

QUESTION: Why use function approximation?  
 State ( $s$ ) can be very high-dimensional (eg, image)

QUESTION: what do we approximate?

RECALL: Linear Bandits

Assumption:  $r(s, a) = \phi(s, a)^T \theta^*$ ,  $\theta^* \in \mathbb{R}^d$

Why? Knowledge of  $r \mapsto$  Knowledge of  $\pi^*$  (optimal policy)

Notice: Knowledge of  $Q^* \mapsto$  Knowledge of  $\pi^*$

$$Q_n^*(s, a) = \mathbb{E}_{\pi^*} \left[ \sum_{t=n}^H r_t \mid (s_n, a_n) = (s, a) \right]$$

Is Linear  $Q^*$  sufficient to learn sample efficiently?

	(LINEAR) BANDITS	$Q^*$ LINEAR
UPPER BOUND	$\tilde{O}(d\sqrt{k})$ (LINUCB) $\leftarrow$ # episodes	$O(e^H)$
LOWER BOUND	$\Omega(d\sqrt{k})$	$\Omega(e^d)$

**NO!**  $\Rightarrow$  Need additional assumptions!

LINEAR MDPs  $\approx$  MDPs that can be decomposed  
as a sequence of Linear Bandits

$$\text{Def: } \begin{cases} r_h(s,e) = \phi(s,e)^T \theta_h^R \\ P_h(s'|s,e) = \phi(s,e)^T \Psi(s') \end{cases}$$

Also called "Low Rank" MDPs:

$$\begin{array}{c} s' \\ \boxed{\begin{array}{c} | \\ \hline p(s'|s,\pi(s)) \\ \hline \end{array}} \\ s \end{array} = \begin{array}{c} \boxed{\phi(s,\pi(s))} \\ \hline \boxed{\Psi(s')} \end{array} \quad \boxed{|\Psi(s')|}$$

$$\forall \pi: P^\pi = \Phi^\pi \Psi^T \quad \left[ \begin{array}{l} TV = R + PV \\ = (\Phi \theta^R + \Phi \Psi V) \in \text{Range } \mathbb{I} \end{array} \right]$$

### FOUNDAMENTAL PROPERTY OF LINEAR MDPs

Def: Bellman Operator  $(T_h V)(s,e) = r_h(s,e) + \sum_{s'} P_h(s'|s,e) V(s')$

PROPOSITION

For any value function  $V$ ,  $(T_h V)(s,e) = \phi(s,e)^T \theta$  for some  $\theta$ . ← Bellman operator

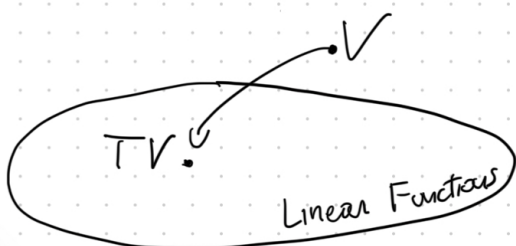
Proof:  $(T_h V)(s,e) \stackrel{\text{def}}{=} r_h(s,e) + \sum_{s'} p(s'|s,e) V(s')$

$$= \phi(s,e)^T \theta_h^R + \sum_{s'} \phi(s,e)^T \Psi(s') V(s')$$

$$= \phi(s,e)^T \left[ \theta_h^R + \underbrace{\sum_{s'} \Psi(s') V(s')}_{\text{def}} \right]$$

$$= \phi(s,e)^T (\theta_h^R + \tilde{\theta}(V))$$

$$= \phi(s,e)^T (\theta_h^R + \tilde{\theta}(V))$$



# LEARNING EFFICIENTLY ON LINEAR MDPs

ALGORITHM: LEAST-SQUARE VALUE ITERATION WITH UCB

for episode  $K=1, \dots, K$

$$D = \left\{ (s_{tK}, e_{tK}, r_{tK}, s'_{t+1,K}) \right\}_{t=1, \dots, H}^{K=1, \dots, K}$$

↖ Dataset

for timestep  $t=1, \dots, H$

covariance  $\rightarrow \Sigma_t = \sum_{i=1}^{K-1} \phi_{ti} \phi_{ti}^T + \lambda I$  (covariance)

$$\hat{\theta}_t = \Sigma_t^{-1} \sum_{i=1}^{K-1} \phi_{ti} [r(s_{ti}, e_{ti}) + \max_{a'} \bar{Q}_{tH}(s_{t+1,i}, e')] ]$$

$$\bar{Q}_t(\cdot, \cdot) = \min \left\{ H, \phi(\cdot, \cdot)^T \hat{\theta}_t + \sqrt{\beta} \|\phi(\cdot, \cdot)\|_{\Sigma_t^{-1}} \right\}$$

bonus  $b_t(\cdot, \cdot)$

end for

Deploy  $\pi$  and add trajectories to the dataset

FIRST STEP ANALYSIS  $\approx$  Total Error = estimation error + bonus + next-state error  
at timestep  $t$

$$(\bar{Q}_{tK} - Q_t^\pi)(s, e) = \phi(s, e)^T \hat{\theta}_{tK} + b_t(s, e) - r(s, e) - \mathbb{E}_{s' \sim p(s, e)} V_{tH}^\pi(s')$$

$$\Sigma_{tK}^{-1} \left[ \sum_{i=1}^{K-1} \phi_{ti} \left\{ r_{ti} + \max_{a'} \bar{Q}_{tH}(s_{t+1,i}, e') \right\} \right]$$

(Transition noise)

$$= \mathbb{E}_{s' \sim p(s_{ti}, e_{ti})} \max_{a'} \bar{Q}_{tH}(s', e') + \gamma r_{ti}$$

Linear MDP Property

$$= \left\{ \phi_{ti}^T \hat{\theta}_t + \phi_{ti}^T \tilde{\theta}(\bar{V}_{tH,i}) + \gamma \right\} \stackrel{\text{def } e_t(s, e)}{\text{estimation error}}$$

$$= \underbrace{\phi(s, e)^T \Sigma_{tH}^{-1} \sum_{i=1}^{K-1} \left[ \theta_t^e + \tilde{\theta}(\bar{V}_{tH,i}) \right]}_{\text{estimation error}} + \underbrace{\left[ \phi(s, e)^T \Sigma_{tH}^{-1} \sum_{i=1}^{K-1} \phi_{ti} \gamma r_{ti} \right]}_{\text{next-state error}} - r(s, e) - \mathbb{E}_{s' \sim p(s, e)} V_{tH}^\pi(s') + b_t(s, e)$$

$$= r(s, e) + \mathbb{E}_{s' \sim p(s, e)} \bar{V}_{tH,i}(s') - r(s, e) - \mathbb{E}_{s' \sim p(s, e)} V_{tH}^\pi(s') + e_t(s, e) + b_t(s, e)$$

$$= \mathbb{E}_{s' \sim p(s, e)} [ \bar{V}_{tH,i} - V_{tH}^\pi ](s') + b_t(s, e) + e_t(s, e)$$

BOUND ON TRANSITION ERROR AND OPTIMISM

LEMMA ↻

wp  $\geq 1-\delta$ :  $|e_t(s, \varrho)| = \varphi(s, \varrho)^T \sum_{i=1}^t \varphi_{t,i} \gamma_{t,i} \leq \|\varphi(s, \varrho)\|_{\Sigma_{tK}^{-1}} \underbrace{\left\| \sum_{i=1}^{t-1} \varphi_{t,i} \gamma_{t,i} \right\|_{\Sigma_{tK}^{-1}}}_{\sqrt{B} = \tilde{O}(dH)}$

$$(\bar{Q}_{tK} - Q^{\pi^*})(s, \pi^*(s)) = \underbrace{E_{s' \sim p(s, \pi^*(s))} [\bar{V}_{t+1, K} - V_{t+1}^{\pi^*}]}_{\geq 0 \text{ by induction}}(s') + \underbrace{b_t(s, \varrho) - e_t(s, \varrho)}_{\text{(this is how the bonus is chosen)}} \geq 0$$

REGRET BOUND

$$\begin{aligned} \sum_{k=1}^K [V^* - V^{\pi_k}] &\leq \sum_{k=1}^K [\bar{V}_k - V^{\pi_k}] = \sum_k \sum_{t=1}^H \underbrace{[b_{tk}(s_{tk}, a_{tk})]}_{\substack{\text{experienced state-actions} \\ \sqrt{B} \|\varphi_{tk}\|_{\Sigma_{tk}^{-1}}}} + \underbrace{e_{tk}(s_{tk}, a_{tk})}_{\leftarrow \text{same order}} \\ &\leq \sum_{t=1}^H \sqrt{K} \sqrt{B} \sqrt{\sum_k \|\varphi_{tk}\|_{\Sigma_{tk}^{-1}}} \\ &\sim H \underbrace{\sqrt{K}}_{\tilde{O}(dH)} \underbrace{\sqrt{B}}_{\tilde{O}(d)} \sqrt{d} \sim H \frac{d}{\sqrt{K}} \sqrt{d} \sqrt{K} \text{ # steps} \end{aligned}$$