# RL with continuous action spaces
## General function approximation

Aarti Singh

Machine Learning 10-734
Nov 6, 2025

Slides courtesy: Wen Sun

# Continuous action spaces

Bandits:

Reward is linear, Lipshitz, GP, NN, …

e.g. $r^*(x) = x^T\theta^*$                         x, $\theta^*$ are d-dimensional

MDP:

Linear MDP - Reward is linear, Transition is low rank

$$r_h(s,a) = w_h^T \phi(s,a), P_h(s'|s,a) = \mu_h(s')^T \phi(s,a)$$

LSVI-UCB algorithm has low regret $\tilde{O}(H^2\sqrt{d^3 N})$

Linear Q* - $Q^*(s,a) = \theta^{*\,T}\phi(s,a)$
                Doesn't work!

# LSVI-UCB: Least Square Value Iteration with UCB

Value iteration at episode n using $\{s_h^i, a_h^i, r_h^i, s_{h+1}^i\}_{h=1,i=1}^{H-1,n-1}$

$$\widehat{V}_H^n(s) = 0, \forall s$$

For h = H-1, H-2, ..., 1

*least square*

*Bellman consistency*

$T\theta$

$\vec{Q}(s,a)$

$$\theta_h^n \leftarrow \underset{\theta}{\arg\min} \sum_{i=1}^{n-1} \left( \langle \phi(s_h^i, a_h^i), \theta \rangle - r_h^i - \widehat{V}_{h+1}^n(s_{h+1}^i) \right)^2 + \lambda \|\theta\|_2^2$$

$$\widehat{Q}_h^n(s,a) = \min \left\{ b_h^n(s,a) + \langle \phi(s,a), \theta_h^n \rangle, \quad H \right\}, \forall s, a$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s,a), \qquad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s,a), \forall s$$

# Bellman error

Consider $f(s, a) = Q(s, a)$.

Bellman error $= f(s, a) - T f(s, a)$    $\equiv Q - TQ$

$$= f(s, a) - \left( r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} f(s', a') \right)$$

annotations: $\phi(s,a)$, $\phi(s,a)$, $\phi(s,a)$

$f$ linear
$Tf$ linear

$Tf = \sum_{s'} P(s' | s, a) \max_{a'} f(s', a')$

If Bellman error $\neq 0$, then $f \neq Q^*$

Why does linear $Q^*$ not suffice?

> Even if $f$ is linear, $T f$ may not be linear $T$ unless transitions $P$ (and reward $r$) is also linear!

# Bellman completeness

**Bellman completeness**: For any Q function in $\mathcal{F}$, its Bellman update is also in $\mathcal{F}$

Implies Bellman error

can be 0 for $f = Q^*$

# LSVI-UCB: Least Square Value Iteration with UCB

Value iteration at episode n using $\{s_h^i, a_h^i, r_h^i, s_{h+1}^i\}_{h=1,i=1}^{H-1,n-1}$

$$\widehat{V}_H^n(s) = 0, \forall s$$

For h = H-1, H-2, …, 1

$Q = T\theta$

$$\theta_h^n \leftarrow \underset{\theta}{\arg\min} \sum_{i=1}^{n-1} \left( \langle \phi(s_h^i, a_h^i), \theta \rangle - r_h^i - \widehat{V}_{h+1}^n(s_{h+1}^i) \right)^2 + \lambda \|\theta\|_2^2$$

$$\widehat{Q}_h^n(s,a) = \min \left\{ b_h^n(s,a) + \langle \phi(s,a), \theta_h^n \rangle, \quad H \right\}, \forall s,a$$

$\dfrac{1}{\sqrt{N_h^n(s,a)}}$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s,a), \qquad \pi_h^n(s) = \arg\max_a \widehat{Q}_h^n(s,a), \forall s$$

# LSVI-UCB does not work under Bellman completeness

$$\widehat{Q}_h^n(s,a) = \min\left\{ b_h^n(s,a) + \langle \phi(s,a), \theta_h^n \rangle, \quad H \right\}, \forall s, a$$

Issue: Adding bonus which may be non-realizable
  (e.g. in linear case, bonus may be nonlinear in s)

Recall $\quad b_h^n(s,a) = \|\phi\|_{\Lambda_h^{n-1}} \beta \longrightarrow$ nonlinear in $\phi(s,a)$

Need different algorithm (no bonus on Q)
  – how to achieve optimism?

# Average Bellman error

Weaker notion of Bellman error:

Evaluate $g$-approximation of $Q$ using a policy $\pi_f$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ g(s_h, a_h) - r(s_h, a_h) - \underbrace{\mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} g(s_{h+1}, a) \right]}_{\text{Bellman error } Q - TQ} \right]$$

$f$: defines roll-in distribution over $s_h, a_h$

We know that $\mathcal{E}(Q^\star; f, h) = 0, \forall f$   if $g = Q^\star$   $g = Tg$   pointwise $(s, a)$

Hence, any $g$ such that $\mathcal{E}(g; f, h) \neq 0$, is an incorrect $Q^\star$ approximator

# Average Bellman error

Evaluate average Bellman wrt *V* function induced by $g$ as well:

$$V - TV$$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h \sim d_h^{\pi_f}} \left[ V_g(s_h) - r(s_h, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P(\cdot | s_h, \pi_g(s_h))} \left[ V_g(s_{h+1}) \right] \right]$$

$f$: defines roll-in distribution over $s_h, a_h$

Again we have $\mathcal{E}(Q^\star; f, h) = 0, \forall f$

$$\text{if } g = Q^\star \quad V_g = V^\star$$

( because: $V_{Q^\star}(s) - r(s, \pi_{Q^\star}(s)) - \mathbb{E}_{s' \sim P_h(\cdot | s, \pi_{Q^\star}(s))} V_{Q^\star}(s') = 0$ )

Hence, any $g$ such that $\mathcal{E}(g; f, h) \neq 0$, is an incorrect $Q^\star$ approximator

# Bellman rank

∃ two mappings $W_h : \mathcal{F} \mapsto \mathbb{R}^d$, $X_h : \mathcal{F} \mapsto \mathbb{R}^d$ ( d = Bellman-rank)

s.t. $\forall f, g \in \mathcal{F} : \mathcal{E}(g; f, h) = \langle W_h(g), X_h(f) \rangle$

$\forall h : \mathcal{E}_h \in \mathbb{R}^{|\mathcal{F}| \times |\mathcal{F}|}$

|  | g | f | approx | | |
|---|---|---|---|---|---|
| | | | | | |
| $\pi_f$ | $\mathcal{E}_{g;f,h}$ | $\mathcal{E}_{f,f,h}$ | | | |
| arg | | | | | |
| | | | | | |
| | | | | | |

Rank of this matrix = Bellman rank

Note: we just assume the existence of W, X, but they are unknown

# Examples of Bellman rank

- **Linear Bellman completeness**: For any linear Q function, its Bellman update is also linear

Given feature $\phi$, take any linear function $\theta^\top \phi(s, a)$:

$$\forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s,a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

$$\theta \leftarrow T\theta$$

**Claim: it has Q-Bellman rank d**

$\forall g(s, a) := \theta^\top \phi(s, a)$, we have:

$$\arg(\theta - T\theta)$$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \max_{a \in \mathcal{A}} \theta^\top \phi(s_{h+1}, a) \right] \right]$$

$$= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ \theta^\top \phi(s_h, a_h) - w^T \phi(s_h, a_h) \right] \quad \text{Bellman completeness}$$

$$= \left\langle \theta - w, \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} [\phi(s_h, a_h)] \right\rangle$$

11

# Examples of Bellman rank

- **Linear Bellman completeness**: For any linear Q function, its Bellman update is also linear.

Given feature $\phi$, take any linear function $\theta^\top \phi(s, a)$:

$$\rightarrow \quad \forall h, \exists w \in \mathbb{R}^d, s.t., w^\top \phi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_h(s,a)} \max_{a'} \theta^\top \phi(s', a'), \forall s, a$$

- **Linear MDP** $\quad$ $r$ linear $\,,\,$ $P$ low rank $\quad$ $r = \theta^\top \phi(s, a) \quad P = \mu(s')^\top \phi(s, a)$

$\Rightarrow$ linear Bellman completeness $\Rightarrow$ Q-Bellman rank d

# Examples of Bellman rank

- **Linear Q\* and V\***    $Q^\star(s,a) = (w^\star)^\top \phi(s,a), \quad V^\star(s) = (\theta^\star)^\top \psi(s), \forall s, a$

**Claim: it has Q-Bellman rank 2d**

$$\mathcal{F}_h = \left\{ (w,\theta) : \max_a \underbrace{w^\top \phi(s,a)}_{\text{Q}(s,a)} = \underbrace{\theta^\top \psi(s)}_{\text{V}(s)}, \forall s \right\}$$

$$\mathcal{E}(g; f, h) = \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - r(s_h, a_h) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \theta^\top \psi(s_{h+1}) \right] \right]$$

avg(Q − TQ)

$$= \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \left[ w^\top \phi(s_h, a_h) - \underbrace{(w^\star)^\top \phi(s_h, a_h) + \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ (\theta^\star)^\top \psi(s_{h+1}) \right]}_{Q^\star(s_h, a_h)} \right.$$

$Q^\star(s_h, a_h)$     $V^\star(s_{h+1})$

$$\left. - \mathbb{E}_{s_{h+1} \sim P_h(\cdot | s_h, a_h)} \left[ \theta^\top \psi(s_{h+1}) \right] \right]$$

$$= \left\langle \begin{bmatrix} w - w^\star \\ \theta - \theta^\star \end{bmatrix}, \; \mathbb{E}_{s_h, a_h \sim d_h^{\pi_f}} \begin{bmatrix} \phi(s_h, a_h) \\ -\mathbb{E}_{s' \sim P_h(s_h, a_h)}[\psi(s')] \end{bmatrix} \right\rangle$$

$\psi(s') \neq \eta^\top \phi(s,a)$

13

# Examples of Bellman rank

- **Linear Q* and V***   $Q^{\star}(s, a) = (w^{\star})^{\top}\phi(s, a), \quad V^{\star}(s) = (\theta^{\star})^{\top}\psi(s), \forall s, a$

**Claim: it has Q-Bellman rank 2d**

Note that $\psi(s_h, a_h) := \mathbb{E}_{s' \sim P_h(s_h, a_h)}[\psi(s')]$   is in general not linear in $\phi(s_h, a_h)$ if transition dynamics are not linear

But V* linear *inherently* implies transition dynamics are linear:

Since V* = TV*, we have

$$\theta^{*T}\psi(s) = \max_a \left( r(s, a) + \theta^{*T}\psi(s, a) \right)$$

$$V^* \qquad \mathbb{E}_{s' \sim P(\cdot|s,a)}[\psi(s')]$$

which implies transition dynamics are linear (given definition of $\psi(s, a)$).

**Linear Q*, V* suffices, though linear Q* doesn't!**

14

# Examples of Bellman rank

- **Low rank MDP**

$$P_h(s' \mid s, a) = \mu_h^\star(s')^\top \phi_h^\star(s, a)$$

$\in \mathbb{R}^d$

Linear $- \phi$ known

(neither $\mu^\star$ nor $\phi^\star$ is known)

**Claim: this model has V-Bellman rank $d$**

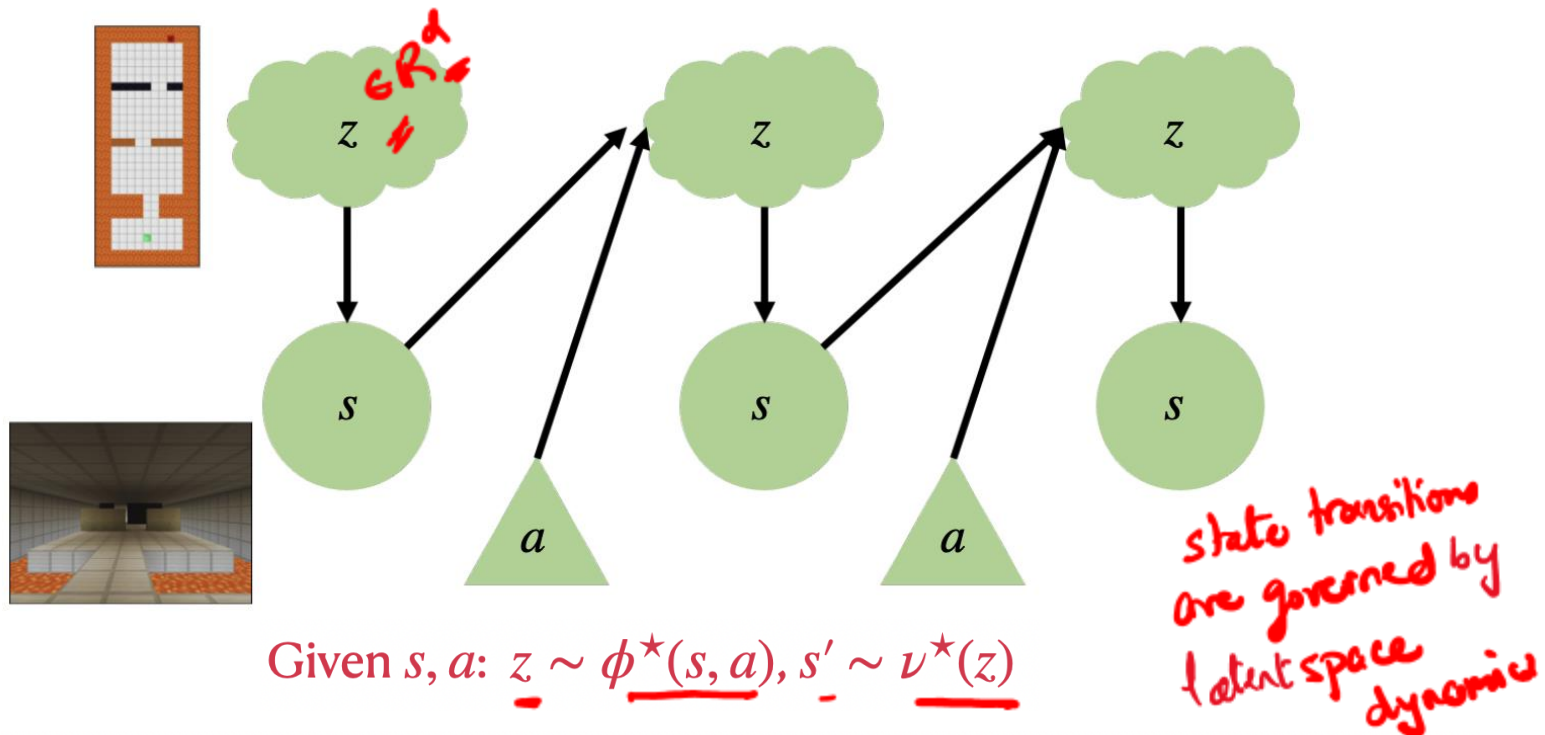$$\mathscr{F}_h = \{\theta^\top \phi(\cdot, \cdot) : \|\theta\|_2 \leq W, \phi \in \Phi\}$$

arg $(V - TV)$

$$\mathbb{E}_{s_h \sim d_h^{\pi_f}}\left[V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot \mid s_h, \pi_g(s_h))}[V_g(s_{h+1})]\right]$$

$$= \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} \mathbb{E}_{s_h \sim P_{h-1}(\cdot \mid \tilde{s}, \tilde{a})}\left[V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot \mid s_h, \pi_g(s_h))}[V_g(s_{h+1})]\right]$$

$$= \mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}} \int_{s_h} \underbrace{\mu_{h-1}^\star(s_h)^\top \phi_{h-1}^\star(\tilde{s}, \tilde{a})}_{P_{h-1}(s_h \mid \tilde{s}, \tilde{a})} \underbrace{\left[V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot \mid s_h, \pi_g(s_h))}[V_g(s_{h+1})]\right]}_{\text{Scalar}} d(s_h)$$

$$= \left\langle \int_{s_h} \mu_{h-1}^\star(s_h)\underbrace{\left[V_g(s_h) - r(s, \pi_g(s_h)) - \mathbb{E}_{s_{h+1} \sim P_h(\cdot \mid s_h, \pi_g(s_h))}[V_g(s_{h+1})]\right]}_{\text{Scalar}} d(s_h), \quad \underbrace{\mathbb{E}_{\tilde{s}, \tilde{a} \sim d_{h-1}^{\pi_f}}[\phi_{h-1}^\star(\tilde{s}, \tilde{a})]}_{} \right\rangle$$

# Examples of Bellman rank

- **Latent variable MDP**        V-Bellman rank = Number of latent states



Given $s, a$: $z \sim \phi^{\star}(s, a)$, $s' \sim \nu^{\star}(z)$

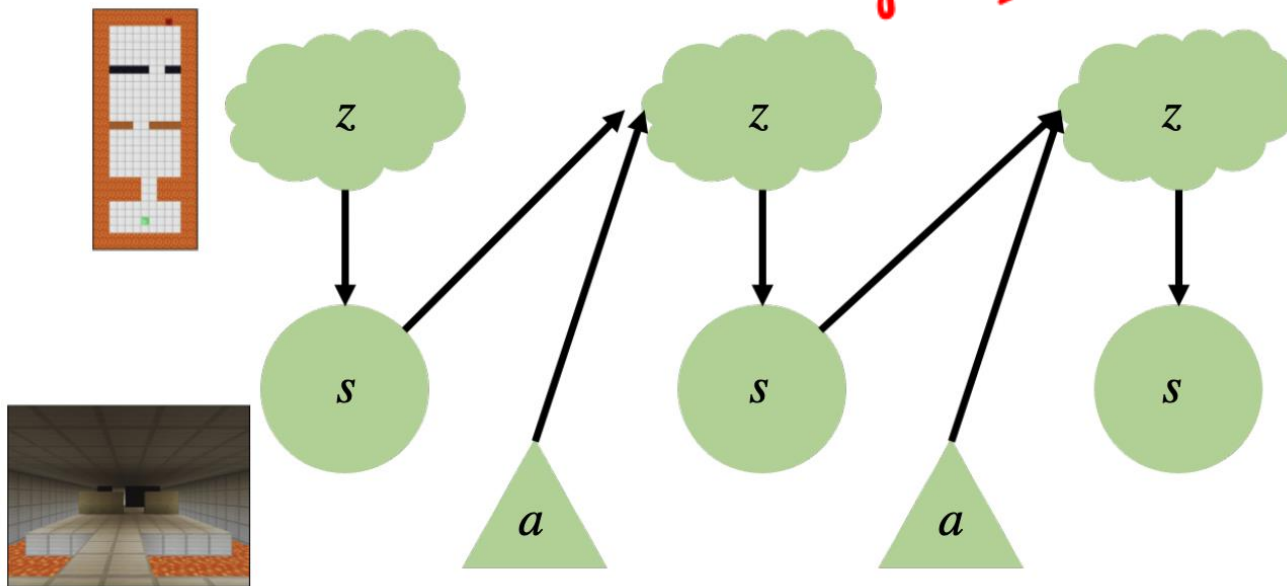state transitions are governed by latent space dynamics

Latent variable MDP is captured by low-rank MDP, so it has small V-Bellman rank…

# Examples of Bellman rank

- **Latent variable MDP**

V-Bellman rank = Number of latent states
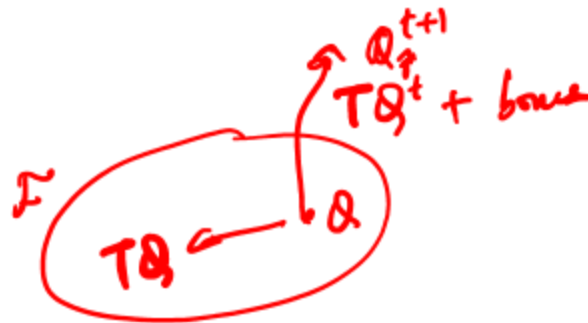
*Dim (z is continuous)*

*(finite z)*



Given $s, a$: $z \sim \phi^\star(s, a)$, $s' \sim \nu^\star(z)$

- **Block MDP -** Special case of latent variable MDP where a state can only be generated from one latent state i.e. one-to-one mapping, hence latent state is deterministically decodable $\quad P(s|z) > 0 \;\Rightarrow\; P(s|z') = 0 \;\; \forall \, z' \neq z$

# Algorithm under Bellman rank

- Works for general function approximation with low Bellman rank

- Gives optimism without adding (nonlinear) bonus

# Bilinear-UCB

A general algorithm under Bellman rank that can learn an $\epsilon$ near optimal policy with number of samples

$$\text{e.g., poly}(H, \text{b-rank}, \ln(|\mathscr{F}|), 1/\epsilon^2)$$

$y \to d \ll |S|, |A|$

# Q-Bellman rank setting

Recall our hypothesis class $\mathscr{F}$, where each $g \in \mathscr{F}$ is in the form of $g(s, a)$

For $Q$-Bellman rank, we define Bellman error loss as:

$$\ell(s_h, a_h, s'_{h+1}, g) = g(s_h, a_h) - r(s_h, a_h) - \max_{a'} g(s_{h+1}, a')$$

*generative*

If we had a dataset $\mathscr{D} := \{s_h, a_h, s_{h+1}\}$ where $s_h, a_h \sim d_h^{\pi_f}, s_{h+1} \sim P_h(\cdot \mid s_h, a_h)$

then $\forall g : \mathbb{E}_{\mathscr{D}}[\ell(s_h, a_h, s_{h+1}, g)]$ is an unbiased est of $\mathscr{E}(g; f, h)$

↓ *empirical average*

# V-Bellman rank setting

Recall our hypothesis class $\mathscr{F}$, where each $g \in \mathscr{F}$ is in the form of $g(s, a)$

For V-Bellman rank, we define Bellman error loss as:

$$V(s) = Q(s, \pi(s))$$

$$\ell(s_h, a_h, s'_{h+1}, g) = \frac{\mathbf{1}\{a_h = \pi_g(s_h)\}}{1/A} \left( g(s_h, a_h) - r(s_h, a_h) - \max_{a'} g(s_{h+1}, a') \right)$$

If we had a dataset $\mathscr{D} := \{s_h, a_h, s_{h+1}\}$ where $s_h \sim d_h^{\pi_f}, a_h \sim U(\mathscr{A}),$

$$s_{h+1} \sim P_h(\cdot \mid s_h, a_h)$$

then $\forall g : \mathbb{E}_{\mathscr{D}}[\ell(s_h, a_h, s_{h+1}, g)]$ is an unbiased est of $\mathscr{E}(g; f, h)$

# Bilinear-UCB

At iteration $t$ :

Select $f_t = \arg\max\limits_{g \in \mathscr{F}} V_g(s_0)$

s.t., $\forall h : \sum\limits_{i=0}^{t-1} \left( \mathbb{E}_{\mathscr{D}_{h,i}}[\ell(s_h, a_h, s_{h+1}, g)] \right)^2 \leq R^2$

*empirical Bellman error*

*Computationally efficient?*

*no bonus required*

*≡ ellipsoid constraint for linear bandits*

# Bilinear-UCB

At iteration $t$ :

Select $f_t = \arg\max\limits_{g \in \mathcal{F}} V_g(s_0)$

s.t., $\forall h : \sum\limits_{i=0}^{t-1} \left( \mathbb{E}_{\mathcal{D}_{h,i}}[\ell(s_h, a_h, s_{h+1}, g)] \right)^2 \leq R^2$

*avg Bellmen error → 0*

*generative setting*

For all h, create $\mathcal{D}_{h,t} = \{s_h, a_h, s_{h+1}\}$ w/ m triples, where:

- For Q-B rank case: $s_h, a_h \sim d_h^{\pi_{f_t}}, s_{h+1} \sim P_h(\cdot \mid s_h, a_h)$

- For V-B rank case: $s_h \sim d_h^{\pi_{f_t}}, a_h \sim U(A), s_{h+1} \sim P_h(\cdot \mid s_h, a_h)$

# Bilinear-UCB

$$\text{Select } f_t = \arg\max_{g \in \mathcal{F}} V_g(s_0) \quad \text{s.t., } \forall h : \sum_{i=0}^{t-1} \left( \mathbb{E}_{\mathscr{D}_{h,i}}[\ell(s_h, a_h, s_{h+1}, g)] \right)^2 \leq R^2$$

UI
bonus

1. When the batch size ($|\mathscr{D}_{h,i}|$) is large,

$$\mathbb{E}_{\mathscr{D}_{h,i}} \ell(s_h, a_h, s_{h+1}, g) \rightarrow \mathscr{E}(g; f_i, h)$$

2. We know that $\sum_{i=1}^{t-1} \mathscr{E}(f^\star; f_i, h) = 0$

ellipsoid

linear bandit : $\theta^\star \in C_t$

3. By properly setting batch size and R, we eliminate wrong hypothesis, but keep $f^\star$

4. This gives optimism: $V_{f_t}(s_0) \geq V_{f^\star}(s_0) := V^\star(s_0)$

Optimism allows explore and exploit tradeoff!

# Analysis of Bilinear-UCB

Uniform convergence style assumption on our hypothesis class $\mathscr{F}$:

Given any distribution $\nu \in \Delta(S \times A \times S)$, and $m$ i.i.d samples $\{s_i, a_i, s_i'\}$ from $\nu$, w/ probability at least $1 - \delta$,

$$\forall g : \left| \mathbb{E}_\nu \ell(s, a, s', g) - \mathbb{E}_{\mathscr{D}} \ell(s, a, s', g) \right| \leq \varepsilon_{gen}(m, \mathscr{F}, \delta)$$

true         empirical

Example: when $\mathscr{F}$ is discrete (for B-rank loss), Hoeffding + union bound over $\mathscr{F}$ implies:

$$\varepsilon_{gen}(m, \mathscr{F}, \delta) := 2H\sqrt{\frac{\ln(|\mathscr{F}|/\delta)}{m}}$$