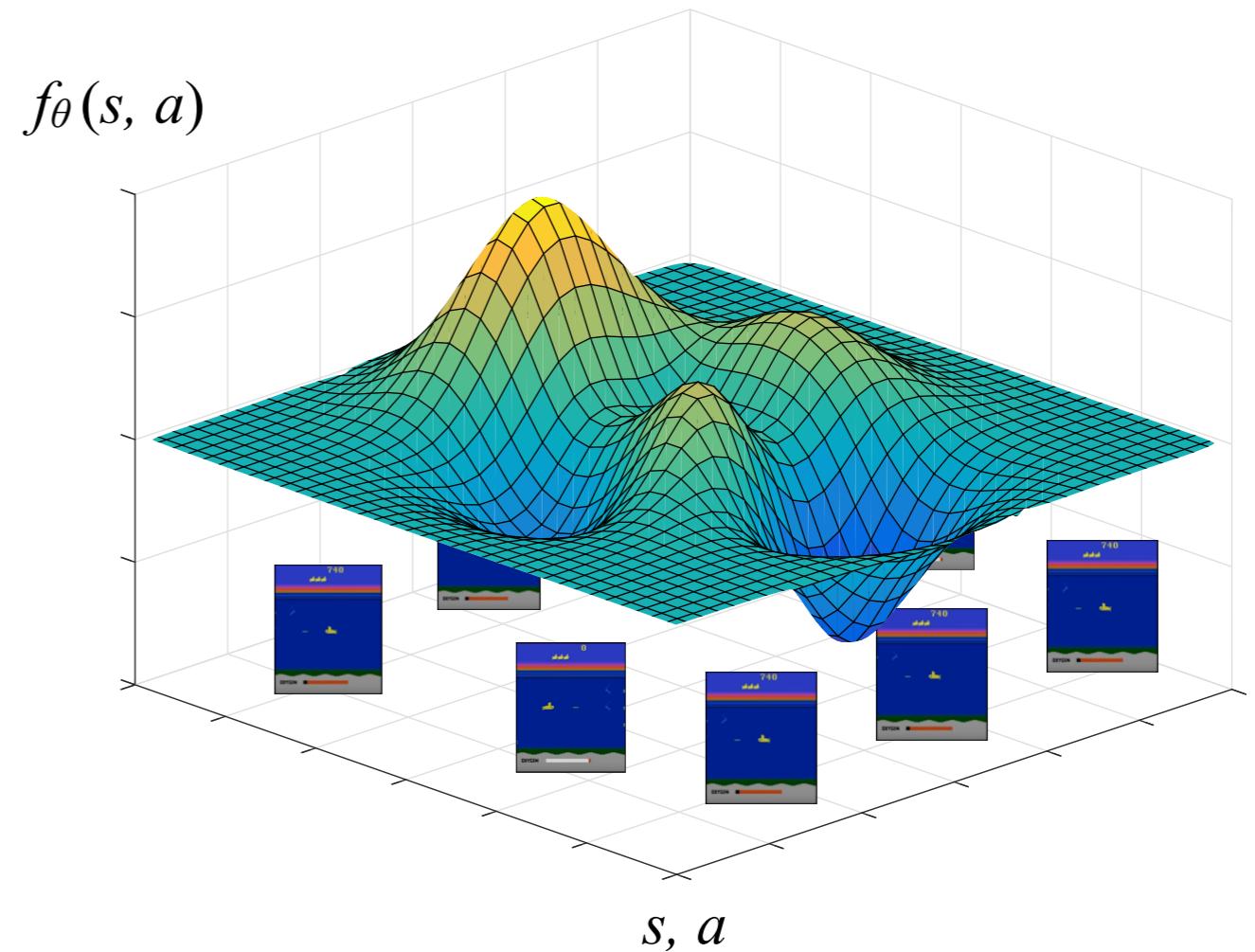
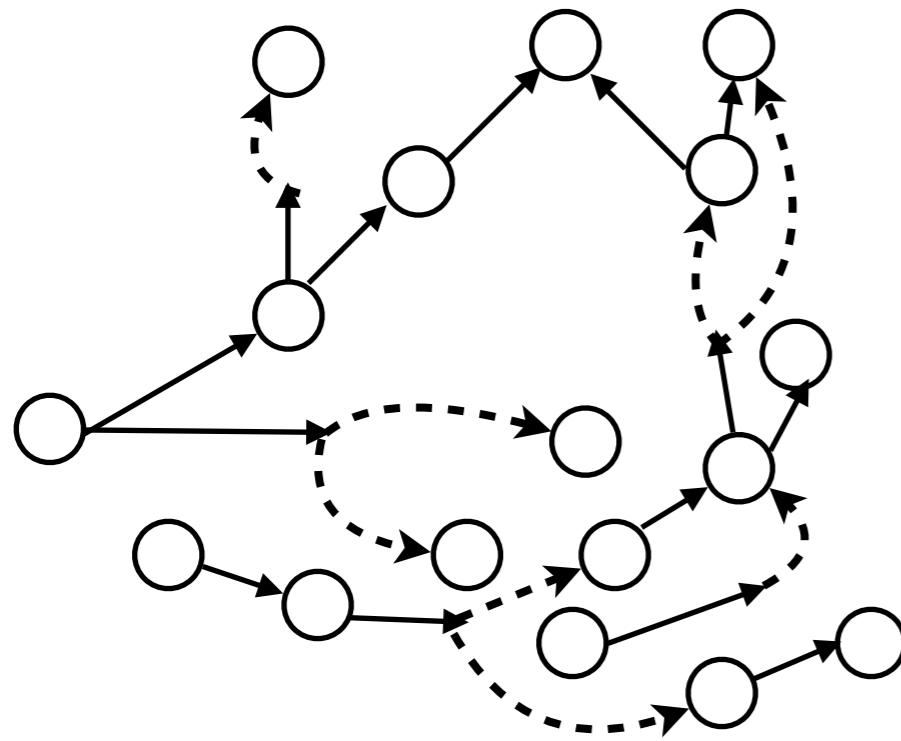
(FQE: learn Q^π)

$$f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma f_{k-1}(s', \pi))^2]$$

iterative

$$\approx \mathcal{T}^\pi f_{k-1}$$

$\mathbb{E}[\cdot | s, a]$



Training: $\hat{f} = f_k$ where

(FQI: learn Q^*)

$$f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma \max_{a'} f_{k-1}(s', a'))^2]$$



Validation:

(FQE: learn Q^π)

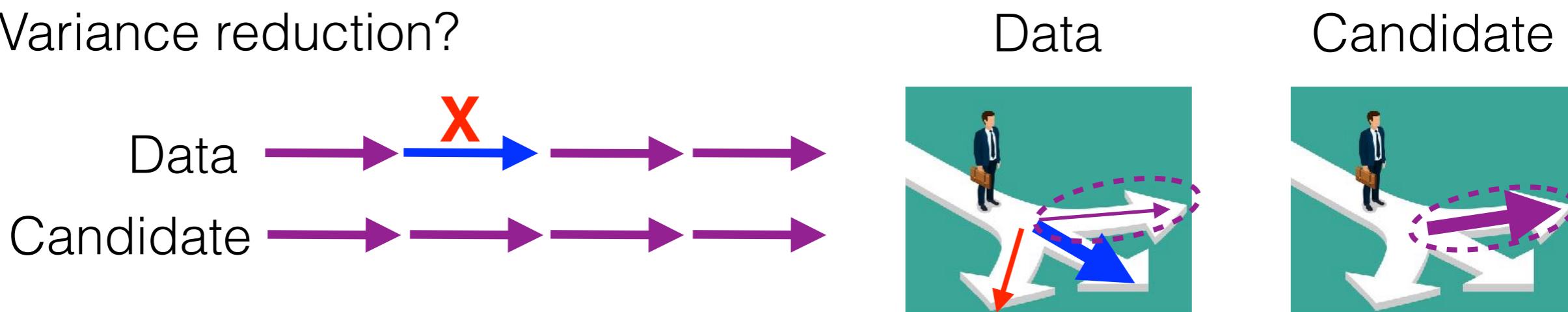
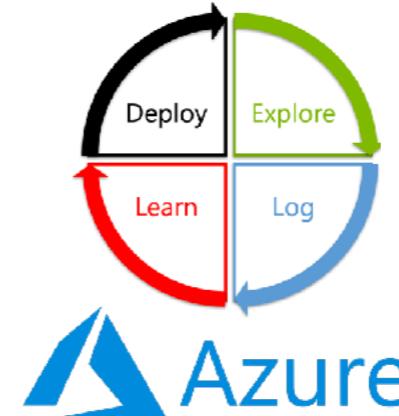
$$f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma f_{k-1}(s', \pi))^2]$$

\downarrow π = greedy w.r.t. \hat{f}

Hyperparameter-free methods?

Importance sampling [Precup'00]

- Hyperparameter-free ✓
- No Markovianity required ✓
- Industry deployment (ctx. bandit, horizon=1)
- **Exponential-in-horizon** variance!
- Variance reduction?

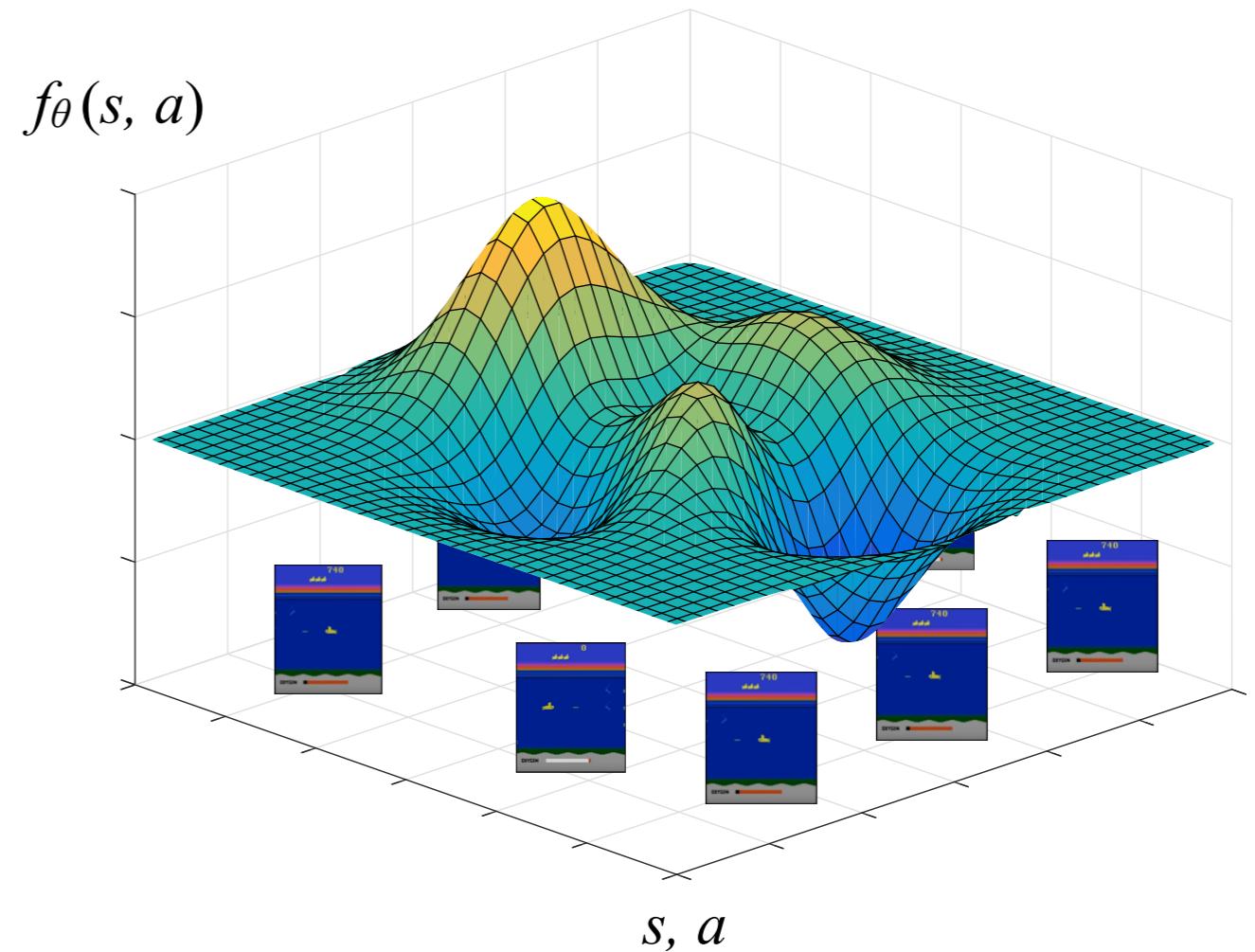
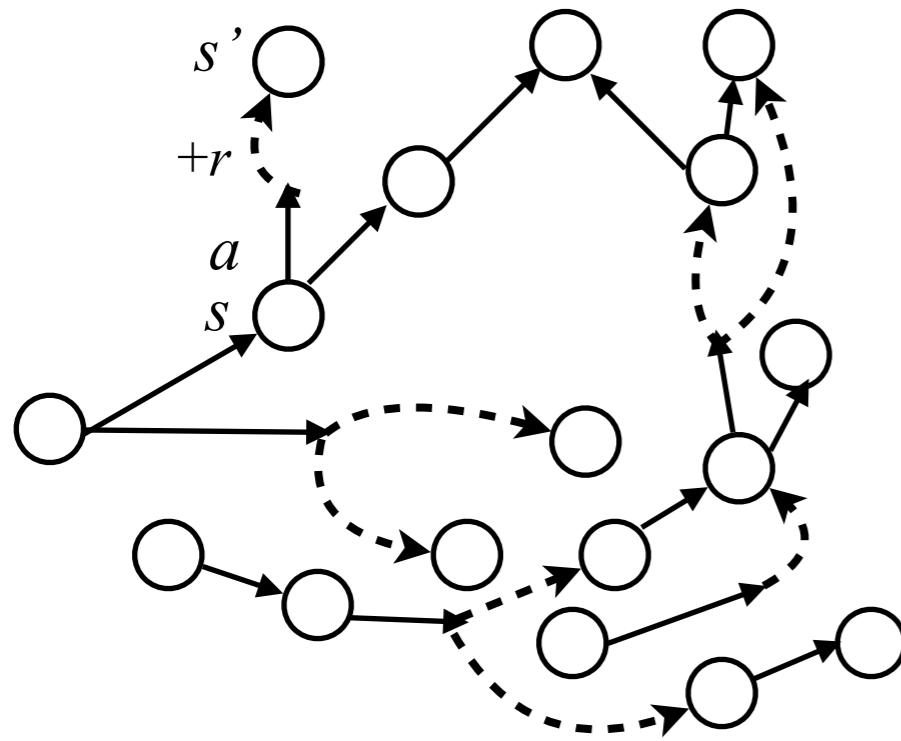


Doubly robust [JL'16]

- Even **perfect** control variate cannot eliminate **exponential** variance!

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

Nan Jiang, Lihong Li. ICML-16. Doubly Robust Off-policy Value Evaluation for Reinforcement Learning.



Training: $\hat{f} = f_k$ where

(FQI: learn Q^*)

$$f_k \leftarrow \arg \min_{f_{\theta}} \mathbb{E}_D [(f_{\theta}(s, a) - r - \gamma \max_{a'} f_{k-1}(s', a'))^2]$$

π = greedy w.r.t. \hat{f}

Reformulation: Value-function Selection

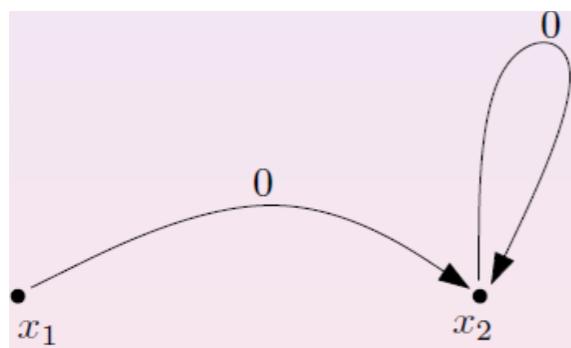
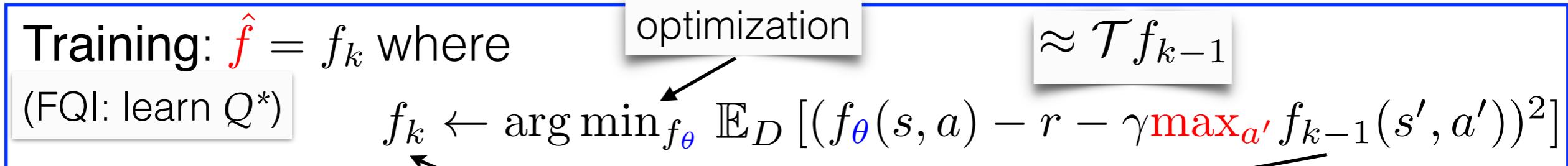
Simple(?) Problem

- Run different training algorithms
- Get candidate value functions f_1, f_2, \dots
- Holdout data $\{(s, a, r, s')\}$
- Select a good approx of Q^* w/ a “small” holdout dataset?
 - “small” = no $|S|$ or exponential-in-horizon
 - & no further function approximation!
- Simpler: identify Q^* out of f_1, f_2

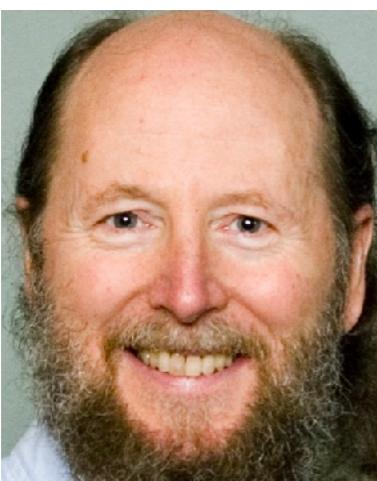


The training perspective

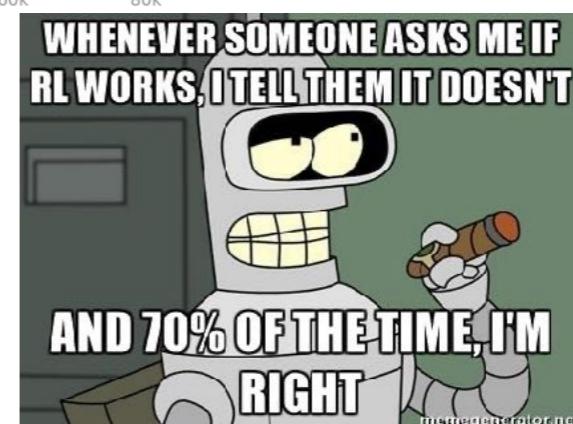
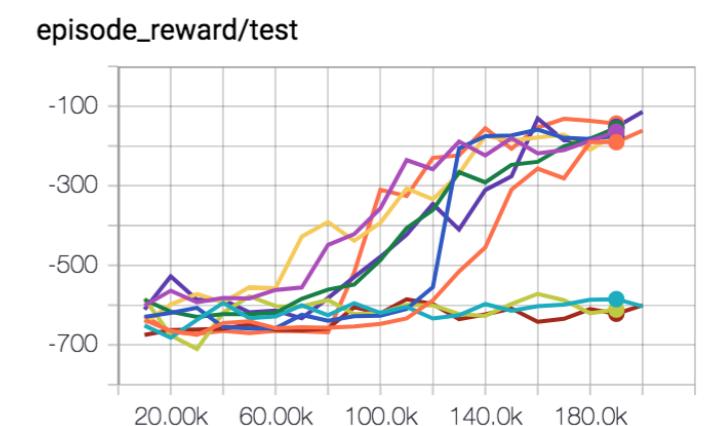
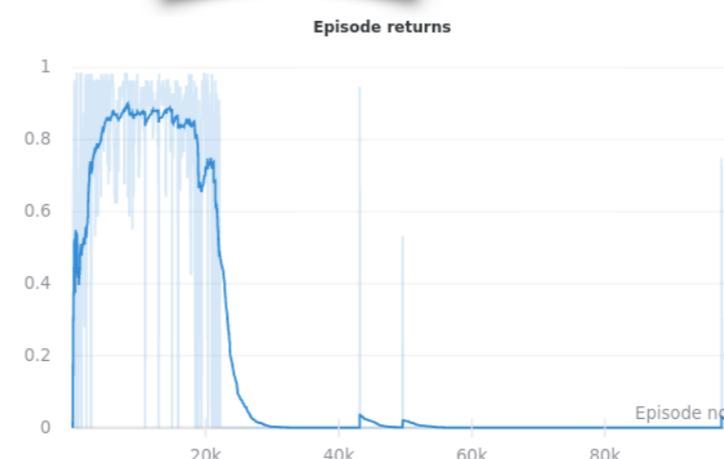
- Baird'95: design L s.t. $Q^* \stackrel{?}{=} \arg \min_{f \in \mathcal{F}} L(f)$
- RL **doesn't** work like that!



Divergence under 1-d linear
[TvR'96]



“Deadly triad”



The training perspective

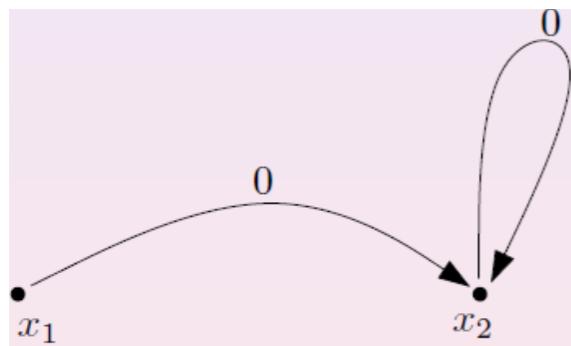
- Baird'95: design L s.t. $Q^* \stackrel{?}{=} \arg \min_{f \in \mathcal{F}} L(f)$
- RL **doesn't** work like that!

Training: $\hat{f} = f_k$ where

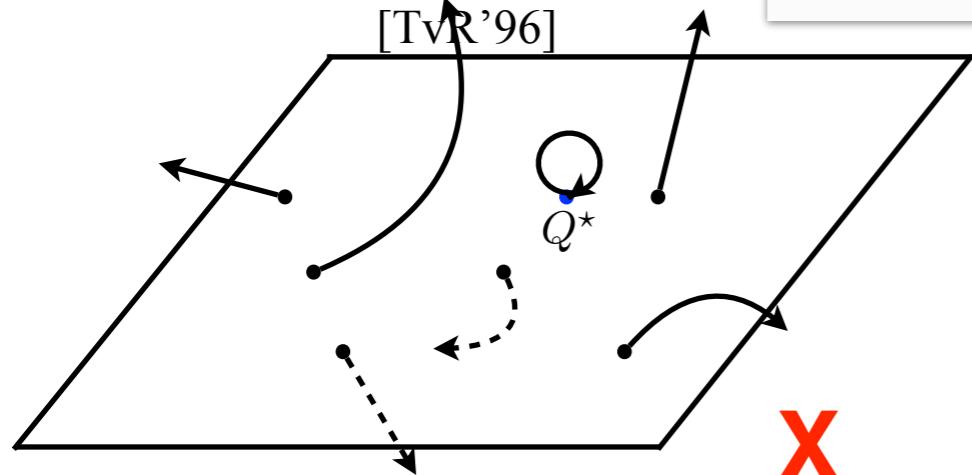
$\approx \mathcal{T}f_{k-1}$

(FQL: learn Q^*)

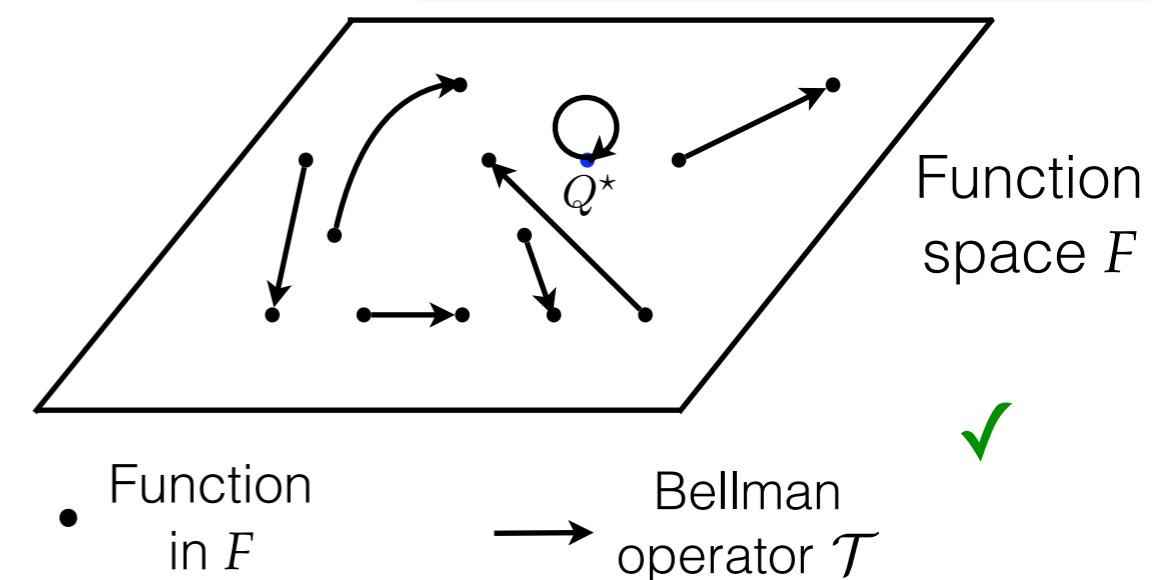
$$f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma \max_{a'} f_{k-1}(s', a'))^2]$$



Divergence under 1-d lin realizability (of Q^*)



“Bellman-completeness”
 $\mathcal{T}f \in \mathcal{F}, \forall f \in \mathcal{F}$



The training perspective

- Baird'95: design L s.t. $Q^* \stackrel{?}{=} \arg \min_{f \in \mathcal{F}} L(f)$
- $f = Q^* \Leftrightarrow f = \mathcal{T}f$, so how about

$$\begin{aligned} & f - \mathcal{T}f \\ &= \mathbb{E}_D [(f(s, a) - \mathbb{E}[r + \gamma \max_{a'} f(s', a') | s, a])^2] \\ &\quad \neq \\ &\quad \mathbb{E}_D [(f(s, a) - (r + \gamma \max_{a'} f(s', a'))))^2] \end{aligned}$$

- Naive “1-sample” estimator is **biased**
 - **debasing** requires **simulator** (“*double sampling*” [Baird’95])
 - or, **helper class** $\mathcal{F}' \ni \mathcal{T}f$ [ASM’08, FS’10]



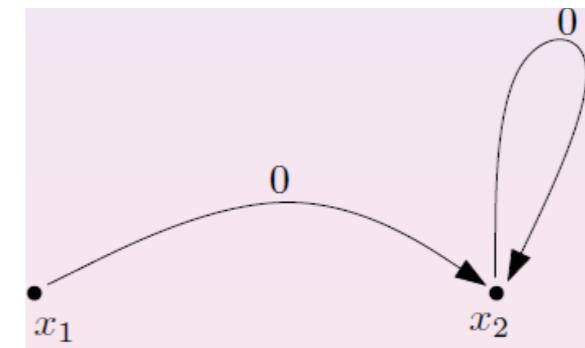
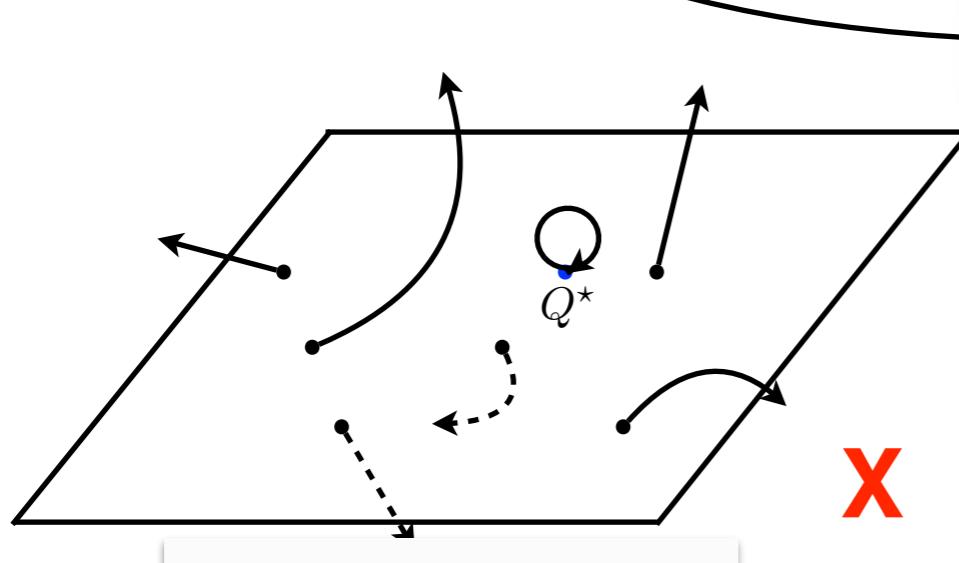
* over-estimate by a Bayes-error-like term: $\mathbb{E}_{d^D} [\mathbb{V}_{s'|s,a} [r + \gamma \max_{a'} f(s', a')]]$

Basis of resolution

Training: $\hat{f} = f_k$ where

(FQI: learn Q^*)

$$f_k \leftarrow \arg \min_{f_\theta \in \mathcal{F}} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma \max_{a'} f_{k-1}(s', a'))^2]$$



Divergence under 1-d linear \mathcal{F}
[TvR'96]

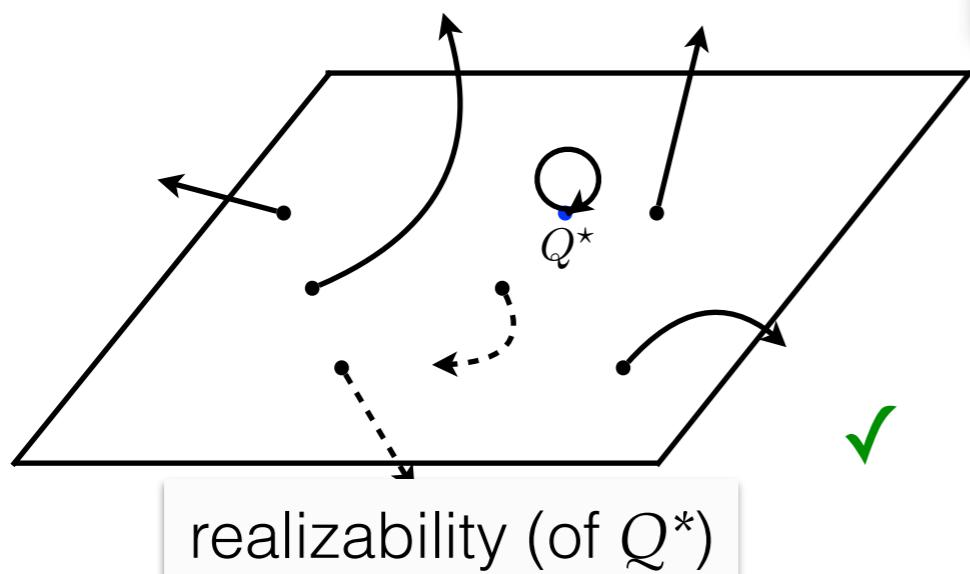
To select b/t f_1, f_2

Basis of resolution

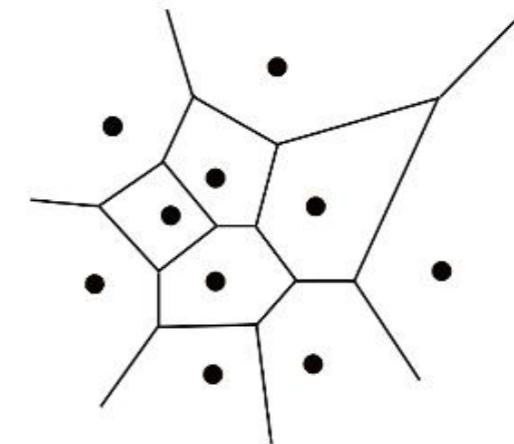
Training: $\hat{f} = f_k$ where

(FQL: learn Q^*)

$$f_k \leftarrow \arg \min_{f_\theta \in \mathcal{F}} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma \max_{a'} f_{k-1}(s', a'))^2]$$



iterative



Convergence under piecewise
constant \mathcal{F} ! [Gordon'95]



same

To select b/t f_1, f_2 , suffices to have class G s.t.

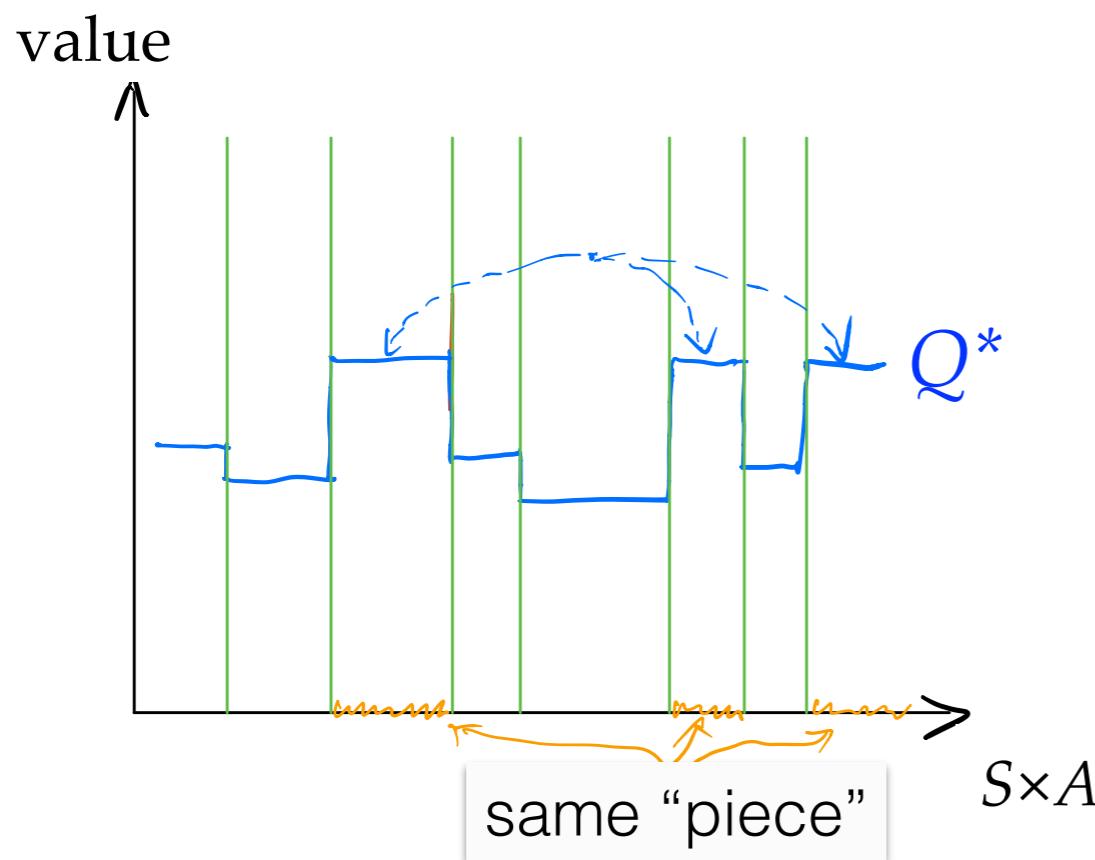
- piecewise constant
- can express Q^*
- small # partitions (bounded complexity)

Then: minimize $\|f - \text{Proj}_G(\mathcal{T}f)\|_{2,D}$

Our method: create such a
magical G “out of nothing”!

Does a **magical** G always exist?

- YES! Just partition $S \times A$ according to **output** of $\boxed{Q^*}$

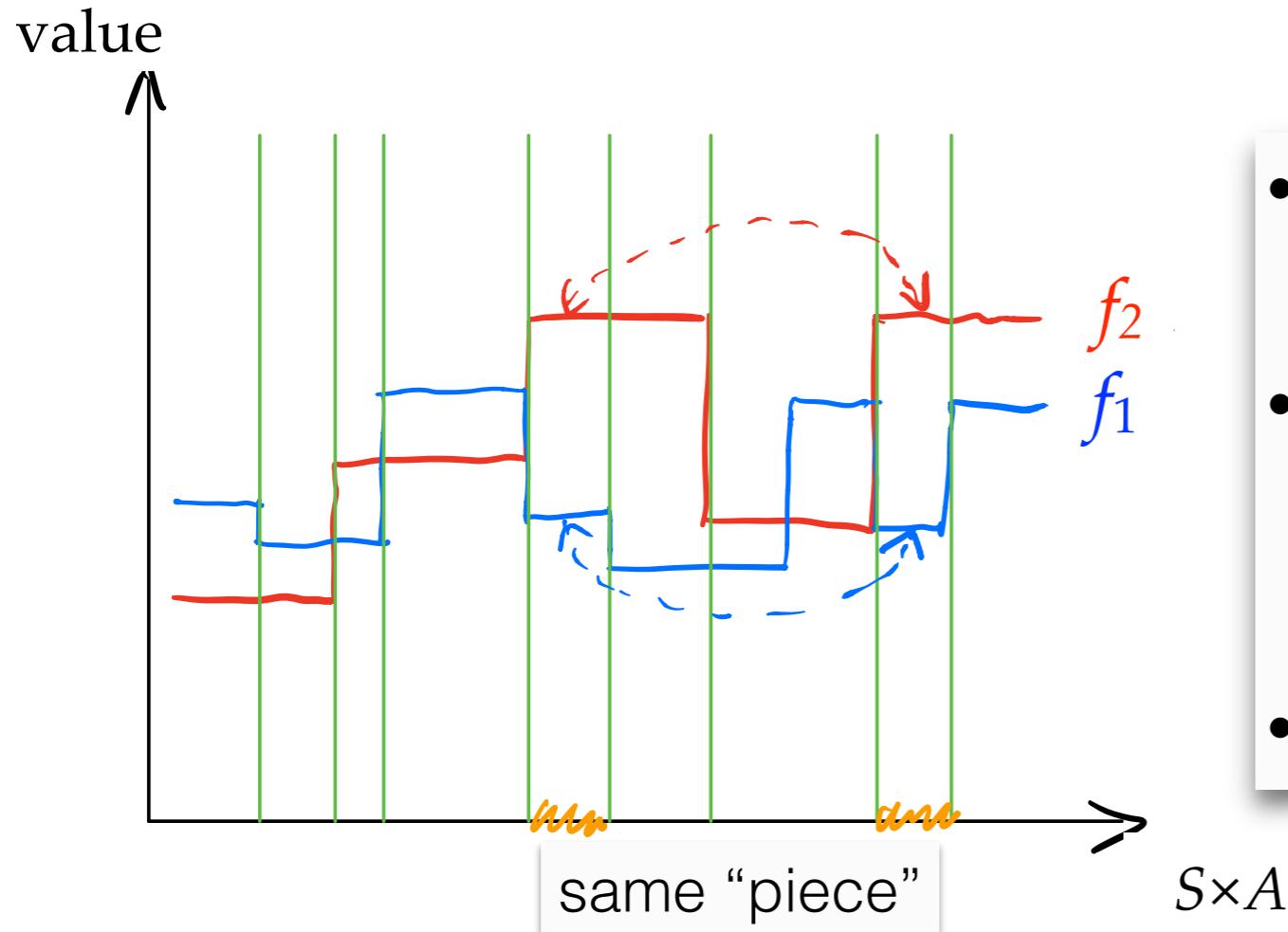


To select b/t f_1, f_2 , suffices to have class G s.t.

- piecewise constant ✓
- can express Q^* ✓
- $\boxed{O(1/\epsilon)}$ partitions (bounded complexity) ✓

Then: minimize $\|f - \text{Proj}_G(\mathcal{T}f)\|_{2,D}$

Batch Value-Function Tournament [XJ, ICML-21]

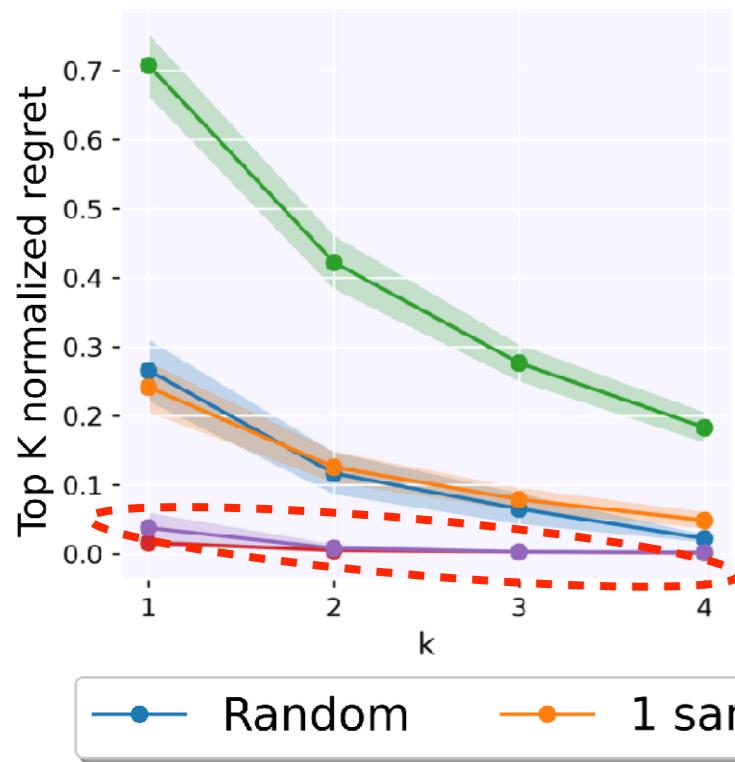


- Algorithm: **BVFT**
 $\arg \min_i \max_j \|f_i - \text{Proj}_{\mathcal{G}_{i,j}}(\mathcal{T}f_i)\|_{2,D}$
- Sample complexity poly in horizon, $1/\varepsilon$, $\log(\#\text{candidates})$, and C (data coverage)
- Computation: $\#\text{data points} * |F|^2$

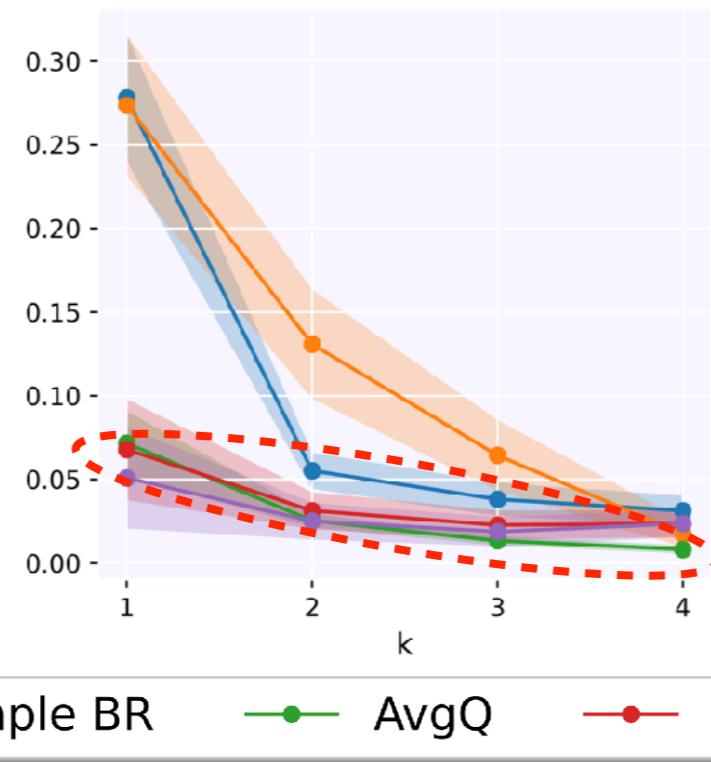
- (Simplified) problem: identify Q^* out of $F = \{f_1, f_2\}$
- Partition $S \times A$ according to both functions simultaneously!
 - Pw-const class $G_{1,2}$ w/ size $O(1/\varepsilon^2)$!!
- Naive extension to >2 functions in F : $O(1/\varepsilon^{|F|})$
 - Pairwise comparison + tournament

Formal
guarantee in
backup slide

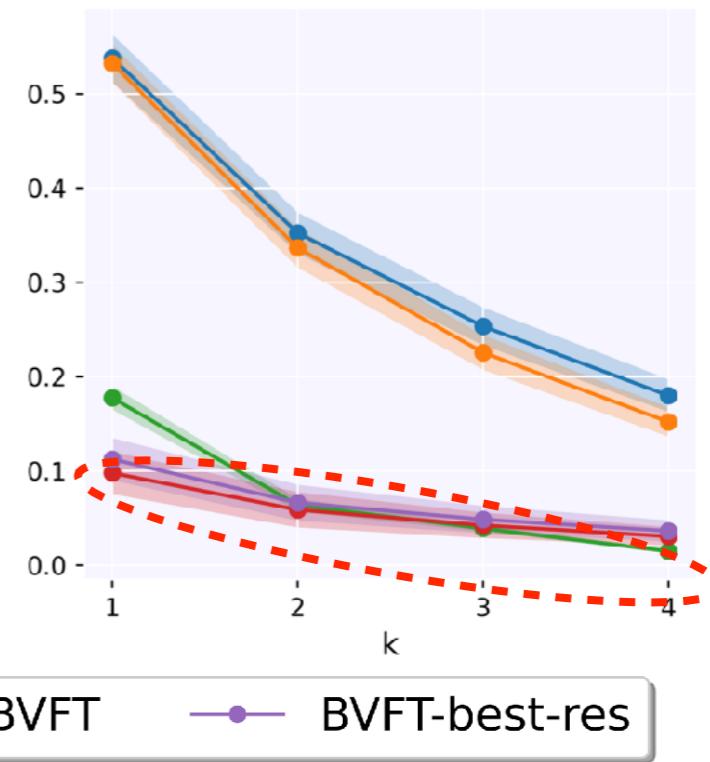
Acrobot-v1



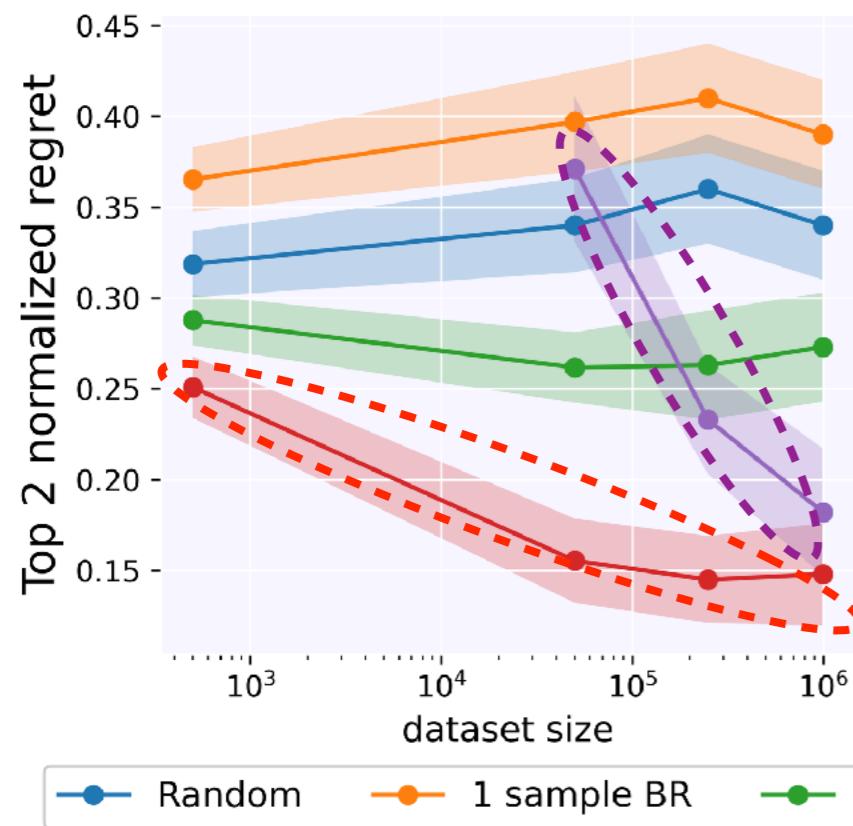
Pendulum-v0



LunarLander-v2



Asterix-v0



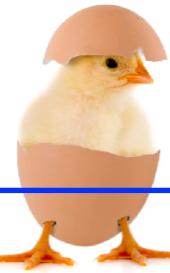


$$f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma f_{k-1}(s', \pi))^2]$$

Neural architecture
designed by “cheating”

Training: $\hat{f} = f_k$ where

(FQI: learn Q^*) $f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma \max_{a'} f_{k-1}(s', a'))^2]$



↓ π = greedy w.r.t. \hat{f}

Validation:

(FQE: learn Q^π) $f_k \leftarrow \arg \min_{f_\theta} \mathbb{E}_D [(f_\theta(s, a) - r - \gamma f_{k-1}(s', \pi))^2]$

- **BVFT**: H-P free solution for value-function selection
- Many open problems in validation
 - Data coverage issues (see lower bound [FKSX'22])
 - Combine with different OPE methods
 - e.g., marginalized importance sampling
[LLTZ'18, NCDL'19, UHJ'20, JH'20, VJY'21, HJ'22]
 - Practical toolkit (cf. for OPE [VLJY'20])