



Bin Luo Contact The Chinese University of Hong Kong Shatin, N.T., Hong Kong SAR binluo@link.cuhk.edu.hk WeChat: RobinBinLuo

Quantum Algorithms for Finite-horizon Markov Decision Processes

Bin Luo¹, Yuwen Huang¹, Jonathan Allcock², Xiaojun Lin¹, Shengyu Zhang², John C.S. Lui¹ ¹The Chinese University of Hong Kong, ²Tencent Quantum Laboratory

Exact Dynamics Setting

Background

- Assumption: Dynamics of the environment are **fully known** to the agent. (\Leftrightarrow The robot has the map of the maze.)
- Classical algorithms assume they have access to a classical oracle $O_{\mathcal{M}}$: $(s, a, h, s') \mapsto (r_h(s, a), P_h(s'|s, a)).$
- Classical value iteration (VI) algorithm can obtain π^* and V_0^* with $O(S^2AH)$ queries to $O_{\mathcal{M}}$.
- Define Bellman optimality value operator $f^{\text{Define } P_{h|s,a}(s') = P_h(s'|s,a)}$ $[\mathcal{T}^h(V_{h+1})]_s \coloneqq \max_{a \in \mathcal{A}} \{r_h(s,a) + P_{h|s,a}^T \tilde{V}_{h+1}\}.$
- VI algorithm repeatedly applies \mathcal{T}^h on V_{h+1} in a backward manner with $V_H = \mathbf{0}$ and updates the policy following $\pi(s,h) = \arg\max \{r_h(s,a) + P_{h|s,a}^T V_{h+1}\}.$

Summary of the Results

Table 1. Classical and quantum query complexities for solving general finite-horizon MDPs

		· · ·		
	Goal	Classical query complexity		
	Coat	Upper bound	Lower bour	
	Optimal π^* , V_0^*	$O(S^2AH)$	$\Omega(S^2A)$	
	ϵ -accurate estimate of π^* and $\{V_h^*\}_{h=0}^{H-1}$	$O(S^2AH)$	$\Omega(S^2A)$	
assical Lower Bound				

Theorem (informal version): Given access to the classical oracle $O_{\mathcal{M}}$, any algorithm, which outputs ϵ -approximations of $\{V_h^*\}_{h=0}^{H-1}$ or π^* with probability at least 0.9, must require at least

- $\Omega(S^2A)$ queries to $O_{\mathcal{M}}$ on the worst case of input \mathcal{M} .
- Note that the above theorem implies that it also requires at least $\Omega(S^2A)$ queries
- to $O_{\mathcal{M}}$ to obtain $\{V_h^*\}_{h=0}^{H-1}$ or π^* .
- Question: whether quantum algorithms can break this barrier in the dependence on action space size $(|\mathcal{A}| \coloneqq A)$ or state space size $(|\mathcal{S}| \coloneqq S)$?

Quantum Oracle of Finite-horizon MDPs

Definition: A quantum oracle of a finite-horizon MDP is a unitary operator O_{OM} such that $O_{QM}: |s\rangle |a\rangle |h\rangle |s'\rangle |0\rangle |0\rangle \mapsto |s\rangle |a\rangle |h\rangle |s'\rangle |r_h(s,a)\rangle |P_h(s'|s,a)\rangle$ for all $(s, a, h, s') \in S \times A \times [H] \times S$.

Quantum Speedup on Action Space Size (A)

Quantum Maximum Searching (QMS) Algorithm [Durr et al., 1999]

- Problem: For an unsorted list $f \in \mathbb{R}^N$, one wants to find the index *i* such that $f(i) = \max_{i \in [N]} f(j)$.
 - Classical algorithm: $\Theta(N)$ queries to the vector f.
 - Quantum algorithm: $\Theta(\sqrt{N})$ queries to the vector f.
 - Suppose N = 1,000,000 and it takes 1 second for each query, then the classical algorithm needs roughly **11.5 days**, but QMS algorithm only needs roughly **17 minutes!**

Quantum Value Iteration QVI-1 Algorithm

- Main idea: apply QMS algorithm when taking the maximum over the whole action space in the classical value iteration algorithm.
- Output: optimal policy π^* and optimal V-value function V_0^* .
- Query complexity: $\tilde{O}(S^2\sqrt{A}H)$ queries to the quantum oracle O_{QM} .

Quantum Speedup on State Space Size (S)

- **New Quantum Subroutine**: Quantum Mean Estimation with Binary Oracles (QMEBO)
- Query complexity: $O\left(\left(\frac{\sqrt{N}}{\epsilon} + \sqrt{\frac{N}{\epsilon}}\right)\log\left(\frac{1}{\delta}\right)\right)$ queries to the function f.

Quantum Value Iteration QVI-2 Algorithm

- Output: ϵ -approximations of π^* and $\{V_h^*\}_{h=0}^{H-1}$.
- Query complexity: $\tilde{O}(S^{1.5}\sqrt{A}H^3/\epsilon)$ queries to the quantum oracle O_{OM} .

Notations

S: state space size A: action space size H: total time horizon ϵ : error term δ : failure probability $\widetilde{O}(\cdot)$ and $\widetilde{\Omega}(\cdot)$ ignore logarithmic factors

References



. [Durr et al., 1999] *A quantum algorithm for finding the minimum*. Christoph Durr, Peter Hoyer.

3. [Montanaro, 2015] Quantum speedup of Monte Carlo methods. Ashley Montanaro.

4. [Sidford et al., 2018] Near-optimal time and sample complexities for solving Markov decision processes with a generative model. Aaron Sidford, Mengdi Wang, Xian Wu, Lin F. Yang, Yinyu Ye. 5. [Wang et al., 2021] Quantum algorithms for reinforcement learning with a generative model. Daochen Wang, Aarthi Sundaram, Robin Kothari, Ashish Kapoor, Martin Roetteler

ative Model	Setting				
tions: Agent Agent Agent Agent $r_{h-1}(s_{h-1}, a_{h-1})$ $r_{h}(s_{h}, a_{h})$ $r_{h}(s_{h}, a_{h})$ $r_{h}(s_{h}, a_{h})$ $r_{h}(s_{h}, a_{h})$ $r_{h}(s_{h}, a_{h})$					
+ $s_h^i(s, a) \stackrel{i.i.d.}{\sim} P_{h s,a}, i = 1,, N$					
sor.) nations of π^* , $\{V_h^*\}_{h=0}^{H-1}$	Sensor				
complexities for solving general finite-horizon MDPs					
iery complexity	Quantum query	complexity			
Lower bound	Upper bound	Lower bound			
$\widetilde{\mathbf{\Omega}}\left(\frac{SAH^3}{\epsilon^2}\right)$	$\widetilde{O}\left(rac{SAH^{2.5}}{\epsilon} ight)$ [QVI-4]	$\widetilde{\mathbf{\Omega}}\left(rac{SAH^{1.5}}{\epsilon} ight)$			
$\widetilde{\Omega}\left(\frac{SAH^3}{\epsilon^2}\right)$	$\widetilde{O}\left(\frac{SAH^{2.5}}{\epsilon}\right) [QVI-4]$ $\widetilde{O}\left(\frac{S\sqrt{A}H^{3}}{\epsilon}\right) [QVI-3]$	$\widetilde{\mathbf{\Omega}}\left(\frac{S\sqrt{A}H^{1.5}}{\epsilon}\right)$			
of Finite-horizon MDPs					
acle of a finite-horizon MDP is a unitary operator G such that					
$\langle H \mapsto s\rangle a\rangle h\rangle \left(\sum_{s' \in S} \sqrt{P_h(s' s,a)} s'\rangle w_{s'}\rangle \right)$ e $ w_{s'}\rangle$ are arbitrary auxiliary states.					
$ONAE \setminus A \mid a a a i + b m a f M a a tana a a a a a a a a a a a a a a a $					
($[0, u]$, one needs to obtain an ϵ -estimate of $\mathbb{E}[X]$. ($[0, u]$, one needs to obtain an ϵ -estimate of $\mathbb{E}[X]$. ($[0, u]$, one needs to obtain an ϵ -estimate of $\mathbb{E}[X]$. ($[0, u]$, one needs to obtain an ϵ -estimate of $\mathbb{E}[X]$. ($[0, u]$, one needs to obtain an ϵ -estimate of $\mathbb{E}[X]$.					
that $O(\sigma^2/\epsilon^2)$ classical samples are required. antum samples. I-3 Algorithm but it applies QME1 to obtain $\frac{\epsilon}{H}$ -estimates of $P_{h s,a}^T V_{h+1}$. $\{V_h^*\}_{h=0}^{H-1}$					
ries to the quantum generative oracle G .					
I-4 Algorithm					
ntum adaptations of "variance reduction" and "total variance" 8] by applying QME1 and QME2 algorithms. ${}_{h}^{*}{}_{h=0}^{H-1}$ and $\{Q_{h}^{*}\}_{h=0}^{H-1}$.					
ies to the quantum generative oracle ${\cal G}_{\cdot}$					
r Bounds					
ⁱ inite-horizon MDP can be reduced to solving an infinite-horizon of solving finite-horizon MDP inherits from those of the 2021].					
proptotically optimal , up to log terms, for computing ϵ - $\{Q_h^*\}_{h=0}^{H-1}$, provided a constant time horizon.					

2. [Li et al., 2020] Breaking the sample size barrier in model-based reinforcement learning with a generative model. Gen Li, Yuting Wei, Yuejie Chi, Yuxin Chen