

Contextual Decision Processes with Low Bellman Rank are PAC-Learnable

Nan Jiang^{1,3}, Akshay Krishnamurthy², Alekh Agarwal³, John Langford³, Robert E. Schapire³
¹University of Michigan, Ann Arbor ²University of Massachusetts, Amherst ³Microsoft Research, NYC

Introduction: 3 challenges of RL

Long-term Planning

Generalization

Contextual Bandits

?

PAC-MDP Theory

Exploration

Our Answer:

- A new measure – **Bellman rank**
 - Captures a wide range of tractable RL problems
- A new algorithm – **OLIVE**
 - Polynomial sample complexity guarantee

Value-based RL in CDPs

Contextual Decision Processes (CDPs): *episodic RL with rich observations*

- Action space A , horizon H .
- Context space X . A context is ...
 - any function of history that expresses a good policy & value function
 - e.g., last 4 frames of images in Atari games
 - e.g., (state, time-step) for finite-horizon tabular MDPs
- An episode: $x_1, a_1, r_1, x_2, \dots, x_H, a_H, r_H$
- Policy $\pi : X \rightarrow A$. Want to maximize $V^\pi = \mathbb{E} \left[\sum_{h=1}^H r_h \mid a_{1:H} \sim \pi \right]$.

In general, $|X|$ is very large \Rightarrow Requires generalization!

Value-based PAC-RL in CDPs

- Input: a function space F which contains Q^*
- Output: π such that, w.p. $\geq 1-\delta$, $V^{\pi^*} - V^\pi \leq \epsilon$ after acquiring $\text{poly}(|A|, H, \log|F|, 1/\epsilon, 1/\delta)$ trajectories.

Need additional condition, otherwise exponential lower bound applies. [1]

Bellman rank

Rank of average Bellman error matrices (maximum over $h=1, \dots, H$)

- Size $|F| \times |F|$
- Q^* has 0 Bellman error on all roll-in policies (col of 0's)
- Sample-efficient to evaluate a row at a time: generate trajectories using $\pi_{f'}$ until h , then random action + importance weighting

candidate value function $f \in F$

roll-in policy $\pi_{f'} : f' \in F$

$$\mathcal{E}(f, \pi_{f'}, h) := \mathbb{E}_{\substack{a_{1:h-1} \sim \pi_{f'} \\ a_h = \pi_{f'}(x_h)}} [f(x_h, a_h) - r_h - \max_{a \in A} f(x_{h+1}, a)]$$

RL problems with low Bellman rank

Problem	Bound	Proof Sketch
<p>Tabular MDP (context = state) PAC Learning: known (e.g., [2])</p>	Bellman rank \leq # states	<p>Bellman error matrix $=$ State distribution induced by π_f \times Bellman error of f on each state</p>
<p>POMDP with rich obs. and reactive value function (context = current obs.) Extends [1]</p>	Bellman rank \leq # hidden states	<p>hidden state \rightarrow rich obs. \rightarrow value</p>
<p>Large MDP with low-rank transition (context = state) New</p>	Bellman rank \leq rank of transition matrix	<p>hidden factor state \rightarrow $\pi_{f'}$</p>
<p>Large MDP with Q^*-irrelevant abstraction (context = abstract state) Known [3]</p>	Bellman rank \leq $\text{poly}(\# \text{ abstract states, \# actions})$	<p>abstract state \rightarrow state \rightarrow $\pi_{f'}$</p>
<p>$P\mathcal{T} _h$ PSRs with rich obs. and reactive value function (context = current obs.) New</p>	Bellman rank \leq $\text{poly}(\text{system dim, \# actions})$	Expressing Bellman error matrix using a submatrix of the System Dynamics Matrix (naturally low-rank for PSRs). (histories = all $(h-1)$ -long seq., tests = length 2 seq.)
<p>Linear Quadratic Regulators (context = state) Known [4]</p>	Bellman rank \leq $\text{poly}(\text{state space dim, action space dim})$	<ul style="list-style-type: none"> Need policy class + state-value function class representation (see Extensions). Crucially depends on the choice of function classes: linear policies + quadratic value functions. Algorithm does not apply as-is due to continuous action space.

References

[1] Krishnamurthy, Agarwal, and Langford. PAC reinforcement learning with rich observations. NIPS 2016.
 [2] Kearns and Singh. Near-Optimal Reinforcement Learning in Polynomial Time. ML 2000.
 [3] Lihong Li. A unifying framework for computational reinforcement learning theory. PhD thesis, 2000.
 [4] Osband and Van Roy. Model-based reinforcement learning and the eluder dimension. NIPS 2014.

OLIVE (Optimism-Led Iterative Value-function Elimination)

Simplified Algorithm (assuming no statistical errors)

- Generate trajectories using $\pi_{f'}$.
- Eliminate all f with non-zero Bellman error.
- Choose a new $\pi_{f'}$ **optimistically**: f is the maximizer of $V_f := \mathbb{E}[f(x_1, \pi_f(x_1))]$ among the surviving functions.
- Repeat until $V^{\pi_{f'}} \geq V^*$ ($\geq V^*$).

Full matrix view

Average Bellman error $\mathcal{E}(f, \pi_{f'}, h)$

Analysis of iteration complexity

- If dark blue vectors are linearly indep., #iterations (for h) \leq Bellman rank.
- Suffices to find a row that contains non-zero entry in surviving columns.
- Optimism** finds the row with a non-zero diagonal entry (for some h).

$$V_{f'} - V^{\pi_{f'}} = \sum_{h=1}^H \mathcal{E}(f', \pi_{f'}, h)$$

Factored matrix view

$\pi_{f'} \times f = 0$

f survives if

Analysis that considers statistical errors (Bellman rank = 2)

[Todd'82]: $|\langle \rightarrow, \rightarrow \rangle| = \Omega(\sqrt{M}\phi) \Rightarrow$ Significant reduction in ellipsoid volume

Sample complexity: $\tilde{O}\left(\frac{M^2 H^3 |A|}{\epsilon^2} \log(|\mathcal{F}|/\delta)\right)$ M : Bellman rank

Extensions

- Can use doubling trick to guess unknown Bellman rank.
- Can compete with functions that have small non-zero Bellman errors.
- Can work with policy class + V -value function class (as opposed to Q).
 - Compete with the best (policy, V -value function) pair that respects *Bellman equation for policy evaluation*.
- Can accommodate infinite classes with bounded statistical complexity.
- Can handle approximately low-rank Bellman error matrices.