# What Structural Conditions Permit Generalization in Reinforcement Learning?



- Joint work with
- Simon Du, Sham Kakade, Jason Lee, Shachar Lovett,
  - Gaurav Mahajan, Ruosong Wang



# Solving Large-scale RL problems requires generalization

At last – a computer program that can beat a champion Go player PAGE 484

nature

**ALL SYSTEMS GO** 

[AlphaZero, Silver et.al, 17]



[OpenAl Five, 18]

[OpenAl, 19]



## Markov Decision Processes: a framework for RL

- A policy:  $\pi$ : States  $\rightarrow$  Actions
- Execute  $\pi$  to obtain a trajectory:  $S_0, a_0, r_0, S_1, a_1, r_1 \dots S_{H-1}, a_{H-1}, r_{H-1}$
- Cumulative *H*-step reward:  $V_{H}^{\pi}(s) = \mathbb{E}_{\pi} \left[ \sum_{t=0}^{H-1} r_{t} \middle| s_{0} = s \right], \quad Q_{H}^{\pi}(s,a) = \mathbb{E}_{\pi} \left[ \sum_{t=0}^{H-1} r_{t} \middle| s_{0} = s, a_{0} = a \right]$
- Goal: Find a policy  $\pi$  that maximizes our value  $V^{\pi}(s_0)$  from  $s_0$ . Episodic setting: We start at  $s_0$ ; act for H steps; repeat...



# Provable Generalization in Supervised Learning (SL)

Generalization is possible in the IID supervised learning setting!

To get  $\epsilon$ -close to best in hypothesis class  $\mathcal{F}$ , we need # of samples that is: • "Occam's Razor" Bound (finite hypothesis class): need  $O(\log |\mathcal{F}|/\epsilon^2)$ 

- Various Improvements:
- VC dim: need only  $O(VC(\mathcal{F})/\epsilon^2)$
- Classification: linearly separable + margin:  $O(\text{margin})/\epsilon^2$ )
- Linear Regression in d dimensions:  $O(d/\epsilon^2)$
- Deep Learning: the algorithm also determines the complexity control

# The key idea in SL: data reuse

With a training set, we can simultaneously evaluate the loss of all hypotheses in our class!

# Sample Efficient RL in the Tabular Case (no generalization here)

- Thm: In the episodic setting,  $poly(S, A, H, 1/\epsilon)$ samples suffice to find an  $\epsilon$ -opt policy with the  $E^3$  algo. [Kearns & Singh '98] also: [Brafman& Tennenholtz '02; K. '03] Key idea: optimism + dynamic programming
- Regret guarantees with model based algos: [Auer+ '09]
- Provable Q-learning (+bonus): [Strehl+ (2006)], [Szita & Szepesvari '10],[Jin
- (asymptotically) optimal reget: Reg(#episod [Azar+ '17],[Dann+'17]

0 2 3 Wall Start 0 Wall 1 Wall 2 Wall 3

$$(+ `18]$$
  
 $(les) = \sqrt{HSA \cdot #episodes}$ 



## I: Provable Generalization in RL Q1: Can we find an $\epsilon$ -opt policy with no S dependence?

- How can we reuse data to estimate the value of all policies in a policy class  $\mathcal{F}$ ? Idea: Trajectory tree algo dataset collection: uniformly at random choose actions for all H steps in an episode. estimation: uses importance sampling to evaluate every  $f \in \mathcal{F}$
- Thm: [Kearns, Mansour, & Ng '00]
  - To find an  $\epsilon$ -best in class policy, the trajectory tree algo uses  $O(A^H \log(|\mathcal{F}|)/\epsilon^2)$  samples • Only  $log(|\mathcal{F}|)$  dependence on hypothesis class size.
  - There are VC analogues as well.
- Can we avoid the  $2^H$  dependence to find an  $\epsilon$ -best-in-class policy? Agnostically, **NO**! **Proof:** Consider a binary tree with  $2^{H}$ -policies and a sparse reward at a leaf node.



# II: Provable Generalization in RL

• Q2: Can we find an  $\epsilon$ -opt policy with no S, A dependence and  $poly(H,1/\epsilon,$  "complexity measure") samples?

### •With various stronger assumptions, yes.

- Linear Bellman Completion: [Munos, '05, Zanette+ '19]
  - Linear MDPs: [Wang & Yang'18]; [Jin+'19] (the transition matrix is low rank)
  - Linear Quadratic Regulators (LQR): standard control theory model
- FLAMBE / Feature Selection: [Agarwal, K., Krishnamurthy, Sun '20]
- Linear Mixture MDPs: [Modi+'20, Ayoub+ '20] •
- Block MDPs [Du+ '19]
- Factored MDPs [Sun+ '19]
- Kernelized Nonlinear Regulator [K.+ '20]





## **Structural conditions Under which Generalization is possible**

## This talk:



Given any  $w \in \mathbb{R}^d$ , there exists a  $\theta \in \mathbb{R}^d$ , such that:  $a) + \mathbb{E}_{s' \sim P(s,a)} \max_{a'} w^{\mathsf{T}} \phi(s',a')$ 

$$\forall s, a : \theta^{\mathsf{T}} \phi(s, a) = r(s, a)$$
$$:= T(w)$$

Linear MDPs, Linear Quadratic Regulator

**Def**: Given feature map  $\phi(s, a) \in \mathbb{R}^d$ , Bellman operator has linear closure

Examples:



- **Generalization is possible here:**
- $\exists$  an algorithm, finding  $\epsilon$ -near optimal policy only needs  $poly(H, d, 1/\epsilon)$  many samples



## what's the structure here that permits generalization?

- We can rewrite the average Bellman error (averaged over any roll-in  $\pi$ ) in a bilinear form:
  - Given a  $Q^*$  candidate

$$\mathbb{E}_{s,a\sim\pi}\left[w^{\mathsf{T}}\phi(s,a) - r(s,a) - \mathbb{E}_{s'\sim P(s,a)}\max_{a'}w^{\mathsf{T}}\phi(s',a')\right]$$
$$= \mathbb{E}_{s,a\sim\pi}\left[w^{\mathsf{T}}\phi(s,a) - T(w)^{\mathsf{T}}\phi(s,a)\right] = \left\langle w - T(w), \mathbb{E}_{s,a\sim\pi}\phi(s,a)\right\rangle$$

## Warm up

$$w^{\mathsf{T}}\phi(s,a)$$
, we have:



### what's the unique structure here that permits generalization?

Define discrepancy  $\ell(s, a, s', w) =$ 

 $\mathbb{E}_{s,a\sim\pi}\ell(s,a,s',w) =$ 

We can also estimate the value of the bilinear form:

$$= w^{\mathsf{T}} \phi(s, a) - r(s, a) - \max_{a'} w^{\mathsf{T}} \phi(s', a')$$

We have:

$$\langle w - T(w), \mathbb{E}_{s,a\sim\pi}\phi(s,a) \rangle$$

Note we have data reuse: given data from  $\pi$ , we can evaluate all w



$$\mathbb{E}_{s,a\sim\pi}\left[w^{\mathsf{T}}\phi(s,a) - r(s,a) - \mathbb{E}_{s'\sim P(s,a)}\max_{a'}w^{\mathsf{T}}\phi(s',a')\right] = \left\langle w - T(w), \mathbb{E}_{\pi}\phi(s,a)\right\rangle$$

$$\mathbb{E}_{s,a\sim\pi}\ell(s,a,s',w) = \left\langle w - T(w), \mathbb{E}_{\pi}\phi(s,a) \right\rangle$$

Note that the analytical form of bilinear structure is unknown

### In summary, it has a **bilinear structure**:

For any roll-in policy  $\pi$ , and any w, we have: (1)

AND (2) there exists a discrepancy function  $\ell$ , s.t.,

# BiLinear Classes: structural properties to enable generalization in RL

- Realizable Hypothesis class:  $\{f \in \mathcal{F}\},\$ 
  - can be model based or model-free class.

**Def:** A ( $\mathcal{F}, \ell$ ) forms **an (implicit) Bilinear class** class if there are  $W_h \in \mathcal{F} \mapsto \mathcal{H}, \& X_h \in \mathcal{F} \mapsto \mathcal{H}$  ( $\mathcal{H}$  being some Hilbert space):

- $\left| \mathbb{E}_{s_{h},a_{h}\sim\pi_{f}} \left[ Q_{f}(s_{h},a_{h}) r(s_{h},a_{h}) V_{f}(s_{h+1}) \right] \right| \leq \left\langle W_{h}(f) W_{h}(f^{\star}), X_{h}(f) \right\rangle$  $\mathbb{E}_{s_h \sim \pi_f, a_h \sim \pi_{est}} \Big[ \ell_f(s_h, a_h, s_{h+1}, g) \Big] = \Big\langle W_h(g) - W_h(f^\star), X_h(f) \Big\rangle, \forall g \in \mathcal{F}$
- Bilinear regret: on-policy difference between claimed reward and true reward • Data reuse: there is discrepancy function  $\ell_f(s, a, s', g)$  & policy  $\pi_{est}$  s.t.

with associated state-action value, (greedy) value and policy:  $Q_f(s, a), V_f(s), \pi_f$ 

Note:  $W_h \& X_h$  are implicit—no need to known them

## **Back to Linear Bellman Complete:**

$$\mathbb{E}_{s_h, a_h \sim \pi_f} \left[ w^{\mathsf{T}} \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s' \sim P(s, a)} \max_{a'} w^{\mathsf{T}} \phi(s', a') \right] = \left\langle \underbrace{(w - T(w))}_{W_h(f)} - \underbrace{(w^{\star} - T(w^{\star}))}_{W_h(f^{\star})}, \underbrace{\mathbb{E}_{\pi} \phi(s, a)}_{X_h(f)} \right\rangle$$

$$\mathbb{E}_{\pi_f} \mathscr{E}(s, a, s', w') = \left\langle (w' - T(w')) - (w^* - T(w^*)), \mathbb{E}_{s, a \sim \pi_f} \phi(s, a) \right\rangle$$

Note  $W_h(f)$  is unknown as the Bellman backup T(w) is unknown

(1) **Bilinear regret:** for any  $f(s, a) := w^{\top} \phi(s, a)$ , we have:

AND (2) **data-reuse**: there exists a discrepancy function  $\ell$ , s.t.,

# The Algorithm: Bilin-UCB

For  $t = 0 \rightarrow T$ :

- Find the "optimistic"  $f_t \in \mathcal{F}$ :  $\underset{f \in \mathscr{F}}{\operatorname{arg\,max}} V_f(s_0), \text{ s.t., } \sigma_h^2(f) \leq R, \forall h$
- Sample *m* trajectories  $\pi_{f}$  and create a batch dataset:  $D = \{(s_h, a_h, s_{h+1}) \in \text{trajectories}\}$
- Update the cumulative discrepancy function  $\sigma_h(\cdot), \forall h$

$$\sigma_h^2(\ \cdot\ ) \leftarrow \sigma_h^2(\ \cdot\ ) +$$



Note here we roughly have:  $\sigma_h(g) \approx \sum_{k=1}^{r} \left( \mathbb{E}_{s_h, a_h \sim \pi_{f_i}} \mathscr{C}_{f_i}(s_h, a_h, s_{h+1}, g) \right)^2$ i=0



# Theorem 2: Generalization in RL

- Theorem: [Du, Kakade., Lee, Lovett, Mahajan, S, Wang '21] Assume  $\mathscr{F}$  is a bilinear class and the class is realizable, i.e.  $f^* \in \mathscr{F}$ . Using  $\gamma_T^3 \cdot poly(H) \cdot \log(1/\delta)/\epsilon^2$  trajectories, the BiLin-UCB algorithm returns an  $\epsilon$ -opt policy (with prob.  $\geq 1 - \delta$ ).
  - $\gamma_T$  is the max. info. gain  $\gamma_T := \max_{h, f_0 \dots f_{T-1} \in \mathscr{F}} \ln \det \left( I + \frac{1}{\lambda} \sum_{t=0}^{T-1} X_h(f_t) X_h(f_t)^T \right)$
  - $\gamma_T \approx d \log T$  for  $X_h$  in *d*-dimensions

• The proof is "elementary" using the elliptical potential function. [Dani, Hayes, K. '08]

# Proof Sketch

- 3. If  $\pi_{f_t}$  is really sub-optimal, i.e.,  $V^*(s_0) V^{\pi_{f_t}}(s_0) \ge \epsilon$ , then  $f_t$  has large bilinear regret:
- 4. Recall in the alg, we have a constraint  $\sigma_h($

 $\sum_{k=1}^{n-1} \left( W_h(f_t) - W_h(f^\star), X_h(f_i) \right)^2 \le R$ i=0

 $\Rightarrow (W_h(f_t) - W_h(f^*))^\top \Sigma_{t:h}(W_h(f_t) - W_h(f^*))^\top \Sigma_{t:h}(W_h(f^*))^\top \Sigma_{t:h}(W_h(f^*))^\top$ 

1. Optimism:  $V^{\star}(s_0) \leq V_{f_t}(s_0)$ ; this can be verified by showing  $f^{\star}$  is always a feasible solution

2. Using optimism, we can upper bound per-episode regret Bilinear form ("simulation" lemma):  $V^{\star}(s_0) - V^{\pi_{f_t}}(s_0) \le V_{f_t}(s_0) - V^{\pi_{f_t}}(s_0) \le \sum_{t=1}^{H-1} |W_h(f_t) - W_h(f^{\star}), X_h(f_t)|$ 

 $\exists h, \text{ s.t., } |W_h(f_t) - W_h(f^*), X_h(f_t)| \ge \epsilon/H$ 

$$f_t \le R$$
, i.e.,  $\sum_{i=0}^{t-1} \left( \mathbb{E}_{s_h, a_h \sim \pi_{f_i}} \left[ \ell_{f_i}(s_h, a_h, s_{h+1}, f_t) \right] \right)^2$ 

$$(\Sigma_{h;t} := \sum_{i=0}^{t-1} X_h(f_i) X_h(f_i)^{\mathsf{T}} + \lambda I)$$





# **Proof Sketch**

- 3. If  $\pi_{f_t}$  is really sub-optimal, i.e.,  $V^*(s_0) V^{\pi_{f_t}}(s_0) \ge \epsilon$ , then  $f_t$  has large bilinear value:
- 5. Finally, combine the two results and using Cauchy-Schwartz, we have:

$$\epsilon/H \leq \left\| W_{h}(f_{t}) - W_{h}(f^{\star}) \right\|_{\Sigma_{h;t}} \left\| X_{h}(f_{t}) \right\|_{\Sigma_{h;t}^{-1}} \leq \sqrt{R + \lambda} \left\| X_{h}(f_{t}) \right\|_{\Sigma_{h;t}^{-1}}$$
$$\Rightarrow \left\| X_{h}(f_{t}) \right\|_{\Sigma_{h;t}^{-1}} \geq \frac{\epsilon}{H\sqrt{R + \lambda}}$$

 $\exists h, \text{ s.t., } \left| W_h(f_t) - W_h(f^*), X_h(f_t) \right| \ge \epsilon/H$ 4. Recall in the alg, we have a constraint  $\sigma(f_t) \le R$ , i.e.,  $\sum_{k=1}^{t-1} \left( \mathbb{E}_{s_h, a_h \sim \pi_{f_\tau}} \mathscr{C}_{f_\tau}(s_h, a_h, s_{h+1}, f_t) \right)^2 \le R$  $\tau = 0$   $(W_h(f_t) - W_h(f^*))^\top \Sigma_{t;h}(W_h(f_t) - W_h(f^*)) \le R + \lambda \qquad \left(\Sigma_{h;t} \coloneqq \sum_{i=1}^{t-1} X_h(f_i) X_h(f_i)^\top + \lambda I\right)$ 

In *d*-dimensional setting, such event cannot happen more than  $O(d/\epsilon^2)$  many times....





## Theorem: [Du, Kakade., Lee, Lovett, Mahajan, S, Wang '21] The following models are bilinear classes for some discrepancy function $\ell(\cdot)$

- Linear Bellman Completion: [Munos, '05, Zanette+ '19]
  - Linear MDPs: [Wang & Yang'18]; [Jin+ '19] (the transition matrix is low rank)
  - Linear Quadratic Regulators (LQR): standard control theory model
- FLAMBE / Feature Selection: [Agarwal, K., Krishnamurthy, Sun '20]
- Linear Mixture MDPs: [Modi+'20, Ayoub+ '20] lacksquare
- Block MDPs [Du+ '19]
- Factored MDPs [Sun+ '19]
- Kernelized Nonlinear Regulator [K.+ '20]
- Linear  $Q^{\star} \& V^{\star}$
- Reactive PSR/POMDP, Generalized linear MDP, and more.....
- two exceptions: deterministic linear  $Q^{\star}$ ;  $Q^{\star}$ -state-action aggregation

 (almost) all "named" models (with provable generalization) are bilinear classes • Bilinear classes generalize the: Bellman rank [Jiang+ '17]; Witness rank [Sun+ '19]





## **Another Example: Feature Selection for Low-rank MDP**

## The feature selection problem:

- Low-rank MDP:  $P(s'|s, a) = \mu^*(s')^\top \phi^*(s, a), \ \mu^* \& \phi^*$  unknown
  - Function approximation for feature:  $\phi^* \in \Psi \subset S \times A \mapsto \mathbb{R}^d$
  - Function approximation for  $Q^*: Q := \{w^\top \phi(s, a) : w \in \text{Ball}_W, \phi \in \Psi\}$
  - Q: can we find  $\epsilon$  optimal policy w/ samples poly $(d, \ln(\Psi), H, 1/\epsilon)$ ?
    - Yes & it's in the Bilinear class!



## **Another Example: Feature Selection for Low-rank MDP**

## The feature selection problem:

$$\forall f \in \mathcal{Q} : \mathbb{E}_{s_{h},a_{h} \sim \pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$\pi_{f} \mathbb{E}_{s_{h} \sim P(\cdot|s_{h-1},a_{h-1}),a_{h} \sim \pi_{f}(\cdot|s_{h})} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$\int_{s_{h}} \phi^{\star}(s_{h-1},a_{h-1})^{\top} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h} \sim \pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$\pi_{f'} \phi^{\star}(s_{h-1},a_{h-1}), \quad \int_{s_{h}} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h} \sim \pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$\forall f \in \mathcal{Q} : \mathbb{E}_{s_{h},a_{h} \sim \pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$= \mathbb{E}_{s_{h-1},a_{h-1} \sim \pi_{f}} \mathbb{E}_{s_{h} \sim P(\cdot|s_{h-1},a_{h-1}),a_{h} \sim \pi_{f}(\cdot|s_{h})} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$= \mathbb{E}_{s_{h-1},a_{h-1} \sim \pi_{f}} \int_{s_{h}} \phi^{\star}(s_{h-1},a_{h-1})^{\top} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h} \sim \pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$= \left\langle \mathbb{E}_{s_{h-1},a_{h-1} \sim \pi_{f}} \phi^{\star}(s_{h-1},a_{h-1}), \int_{s_{h}} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h} \sim \pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right] \right\rangle$$

$$\forall f \in \mathcal{Q} : \mathbb{E}_{s_{h},a_{h}\sim\pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$= \mathbb{E}_{s_{h-1},a_{h-1}\sim\pi_{f}} \mathbb{E}_{s_{h}\sim P(\cdot|s_{h-1},a_{h-1}),a_{h}\sim\pi_{f}(\cdot|s_{h})} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$= \mathbb{E}_{s_{h-1},a_{h-1}\sim\pi_{f}} \int_{s_{h}} \phi^{\star}(s_{h-1},a_{h-1})^{\mathsf{T}} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h}\sim\pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right]$$

$$= \left\langle \underbrace{\mathbb{E}_{s_{h-1},a_{h-1}\sim\pi_{f}} \phi^{\star}(s_{h-1},a_{h-1})}_{\mathbb{E}_{s_{h}}} \int_{s_{h}} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h}\sim\pi_{f}} \left[ Q_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} Q_{f}(s_{h+1},a') \right] \right\rangle$$

$$\begin{aligned} \forall f \in \mathcal{Q} : \mathbb{E}_{s_{h},a_{h}\sim\pi_{f}} \left[ \mathcal{Q}_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} \mathcal{Q}_{f}(s_{h+1},a') \right] \\ &= \mathbb{E}_{s_{h-1},a_{h-1}\sim\pi_{f}} \mathbb{E}_{s_{h}\sim P(\cdot|s_{h-1},a_{h-1}),a_{h}\sim\pi_{f}(\cdot|s_{h})} \left[ \mathcal{Q}_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} \mathcal{Q}_{f}(s_{h+1},a') \right] \\ &= \mathbb{E}_{s_{h-1},a_{h-1}\sim\pi_{f}} \int_{s_{h}} \phi^{\star}(s_{h-1},a_{h-1})^{\top} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h}\sim\pi_{f}} \left[ \mathcal{Q}_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} \mathcal{Q}_{f}(s_{h+1},a') \right] \\ &= \left\langle \underbrace{\mathbb{E}_{s_{h-1},a_{h-1}\sim\pi_{f}} \phi^{\star}(s_{h-1},a_{h-1})}_{X_{h}(f)}, \underbrace{\int_{s_{h}} \mu^{\star}(s_{h}) \mathbb{E}_{a_{h}\sim\pi_{f}} \left[ \mathcal{Q}_{f}(s_{h},a_{h}) - r_{h} - \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} \mathcal{Q}_{f}(s_{h+1},a') \right]}_{A'} \right\rangle \end{aligned}$$

1. Claim: on-policy Bellman error of  $f := w^{\top} \phi$  has Bilinear form:



## **Another Example: Feature Selection for Low-rank MDP**

### The feature selection problem:

2. Claim: The bilinear regret is estimable using some  $\ell$ 





- We have linear MDP, Linear Bellman complete...
- The most natural one should just be  $Q^*$  being linear-realizable:
  - $Q^{\star}(s,a) = (w^{\star})^{\top} \phi(s,a)$ , under known feature  $\phi$ 
    - However generalization is **impossible** here:
- optimal policy

Another Example: Linear  $Q^* \& V^*$ 

What's the role of linear function approximation in RL?

• Theorem [Weisz, Amortila, Szepesvári '21]: There exists an MDP w/ linear  $Q^*$ , s.t any online RL algorithm requires  $\Omega(\min(2^d, 2^H))$  samples to output a near



## What's the role of linear function approximation in RL?

 Theorem [Du, Kakade., Lee, Lovett, Mahajan, S, Wang '21] For any MDP with both  $Q^{\star}$  and  $V^{\star}$  being realizable (under some known)  $d^{3}$ poly(H)  $\frac{1}{\epsilon^2}$ 

Additional Example: Linear Q\* & V\*

However, if we further assume  $V^{\star}(s) = (\theta^{\star})^{\top} \psi(s)$ , then we will be ok:

features, e.g., RKHS), there exists an algorithm that learns with # of samples:



# Additional Example: Linear Q\* & V\*

Key step: pre-process function class

$$\left\{Q,V\right\} := \left\{(w,\theta) : \forall s, \max_{a} w^{\mathsf{T}}\phi(s,a) = \theta^{\mathsf{T}}\psi(s)\right\}$$

 $v(s') \bigg| = \langle W_h([w,\theta]) - W_h([w^{\star},\theta^{\star}]), X_h([w,\theta]) \rangle$ 

$$\mathbb{E}_{s_h,a_h\sim\pi_f}\left[w^{\mathsf{T}}\phi(s_h,a_h)-r_h-\mathbb{E}_{s'\sim P(s_h,a_h)}\theta^{\mathsf{T}}\psi(s_h,a_h)-r_h-\mathbb{E}_{s'\sim P(s_h,a_h)}\theta^{\mathsf{T}}\psi(s_h,a_h)\right]$$

(1) on-policy Bellman error of any  $f := (w, \theta)$  has bilinear form: (2) there exists  $\ell(s, a, s', [w, \theta]) = w^{\top} \phi(s, a) - r(s, a) - \theta^{\top} \psi(s')$ , s.t.,

 $\mathbb{E}_{s_h,a_h\sim\pi_f}\left[\ell(s_h,a_h,s_{h+1}',[w',\theta'])\right] = \langle W_h([w',\theta']) - W_h([w^\star,\theta^\star]), X_h(f) \rangle, \forall [w',\theta'] \rangle$ 

### Linear $Q^* \& V^*$ has Bilinear Structure

Function classes for  $Q^* \in Q := \{ w^\top \phi(s, a) : w \in \text{Ball}_w \}, V^* \in \mathcal{V} := \{ \theta^\top \psi(s) : \theta \in \text{Ball}_B \}$ 



# Final Example: Linear Mixture Model (model-based) **The linear Mixture Model:** Function class: $\mathscr{F} = \{ P : P(s' | s, a; \theta) = \theta^{\mathsf{T}} \phi(s, a, s'), \theta \in \text{Ball}_W \}$

- Imagine we have d simulators,  $P^{i}(s'|s, a), i \in [d]$ ,
- we assume the ground truth is the linear mixture of d simulators:

$$P^{\star}(s'|s,a) = \sum_{i=1}^{d} \theta^{\star}[i]P^{i}(s'|s,a)$$

Why this model is ever interesting?

# Final Example: Linear Mixture Model (model-based) **The linear Mixture Model:**

Function class:  $\mathcal{F} = \{P : P(s')\}$ 

(Notation:  $f \in \mathcal{F}$  is a potential transition, and  $f^* := P^*$ )

Claim 1: For any  $f \in \mathcal{F}$ , the on-policy Bellman error of  $Q_f$  under  $\pi_f$  has bilinear form:  $\mathbb{E}_{\pi_f}\left[Q_f(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s' \sim f^\star(s_h, a_h)}V_f(s')\right] = \left\langle W_h(f) - W_h(f^\star), X_h(f) \right\rangle$ 

$$|s, a; \theta) = \theta^{\mathsf{T}} \phi(s, a, s'), \theta \in \mathsf{Ball}_W \}$$

## Final Example: Linear Mixture Model

### **The linear Mixture Model:**

- Function class:  $\mathcal{F} = \{P : P(s')\}$

$$\ell_f(s, a, s', g) = \mathbb{E}_s$$

s.t., 
$$\mathbb{E}_{\pi_f} \mathscr{C}_f(s_h, a_h, s'_{h+1}, g)$$

$$|s, a; \theta) = \theta^{\mathsf{T}} \phi(s, a, s'), \theta \in \mathsf{Ball}_W$$

(Notation:  $f \in \mathcal{F}$  is a potential transition, and  $f^* := P^*$ )

Claim 2: For any  $f \in \mathcal{F}$ , there exists discrepancy  $\ell_f(s, a, s', g)$  to measure bilinear form:  $E_{s' \sim g(s,a)} \left| V_f(s') - V_f(s') \right|$  $= \langle W_h(g) - W_h(f^{\star}), X_h(f) \rangle$