

# Nearly Minimax Optimal Reinforcement Learning for Linear Mixture Markov Decision Processes

Dongruo Zhou, Quanquan Gu, Csaba Szepesvári

Presenter: Hao Qin

## Outline

### **Introduction to the linear mixture model setting**

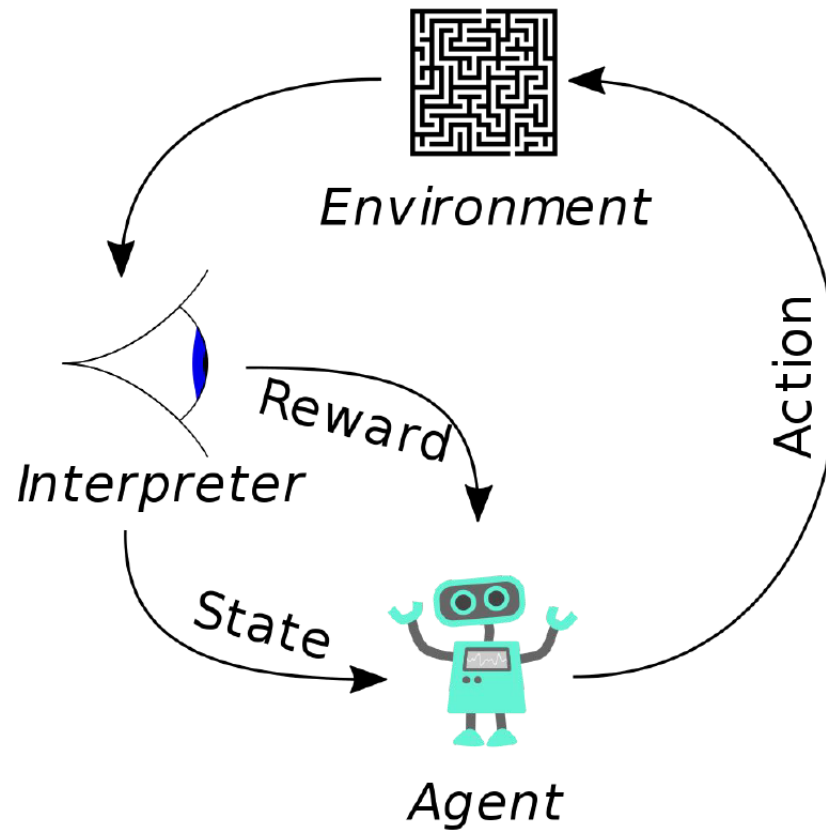
Value target regression and UCRL-VTR algorithm

Weighted regression model

Final regret analysis

Summary

# Online Reinforcement Learning



# Online Reinforcement Learning

Application: game, autonomous driving, dialogue system (Siri) ...



# Online learning Time-inhomogeneous Episodic MDPs

$$M( \underbrace{s \in \mathcal{S}}_{\text{state space}}, \underbrace{a \in \mathcal{A}}_{\text{action space}}, \underbrace{H}_{\text{episode length}}, \underbrace{(r_h(s, a))_h}_{\text{reward functions}}, \underbrace{(\mathbb{P}_h(s'|s, a))_h}_{\text{transition dynamics}} )$$

For  $k$  from 1 to  $K$ , starting from the initial state  $s_0^k$ ,

Deterministic policy  $\pi^k = (\pi_h^k)_{h=1}^H$

For  $h$  from 1 to  $H$ ,

Select action  $a_h^k \leftarrow \pi_h^k(s_h^k)$

Observe returned  $r_h^k$  reward and next-state  $s_{h+1}^k$

# Online learning Time-inhomogeneous Episodic MDPs

## Objective

Minimize the regret

$$\text{Regret}(M, K) = \sum_{k=1}^K [V_1^*(s_1^k) - V_1^{\pi^k}(s_1^k)]$$

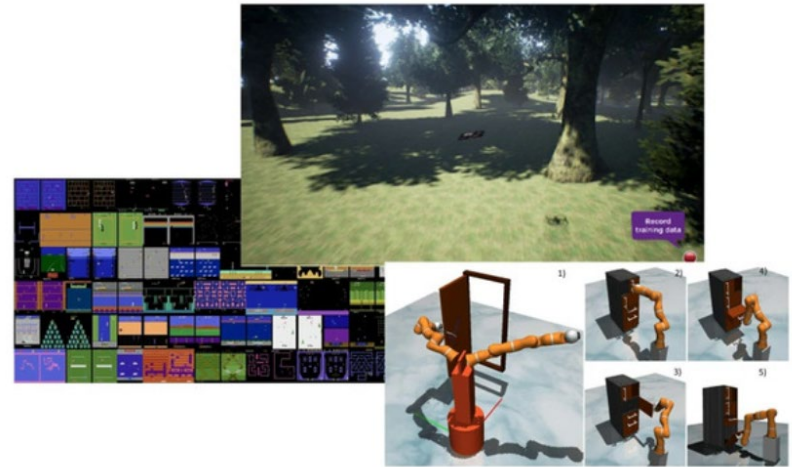
Where  $V_h^\pi(s)$  and  $V_h^*(s)$  have been defined as

$$V_h^\pi(s) = Q_h^\pi(s, \pi_h(s)), \quad Q_h^\pi(s, a) = \mathbb{E}_{\pi, h, s, a} \left[ \sum_{h'=h}^H r_h(s_{h'}, a_{h'}) \right], \quad V_{H+1}^\pi(s) = 0,$$

$$V_h^*(s) = \sup_{\pi} V_h^\pi(s), \quad Q_h^*(s, a) = \sup_{\pi} Q_h^\pi(s, a)$$

# RL with Linear Function Approximation

In reality, the feature space ( $space \times action$ ) can be quite large



Linear mixture MDPs

$$P_h(s'|s, a) = \sum_{j=1}^d \theta_{h,j} \phi_j(s'|s, a)$$

# RL with Linear Function Approximation

Linear mixture MDPs model assumption

1. Linear Mixture model

$$P_h(s'|s, a) = \sum_{j=1}^d \theta_{h,j} \phi_j(s'|s, a)$$

2. Reward function  $r_h(s, a)$  is known.

3. Known transition model family  $\mathcal{P}$ , where  $\phi_1, \phi_2, \dots, \phi_d \in \Phi$



# RL with Linear Function Approximation

## Linear mixture MDPs model assumption

### 1. Linear Mixture model

$$\mathbb{P}_h(s'|s, a) = \langle \phi(s'|s, a), \theta_h^* \rangle \quad \|\theta_h^*\|_2 \leq B$$
$$\phi_V(s, a) = \sum_{s' \in \mathcal{S}} \phi(s'|s, a) V(s'), \quad \forall V: \mathcal{S} \rightarrow [0, 1], (s, a) \in \mathcal{S} \times \mathcal{A}, \|\phi_V(s, a)\|_2 \leq 1$$

### 2. Rewa

$$[\mathbb{P}_h V_{k, h+1}](s_h^k, a_h^k) = \left\langle \sum_{s' \in \mathcal{S}} \phi(s'|s_h^k, a_h^k) V_{k, h+1}(s'), \theta_h^* \right\rangle = \langle \phi_{V_{k, h+1}}(s'|s, a), \theta_h^* \rangle$$

### 3. Known transition model family $\Phi$ , where $\phi_1, \phi_2, \dots, \phi_d \in \Phi$

Example:

$$\text{Tabular MDPs } d = S^2 A, \phi(s'|s, a) = e_{s, a, s'} \in \mathbb{R}^d, \theta_h^* = [\mathbb{P}_h(s'|s, a)]_{s, a, s'} \in \mathbb{R}^d$$

## Outline

Introduction to the linear mixture model setting

**Value target regression and UCRL-VTR algorithm**

Weighted regression model

Final regret analysis

Summary

# Upper Bound: From Linear Bandits to Linear Mixture MDPs

- For linear mixture MDPs, only need to estimate the unknown  $\theta_h^*$
- Key observation For function  $V$ ,  $\mathbb{E}_{s' \sim \mathbb{P}_h(\cdot|s, a)} V(s') = \langle \theta_h^*, \phi_V(s, a) \rangle$

Value-targeted regression (VTR):

With enough number of pairs  $(s_h^k, a_h^k, s_{h+1}^k, V_{k,h+1})$ ,  $s_{h+1}^k \sim \mathbb{P}_h(\cdot|s_h^k, a_h^k)$  estimating  $\theta_h^*$  is possible by doing regression over  $(\phi_{V_{k,h+1}}(s_h^k, a_h^k), V_{k,h+1}(s_{h+1}^k))$

$$[\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) = \langle \phi_{V_{k,h+1}}(s' | s, a), \theta_h^* \rangle$$

- Example: UCRL-VTR (Jia et al., 2020) : Value-target regression + OFUL
- Theorem (Jia et al., 2020): UCRL-VTR enjoys regret

$$\text{Regret}(M, K) = \tilde{O}(dH^2\sqrt{K})$$

# RL with Linear Function Approximation

Notation

$$\phi_{V_{k,h+1}}(s_h^k, a_h^k) = \phi_{k,h+1}$$

$$\phi_{V_{k,h+1}^2}(s_h^k, a_h^k) = \chi_{k,h+1}$$

For example

$$\mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s_h^k, a_h^k)} V_{k,h+1}(s') = \langle \theta_h, \phi_{V_{k,h+1}}(s_h^k, a_h^k) \rangle = \langle \theta_h, \phi_{k,h+1} \rangle$$

# UCRL-VTR algorithm

At  $k$ th episode, the agent

- Receive the initial state. At  $h$ th stage,
- Solution to ridge regression over pairs:

$$\left( \phi_{k,h+1}, V_{k,h+1}(s_{h+1}^k) \right)_k \rightarrow \hat{\theta}_{k,h}$$

$$\hat{\theta}_{k,h} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \lambda \|\theta\|_2^2 + \sum_{j=1}^{k-1} \left[ \langle \phi_{j,h+1}, \theta \rangle - V_{j,h+1}(s_{h+1}^j) \right]^2$$

- Construct confidence sets

$$\hat{\mathcal{C}}_{k,h} = \left\{ \theta : \left\| \hat{\Sigma}_{k,h}^{-\frac{1}{2}} (\theta - \hat{\theta}_{k,h}) \right\|_2 \leq \hat{\beta}_k \right\}, \hat{\Sigma}_{k,h} = \lambda \mathbf{I} + \sum_{j=1}^{k-1} \phi_{j,h+1} \phi_{j,h+1}^T$$

- Estimate previous value function by Bellman equation

$$Q_{k,h}(\cdot, \cdot) = \left[ r_h(\cdot, \cdot) + \max_{\theta \in \hat{\mathcal{C}}_{k,h}} \left\langle \theta, \phi_{V_{k,h+1}}(\cdot, \cdot) \right\rangle \right]_{[0,H]}, V_{k,h}(\cdot) = \max_{a \in \mathcal{A}} Q_{k,h}(\cdot, a)$$

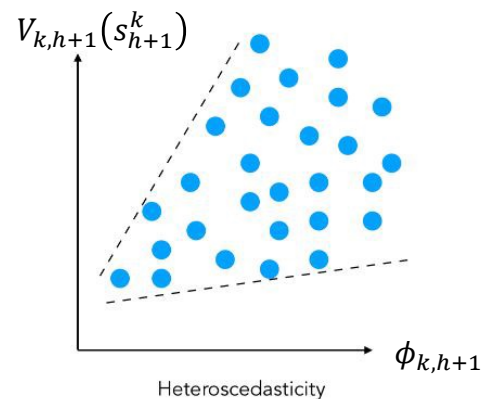
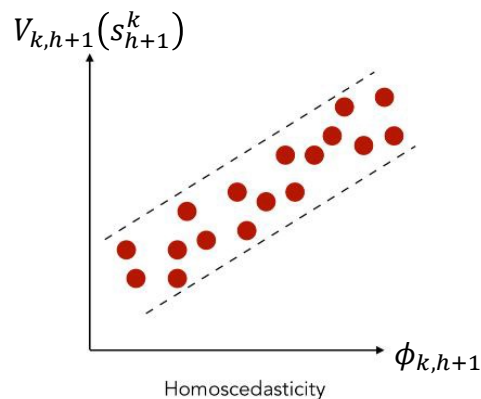
- Select action  $a_h^k = \arg \max_{a \in \mathcal{A}} Q_{k,h}(s_h^k, a), h = 1, \dots$

# UCRL-VTR algorithm

## Shortcomings of UCRL-VTR

$$\langle \theta_h^*, \phi_{k,h+1} \rangle + \varepsilon_k = V_{k,h+1}(s_{h+1}^k)$$

- Choose  $\beta_k$  (the radius of confidence set  $\hat{\mathcal{C}}_{k,h}$ ) proportional to the magnitude of the value function  $V_{k,h+1}(\cdot)$  rather than its variance  $[\mathbb{V}_h V_{k,h+1}](\cdot; \cdot)$
- In the heteroscedastic model, value functions  $[\mathbb{V}_h V_{k,h+1}](\cdot; \cdot)$  cannot be bounded uniformly at different stage, hence we need to make some adjustments such as using weighted least-square estimator.



# Bernstein-Type Self Normalized Inequality

- **Theorem 2 (Bernstein inequality for vector-valued martingales)**

If there is stochastic process  $\{x_t, \eta_t\}_{t \geq 1}$  and a linear relationship that  $y_t = \langle \mu^*, x_t \rangle + \eta_t$ . Also,  $\eta_t, x_t$  satisfy that

$$|\eta_t| \leq R, \quad \mathbb{E}[\eta_t | \mathcal{G}_t] = 0, \quad \|x_t\|_2 \leq X \quad \text{and} \quad \mathbb{E}[\eta_t^2 | \mathcal{G}_t] \leq \sigma^2$$

We have

$$\forall T, \quad \|\mu^* - \hat{\mu}_t\|_{A_t} \leq \hat{\beta}_t + \sqrt{\lambda} \|\mu^*\|_2, \quad \hat{\beta}_t = \tilde{O}(\sigma\sqrt{d} + R)$$

Where  $A_t = \lambda I + \sum_{i=1}^t x_i x_i^T$ ,  $\mu^t = A_t^{-1} (\sum_{i=1}^t y_i x_i)$  and  $\lambda > 0$ .

$$\hat{\beta}_t = \tilde{O}(R\sqrt{d}) \rightarrow \tilde{O}(\sigma\sqrt{d} + R)$$

- Can be extended to sub-exponential random variable case
- Strict improvement from Abbasi-Yadkori et al. (2011):  $\hat{\beta}_t = \tilde{O}(R\sqrt{d})$
- Following induction proof by Dani et al. (2008)

## Outline

Introduction to the linear mixture model setting

Value target regression and UCRL-VTR algorithm

## **Weighted regression model**

Final regret analysis

Summary



# Warmup: Heteroscedastic Linear Bandits

For  $t$  from 1 to  $T$ ,

Receive the decision set:  $\mathcal{D}_t$

Select an action  $\alpha_t \in \mathcal{D}_t$  based on past observation

Observe reward  $r_t$  and **variance**  $\sigma_t$



$$r_t = \langle \boldsymbol{\mu}^*, \mathbf{a}_t \rangle + \epsilon_t, \quad |\epsilon_t| \leq R, \quad \mathbb{E}[\epsilon_t | \mathbf{a}_{1:t}, \epsilon_{1:t-1}] = 0, \quad \mathbb{E}[\epsilon_t^2 | \mathbf{a}_{1:t}, \epsilon_{1:t-1}] \leq \sigma_t^2$$

Objective: Minimize the regret

$$R(T) = \sum_{t=1}^T \sup_{\mathbf{a} \in \mathcal{D}_t} \langle \mathbf{a}, \boldsymbol{\mu}^* \rangle - \langle \mathbf{a}_t, \boldsymbol{\mu}^* \rangle$$

# OFUL and weighted OFUL methods

OFUL (Abbasi-Yadkori et al., 2011): ‘Optimistic’ estimation of reward  
Compute the ridge regression over predictor-response pairs:



$$\hat{\mu}_t \leftarrow \arg \min_{\mu \in \mathbb{R}^d} \lambda \|\mu\|_2^2 + \sum_{i=1}^t [\langle \mu, \alpha_i \rangle - r_i]^2$$

Weighted OFUL:

Solution to ridge regression over **weighted** predictor-response pairs  $\left(\frac{\alpha_t}{\bar{\sigma}_t}, \frac{r_t}{\bar{\sigma}_t}\right) \rightarrow \hat{\mu}_t$

$$\hat{\mu}_t \leftarrow \arg \min_{\mu \in \mathbb{R}^d} \lambda \|\mu\|_2^2 + \sum_{i=1}^t \frac{[\langle \mu, \alpha_i \rangle - r_i]^2}{\bar{\sigma}_i^2}$$

$\bar{\sigma}_t = \max\{\alpha, \sigma_t\} \neq 0$

# OFUL and weighted OFUL methods

OFUL (Abbasi-Yadkori et al., 2011):

**Theorem** : *the regret of OFUL is*

$$\tilde{O}(dR\sqrt{T})$$

Weighted OFUL:

Theorem 3 *Weighted OFUL enjoys the regret*

$$R(T) = \tilde{O}\left(R\sqrt{dT} + d\sqrt{\sum_{t=1}^T \sigma_t^2}\right)$$

'Nearly' independent of T when  $\sum_{t=1}^T \sigma_t^2 \geq \frac{TR^2}{d}$

Strictly improves  $\tilde{O}(dR\sqrt{T})$ .



# UCRL-VTR+ algorithm

At  $k$ th episode, the agent

- Receive the initial state. At  $h$ th stage,
- Solution to ridge regression over weighted pairs:

$$\left( \frac{\phi_{k,h+1}}{\bar{\sigma}_{k,h}}, \frac{V_{k,h+1}}{\bar{\sigma}_{k,h}} \right)_k \rightarrow \hat{\theta}_{k,h}$$

$$\hat{\theta}_{k,h} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \lambda \|\theta\|_2^2 + \sum_{j=1}^{k-1} \frac{[\langle \phi_{j,h+1}, \theta \rangle - V_{j,h+1}(s_{h+1}^j)]^2}{\bar{\sigma}_{k,h}^2}$$

- Construct confidence sets  $\hat{\mathcal{C}}_{k,h}$  with diameter  $\hat{\beta}_k$
- Estimate previous value function by Bellman equation
 
$$Q_{k,h}(\cdot, \cdot) = \left[ r_h(\cdot, \cdot) + \max_{\theta \in \hat{\mathcal{C}}_{k,h}} \left\langle \theta, \phi_{V_{k,h+1}}(\cdot, \cdot) \right\rangle \right]_{[0,H]}, \quad V_{k,h}(\cdot) = \max_{a \in \mathcal{A}} Q_{k,h}(\cdot, a)$$
- Select action  $a_h^k = \arg \max_{a \in \mathcal{A}} Q_{k,h}(s_h^k, a), h = 1, \dots$

# UCRL-VTR+ algorithm

## Three major problems

- Q1: How to calculate the empirical variance  $[\widehat{V}_{k,h}V_{k,h+1}](s_h^k, a_h^k)$ ?

The variance  $[V_{k,h}V_{k,h+1}](s_h^k, a_h^k)$  has been defined as

$$[V_{k,h}V_{k,h+1}](s_h^k, a_h^k) := \mathbb{E}_{s' \sim P(\cdot | s_h^k, a_h^k)} [V_{k,h+1}(s') - (\mathbb{P}V_{k,h+1})(s_h^k, a_h^k)]^2$$

- Q2 How to select  $E_{k,h}$  to guarantee that  $[\widehat{V}_{k,h}V_{k,h+1}](s_h^k, a_h^k) + E_{k,h}$  upper bound  $\sigma_{k,h}^*$  w.h.p.?
- Q3 : How to choose an appropriate  $\widehat{\beta}_k$ , such that  $\widehat{C}_{k,h}$  contains  $\theta_h^*$  w.h.p.?

# UCRL-VTR+ algorithm

## A1: Variance Estimator

- Unlike bandit setting, we need to estimate the variance

$$\begin{aligned} [\mathbb{V}V](s, a) &= \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a)} V^2(s') - \left[ \mathbb{E}_{s' \sim \mathbb{P}_h(\cdot | s, a)} V(s') \right]^2 \\ &= \underbrace{\langle \theta_h^*, \chi \rangle}_{\langle \tilde{\theta}_h, \chi \rangle} - \underbrace{[\langle \theta_h^*, \phi \rangle]^2}_{\langle \hat{\theta}_h, \phi \rangle} \end{aligned}$$

- Regression over predictor-response  $(\chi_{k,h+1}, V_{k,h+1}^2(s_{h+1}^k))$

$$\tilde{\theta}_{k,h} = \arg \min_{\theta \in \mathbb{R}^d} \lambda \|\theta\|_2^2 + \sum_{j=1}^{k-1} \left[ \langle \chi_{j,h+1}, \theta \rangle - V_{j,h+1}^2(s_{h+1}^j) \right]^2$$

- Variance estimator with clip

$$\bar{\mathbb{V}}_{k,h} V_{k,h+1}(s_h^k, a_h^k) = [\langle \chi_{k,h+1}, \tilde{\theta}_{k,h+1} \rangle]_{[0, H^2]} - [[\langle \phi_{k,h+1}, \hat{\theta}_{k,h+1} \rangle]^2]_{[0, H^2]}$$

# UCRL-VTR+ algorithm

## A2: Bonus for variance estimation

- Bonus  $E_{k,h}$  has been defined as follows:

$$E_{k,h} = \min \left\{ H^2, 2H\check{\beta}_k \left\| \widehat{\Sigma}_{k,h}^{-\frac{1}{2}} \phi_{k,h+1} \right\|_2 \right\} + \min \left\{ H^2, \tilde{\beta}_k \left\| \widehat{\Sigma}_{k,h}^{-\frac{1}{2}} \chi_{k,h+1} \right\|_2 \right\}$$
$$\hat{\beta}_k = \tilde{O}(\sqrt{d}), \check{\beta}_k = \tilde{O}(d), \tilde{\beta}_k = \tilde{O}(H^2\sqrt{d})$$

Final variance upper bound:  $\bar{\sigma}_{k,h} \leftarrow \sqrt{\max\left\{\frac{H^2}{d}, \bar{\mathbb{V}}_{k,h}V_{k,h+1} + E_{k,h}\right\}}$

Then w.h.p., for all  $\mathbf{k} \in [\mathbf{K}]$  and  $\mathbf{h} \in [\mathbf{H}]$

$$|\bar{\mathbb{V}}_{k,h}V_{k,h+1} - \mathbb{V}_{k,h}V_{k,h+1}| \leq E_{k,h}, \sum_{k=1}^K \sum_{h=1}^H E_{k,h} = \tilde{O}(d^{\frac{3}{2}}H^3\sqrt{K})$$

# UCRL-VTR+ algorithm

## A3: Confidence set

- Lemma 5

The confidence set has been defined as

$$\hat{\mathcal{C}}_{k,h} = \left\{ \theta: \left\| \hat{\Sigma}_{k,h}^{\frac{1}{2}} (\theta - \hat{\theta}_{k,h}) \right\|_2 \leq \hat{\beta}_k \right\}, \hat{\Sigma}_{k,h} = \lambda \mathbf{I} + \sum_{i=1}^{k-1} \frac{\phi_{j,h+1} \phi_{j,h+1}^T}{\bar{\sigma}_{j,h}^2}$$

Where  $\hat{\beta}_k = \tilde{O}(\sqrt{d})$



# UCRL-VTR+ algorithm

At  $k$ th episode, the agent

- Receive the initial state. At  $h$ th stage,
- Solution to ridge regression over weighted pairs:

$$\left( \frac{\phi_{k,h+1}}{\bar{\sigma}_{k,h}}, \frac{V_{k,h+1}(s_{h+1}^k)}{\bar{\sigma}_{k,h}} \right)_k \rightarrow \hat{\theta}_{k,h}$$

$$\hat{\theta}_{k,h} = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \lambda \|\theta\|_2^2 + \sum_{j=1}^{k-1} \left[ \langle \phi_{j,h+1}, \theta \rangle - V_{j,h+1}(s_{h+1}^j) \right]^2 / \bar{\sigma}_{j,h}^2$$

- Construct confidence sets

$$\hat{\mathcal{C}}_{k,h} = \left\{ \theta : \left\| \hat{\Sigma}_{k,h}^{-\frac{1}{2}} (\theta - \hat{\theta}_{k,h}) \right\|_2 \leq \hat{\beta}_k \right\}, \hat{\Sigma}_{k,h} = \lambda \mathbf{I} + \sum_{j=1}^{k-1} \phi_{j,h+1} \phi_{j,h+1}^T / \bar{\sigma}_{j,h}^2$$

- Estimate previous value function by Bellman equation

$$Q_{k,h}(\cdot, \cdot) = \left[ r_h(\cdot, \cdot) + \max_{\theta \in \hat{\mathcal{C}}_{k,h}} \left\langle \theta, \phi_{V_{k,h+1}}(\cdot, \cdot) \right\rangle \right]_{[0,H]}, V_{k,h}(\cdot) = \max_{a \in \mathcal{A}} Q_{k,h}(\cdot, a)$$

- Select action  $a_h^k = \arg \max_{a \in \mathcal{A}} Q_{k,h}(s_h^k, a), h = 1, \dots$

## Outline

Introduction to the linear mixture model setting

Value target regression and UCRL-VTR algorithm

Weighted regression model

**Final regret analysis**

Summary

# UCRL-VTR+ algorithm

## Upper Bound: Regret Decomposition

$$\begin{aligned}
 & V_{k,h}(s_h^k) - V_h^{\pi^k}(s_h^k) \\
 & \leq \langle \hat{\theta}_{k,h}, \phi_{k,h+1} \rangle + \hat{\beta}_k \left\| \hat{\Sigma}_{k,h}^{-\frac{1}{2}} \phi_{k,h+1} \right\|_2 - [\mathbb{P}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k) \\
 & \leq [\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k) + 2\hat{\beta}_k \left\| \hat{\Sigma}_{k,h}^{-\frac{1}{2}} \phi_{k,h+1} \right\|_2 \\
 & \leq [\mathbb{P}_h V_{k,h+1}](s_h^k, a_h^k) - [\mathbb{P}_h V_{k,h+1}^{\pi^k}](s_h^k, a_h^k) + 2\hat{\beta}_k \bar{\sigma}_{k,h} \left\| \hat{\Sigma}_{k,h}^{-\frac{1}{2}} \frac{\phi_{k,h+1}}{\bar{\sigma}_{k,h}} \right\|_2
 \end{aligned}$$

Summing over  $h \in ]H]$  and  $k \in [K]$

$$\begin{aligned}
 & \sum_{k=1}^K V_{k,1}(s_1^k) - V_1^{\pi^k}(s_1^k) \\
 & \leq \sum_{k=1}^K \sum_{h=1}^H 2\hat{\beta}_k \bar{\sigma}_{k,h} \left\| \hat{\Sigma}_{k,h}^{-\frac{1}{2}} \frac{\phi_{k,h+1}}{\bar{\sigma}_{k,h}} \right\|_2 + \tilde{O}(H^2 \sqrt{dK}) \\
 & \leq 2\hat{\beta}_K \sqrt{\sum_{k=1}^K \sum_{h=1}^H \bar{\sigma}_{k,h}^2} \cdot \sqrt{\sum_{k=1}^K \sum_{h=1}^H \left\| \hat{\Sigma}_{k,h}^{-\frac{1}{2}} \frac{\phi_{k,h+1}}{\bar{\sigma}_{k,h}} \right\|_2^2} + \tilde{O}(H^2 \sqrt{dK})
 \end{aligned}$$

# UCRL-VTR+ algorithm

## Naïve approach:

- Bound variance by  $H^2$

$$\begin{aligned} \text{Regret}(M, K) &= \tilde{O}\left(H^2\sqrt{dK} + d\sqrt{H} \left[ \sum_{k=1}^K \sum_{h=1}^H \widehat{V}_{h+1}^{\pi^k}(s_h^k, a_h^k) \right]^{\frac{1}{2}}\right) \\ &\leq \tilde{O}\left(H^2\sqrt{dK} + d\sqrt{H} \cdot (HK \cdot H^2)^{\frac{1}{2}}\right) \end{aligned}$$

## Law of total variance

- Weighted OFUL suggests a regret of UCRL-VTR+

$$\begin{aligned} \text{Regret}(M, K) &= \tilde{O}\left(H^2\sqrt{dK} + d\sqrt{H} \left[ \sum_{k=1}^K \sum_{h=1}^H \widehat{V}_{h+1}^{\pi^k}(s_h^k, a_h^k) \right]^{\frac{1}{2}}\right) \\ &\leq \tilde{O}\left(H^2\sqrt{dK} + d\sqrt{H} \cdot (H^2K)^{\frac{1}{2}}\right) \\ &\leq \tilde{O}(\sqrt{d^2H^3 + dH^4\sqrt{K}}) \end{aligned}$$

- Use Law of total variance (Lattimore and Hutter, 2012; Azar et al., 2017)

$$\sum_{k=1}^K \sum_{h=1}^H \mathbb{V}V_{h+1}^{\pi^k}(s') = \tilde{O}(H^2K)$$

# UCRL-VTR+ algorithm

## Regret of UCRL-VTR+

**Theorem 6 (Upper bound)** *w.h.p., UCRL-VTR+ attains regret*

$$\text{Regret}(M, K) = \tilde{O}\left(\sqrt{d^2 H^3 + d H^4 \sqrt{K}} + d^2 H^3 + d^3 H^2\right)$$

**Theorem 8 (Lower bound)** *Let  $B > 1$ , for any RL algorithm, there exists a time-inhomogeneous episodic  $B$ -linear mixture MDP  $M$ , where*

$$\mathbb{E}[\text{Regret}(M, K)] \geq \Omega(d H^{3/2} \sqrt{K})$$

The regret of UCRL-VTR+ matches the lower bound when  $d \geq H$ ,  $K \geq d^4 H + d^3 H^2$

## Outline

Introduction to the linear mixture model setting

Value target regression and UCRL-VTR algorithm

Weighted regression model

Final regret analysis

**Summary**

# Summary

- Nearly minimax optimal RL utilize a variance-dependent weighted linear regression
- Optimistic estimation of value function
- Combining law of total variance to reach a lower bound



Thank you