

Optimizing Average Reward MDPs with a Generative Model

Aaron Sidford

Stanford University

Departments of Management Science
& Engineering and Computer Science

Contact Info:

- email: sidford@stanford.edu
- website: www.aaronsidford.com

Efficiently Solving MDPs with Stochastic Mirror Descent

joint with Yujia Jin (ICML 2020, arXiv: 2008.12776)

Towards Tight Bounds on the Sample Complexity of Average-reward MDPs

joint with Yujia Jin (ICML 2021, arXiv: 2106.07046)



Yujia Jin

Thank you for slide material!

This Talk

Part 1
Problem and Results

Part 2
Approach #1

[JinS20]

Part 3
Approach #2

[JinS21]

Part 4
Lower Bound

[JinS21]

Take-aways

- Algorithmic tools for solving MDPs!
- Open problem!

Markov Decision Process (MDPs)

Setup

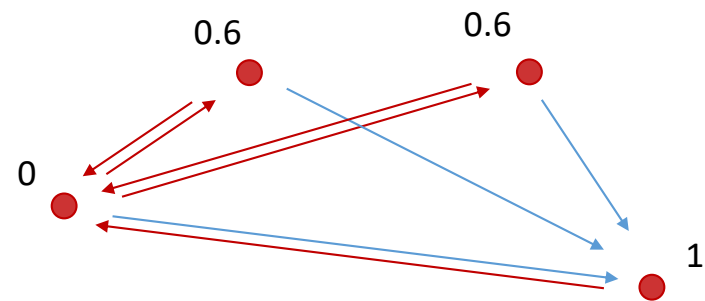
- **States:** finite set S
- **Actions:** finite set A_s for each $s \in S$
- **Transition probabilities:** $p_{s,a} \in \Delta^S$ for each $s \in S, a \in A_s$
- **Rewards:** $r \in [-1,1]^A$ for $A = \cup_{s \in S} A_s$

Goal: compute an ϵ -optimal policy

- Randomized policy: $\pi(s) \in \Delta^{A_s}$ for all s
- Deterministic policy: $\pi(s) \in A_s$ for all s

Won't always distinguish between the two but may mention open problems.

There are other functions, e.g. finite horizon, see next talk!

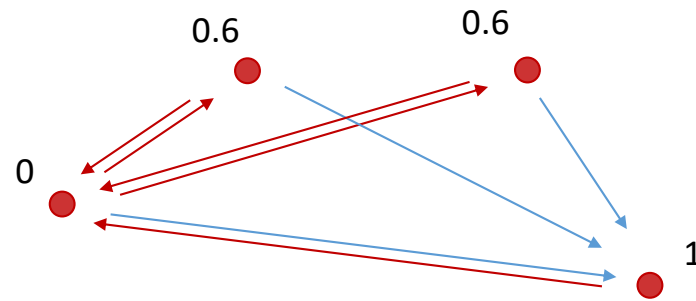


What reward function?

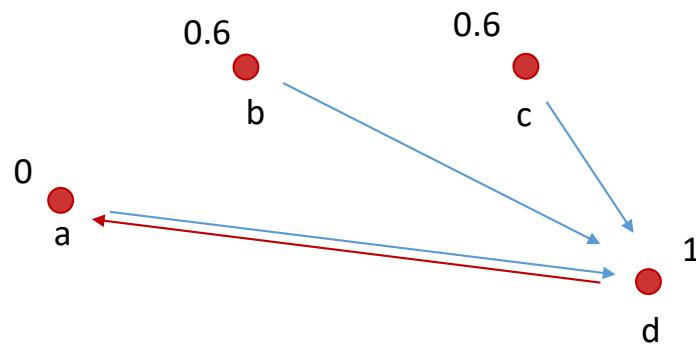
- **Discounted reward (DMDP):** $\gamma \in (0,1)$ and $q \in \Delta^S$
 $v_{\gamma,q}^{\pi} \stackrel{\text{def}}{=} \mathbb{E}_{s_t, \pi(s_t)} \sum_t \gamma^t r_{s_t, \pi(s_t)}$ for $s_0 \sim q$
- **Average reward (AMDP):** $\gamma \rightarrow 1$
 $v^{\pi} \stackrel{\text{def}}{=} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{s_t, \pi(s_t)} \sum_{t \in [T]} r_{s_t, \pi(s_t)}$

For discounted approximate policies there is a difference between whether reward is for specific q or all $q \in \Delta^S$.

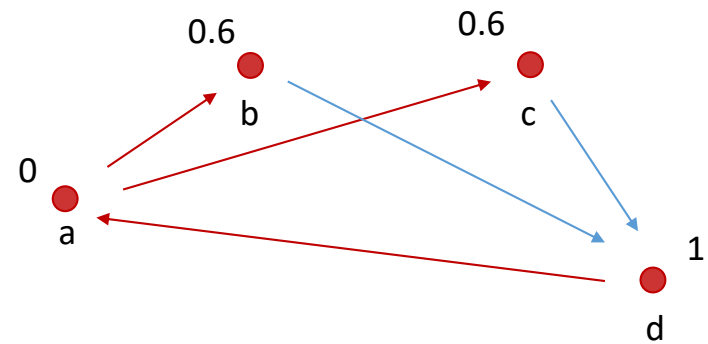
Reward Functions



DMDP optimal policy
(small γ)



AMDP Optimal Policy
($\gamma \rightarrow 1$)



The Problem

Setup

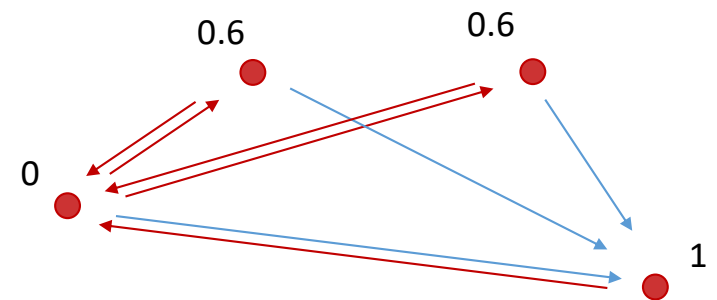
- **States:** finite set S
- **Rewards:** $r \in [-1,1]^S$
- **Actions:** finite set A_s for each $s \in S$
- **Transition probabilities:** $p_{s,a} \in \Delta^S$ for each $s \in S, a \in A_s$

Goal: compute an ϵ -optimal policy

- Randomized policy: $\pi(s) \in \Delta^{A_s}$ for all s
- Deterministic policy: $\pi(s) \in A_s$ for all s

Reward Function

- Discounted: $v_{\gamma,q}^{\pi} = \mathbb{E}_{s_t, \pi(s_t)} \sum_t \gamma^t r_{s_t, \pi(s_t)}$ for $s_0 \sim q$
- Average: $v^{\pi} = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{s_t, \pi(s_t)} \sum_{t \in [T]} r_{s_t, \pi(s_t)}$



Generative Model Sample Complexity

- States, actions, and rewards are known
- Transition probabilities unknown
- Given any s, a can *query* generative model for a sample from $p_{s,a}$
- **Question:** how many samples needed to compute an ϵ -optimal policy?

State-of-the-art

sample complexities


- States S
- A_{tot} total state action pairs
- Discount factor γ
- Max ratio of stationary probability τ
- Largest mixing time of any policy t_{mix}

	<u>Upper Bound</u>	<u>Lower Bound</u>	
Discounted Reward (DMDP)	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13,SWWYY18, W19,AKY20,LWCGC20]	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13]	
Average Reward (AMDP)	$\frac{A_{\text{tot}} t_{\text{mix}}^2 \tau^4}{\epsilon^2}$ [W17]	?	<u>Open Problem</u> What is the optimal sample complexity for $\epsilon = \tilde{O}(1)$?
Our AMDP Results	$\frac{A_{\text{tot}} t_{\text{mix}}^2}{\epsilon^2}$ [JS20]	$\frac{A_{\text{tot}} t_{\text{mix}}}{\epsilon^3}$ [JS21]	Optimal sample complexities for $\epsilon = \tilde{\Omega}(1)$!
		↕ <i>oblivious samples</i> ↕	

Ignoring log factors, differences between $|S||A|$ and A_{tot} , the domain of r , etc.

This Talk

Part 1
Problem and Results



Part 2
Approach #1

[JinS20]

Part 3
Approach #2

[JinS21]

Part 4
Lower bound

[JinS21]

Take-aways

- Algorithmic tools for solving MDPs!
- Open problem!

State-of-the-art

sample complexities

- States S
- A_{tot} total state action pairs
- Discount factor γ
- Max ratio of stationary probability τ
- Largest mixing time of any policy t_{mix}

	<u>Upper Bound</u>	<u>Lower Bound</u>	
Discounted Reward (DMDP)	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13,SWWYY18, W19,AKY20,LWCGC20]	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13]	
Average Reward (AMDP)	$\frac{A_{\text{tot}} t_{\text{mix}}^2 \tau^4}{\epsilon^2}$ [W17]	?	<u>Open Problem</u> What is the optimal sample complexity for $\epsilon = \tilde{O}(1)$?
Our AMDP Results	$\frac{A_{\text{tot}} t_{\text{mix}}^2}{\epsilon^2}$ [JS20]	$\frac{A_{\text{tot}} t_{\text{mix}}}{\epsilon^3}$ [JS21]	Optimal sample complexities for $\epsilon = \tilde{\Omega}(1)$!
		$\frac{A_{\text{tot}} t_{\text{mix}}}{\epsilon^2}$ [JS21]	

↕ oblivious samples ↕

Ignoring log factors, differences between $|S||A|$ and A_{tot} , the domain of r , etc.

Approach #1: Convex Optimization

MDP

- state space $s \in S$
- actions A_s for $a \in S$ and $A \stackrel{\text{def}}{=} \bigcup_{s \in S} A_s$
- transition probabilities $p_{s,a} \in \Delta^S$
- rewards $r_{s,a} \in [-1,1]$

Convex formulation

Same / similar to [W17]

$$\min_{v \in t_{\text{mix}} \cdot [-1,1]^S} \max_{\mu \in \Delta^A} \mu^\top [(P - E)v + r]$$

- row $a \in A_s$ of P is $p_{s,a}$
- row $a \in A_s$ of E is e_a
- ℓ_1 norm of each row of $P - E$ is ≤ 2

Box Simplex Game!

[S17,JST19,CST21]

Solver

Related to [W17,CJST19,CJST20]

- Stochastic mirror descent with careful local norm analysis
- $\tilde{O}(A_{\text{tot}} t_{\text{mix}}^2 / \epsilon^2)$ steps and $\tilde{O}(1)$ per step
- General result about box simplex games!

Rounding

Similar observation / approach taken in [CCBG20] for DMDP

- Scale μ across each A_s so probability distribution
- Lemma: ϵ -approximate $\mu \Rightarrow O(\epsilon)$ -approximate policy
- $\Rightarrow \tilde{O}(A_{\text{tot}} t_{\text{mix}}^2 / \epsilon^2)$ samples to solve an AMDP [JS20]

Discussion

Theorem [JS20]: Can compute ϵ -optimal policy to AMDP using $\tilde{O}\left(\frac{A_{\text{tot}} t_{\text{mix}}^2}{\epsilon^2}\right)$ queries

Properties of resulting algorithm

- Queries are dynamic, which state chosen to query depends on algorithm
- Outputs a randomized property

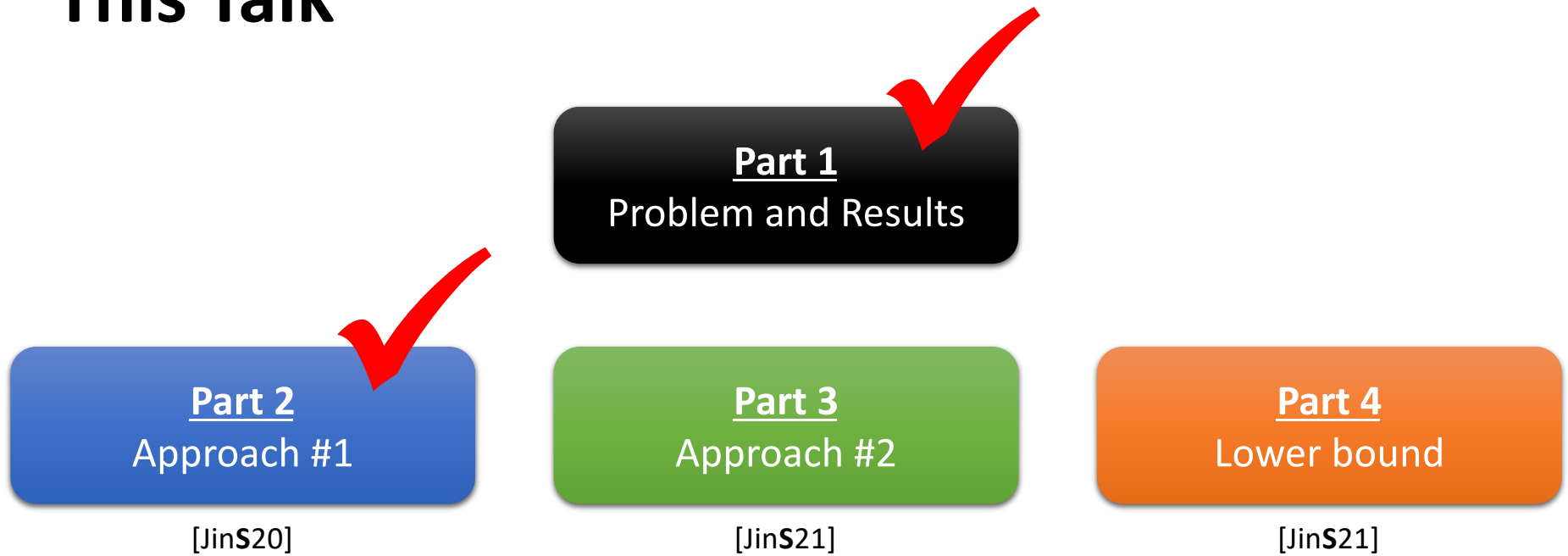
Generalizations

- General sublinear box-simplex solver!
- Recovers $\tilde{O}\left(\frac{A_{\text{tot}}}{(1-\gamma)^4 \epsilon^2}\right)$ sample bound of other convex optimization approach to solving DMDP [CCBG20] (for a fixed initial distribution on vertices)
- Generalizes to constrained MDPs!

Open Problems

- Convex approach completely matching state-of-the-art for DMDPs?

This Talk



Take-aways

- Algorithmic tools for solving MDPs!
- Open problem!

State-of-the-art

sample complexities

- States S
- A_{tot} total state action pairs
- Discount factor γ
- Max ratio of stationary probability τ
- Largest mixing time of any policy t_{mix}

	<u>Upper Bound</u>	<u>Lower Bound</u>	
Discounted Reward (DMDP)	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13,SWWYY18, W19,AKY20,LWCGC20]	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13]	
Average Reward (AMDP)	$\frac{A_{\text{tot}} t_{\text{mix}}^2 \tau^4}{\epsilon^2}$ [W17]	?	<u>Open Problem</u> What is the optimal sample complexity for $\epsilon = \tilde{O}(1)$?
Our AMDP Results	$\frac{A_{\text{tot}} t_{\text{mix}}^2}{\epsilon^2}$ [JS20]	$\frac{A_{\text{tot}} t_{\text{mix}}}{\epsilon^3}$ [JS21]	Optimal sample complexities for $\epsilon = \tilde{\Omega}(1)$!

↕ oblivious samples ↕

Ignoring log factors, differences between $|S||A|$ and A_{tot} , the domain of r , etc.

Approach #2: Reduction

MDP

- state space $s \in S$
- actions A_s for $a \in S$ and $A \stackrel{\text{def}}{=} \cup_{s \in S} A_s$
- transition probabilities $p_{s,a} \in \Delta^S$
- rewards $r_{s,a} \in [-1,1]$

Problem: Given an MDP, find ϵ -optimal policy given generative model access.

- **Discounted reward (DMDP):** $\gamma \in (0,1)$ and $q \in \Delta^S$: $v_{\gamma,q}^\pi \stackrel{\text{def}}{=} \mathbb{E}_{s_t, \pi(s_t)} \sum_t \gamma^t r_{s_t, \pi(s_t)}$ for $s_0 \sim q$
- **Average reward (AMDP):** $v^\pi \stackrel{\text{def}}{=} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{s_t, \pi(s_t)} \sum_{t \in [T]} r_{s_t, \pi(s_t)}$

Lemma: $|v^\pi - (1 - \gamma)v_{\gamma,q}^\pi| \leq 3(1 - \gamma)t_{\text{mix}}$ for all $\gamma \in (0,1)$

Implication: suffices to compute $\epsilon = \Theta\left(\frac{\epsilon}{1-\gamma}\right)$ -approximate policy to DMDP with $\gamma = 1 - \Theta\left(\frac{\epsilon}{t_{\text{mix}}}\right)$

DMDP Bound

$$\frac{A_{\text{tot}}}{(1 - \gamma)^3 \epsilon^2}$$

[AMK13,SWWYY18,
W19,AKY20,LWCGC20]



AMDP Bound

$$\frac{A_{\text{tot}} t_{\text{mix}}}{\epsilon^3}$$

[JS21]

Discussion

Theorem: any $\varepsilon = \Theta\left(\frac{\epsilon}{1-\gamma}\right)$ -optimal policy to DMDP with $\gamma = 1 - \Theta\left(\frac{\epsilon}{t_{\text{mix}}}\right)$ is an ϵ -optimal policy to the AMDP

$\Rightarrow \tilde{O}\left(\frac{A_{\text{tot}} t_{\text{mix}}}{\epsilon^3}\right)$ samples suffice

Note

- To improve on $\tilde{O}(A_{\text{tot}} t_{\text{mix}}^2 \epsilon^{-2})$ need to set $\epsilon = \tilde{\Omega}(t_{\text{mix}}^{-1}) = \Theta((1-\gamma)\epsilon^{-1})$
- Consequently, $\varepsilon = \Omega((1-\gamma)^{-1/2})$ and need [LWCGC20] to improve

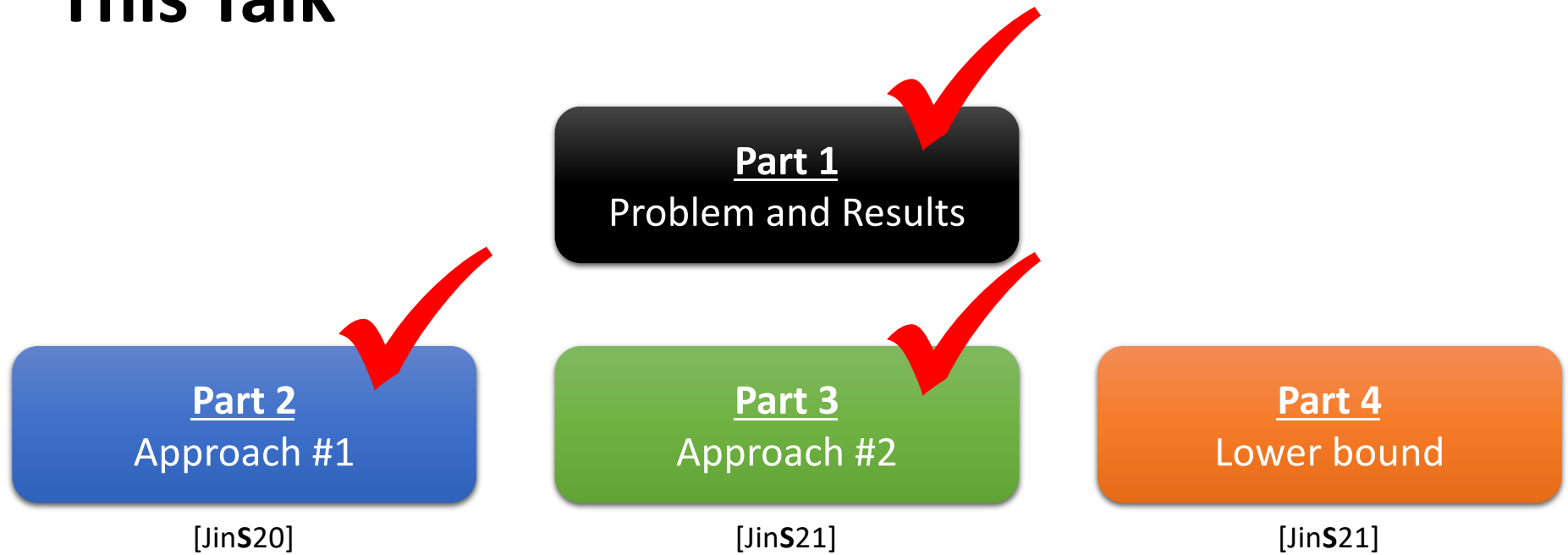
Properties of Resulting Algorithm

- $\tilde{O}(t_{\text{mix}} \epsilon^{-2})$ oblivious samples per state
- Computes deterministic policy
- Only depend on mixing time of deterministic policies

Algorithmic Implication

- Combining with [BLLLSSW21] obtain $\tilde{O}(A_{\text{tot}} |S| + |S|^{2.5})$ time algorithm

This Talk



Take-aways

- Algorithmic tools for solving MDPs!
- Open problem!

State-of-the-art

sample complexities

- States S
- A_{tot} total state action pairs
- Discount factor γ
- Max ratio of stationary probability τ
- Largest mixing time of any policy t_{mix}

	<u>Upper Bound</u>	<u>Lower Bound</u>	
Discounted Reward (DMDP)	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13,SWWYY18, W19,AKY20,LWCGC20]	$\frac{A_{\text{tot}}}{(1-\gamma)^3 \epsilon^2}$ [AMK13]	
Average Reward (AMDP)	$\frac{A_{\text{tot}} t_{\text{mix}}^2 \tau^4}{\epsilon^2}$ [W17]	?	<u>Open Problem</u> What is the optimal sample complexity for $\epsilon = \tilde{O}(1)$?
Our AMDP Results	$\frac{A_{\text{tot}} t_{\text{mix}}^2}{\epsilon^2}$ [JS20]	$\frac{A_{\text{tot}} t_{\text{mix}}}{\epsilon^3}$ [JS21]	Optimal sample complexities for $\epsilon = \tilde{\Omega}(1)$!

↕ oblivious samples ↕

Ignoring log factors, differences between $|S||A|$ and A_{tot} , the domain of r , etc.

Lower Bound Approach

Modify the DMDP construction in [AMK13]

Note

- Essentially reducing AMDP lower bound to DMDP lower bound to best arm identification.
- Proved for oblivious queries. Open / TODO: prove for arbitrary dynamic queries.

MDP

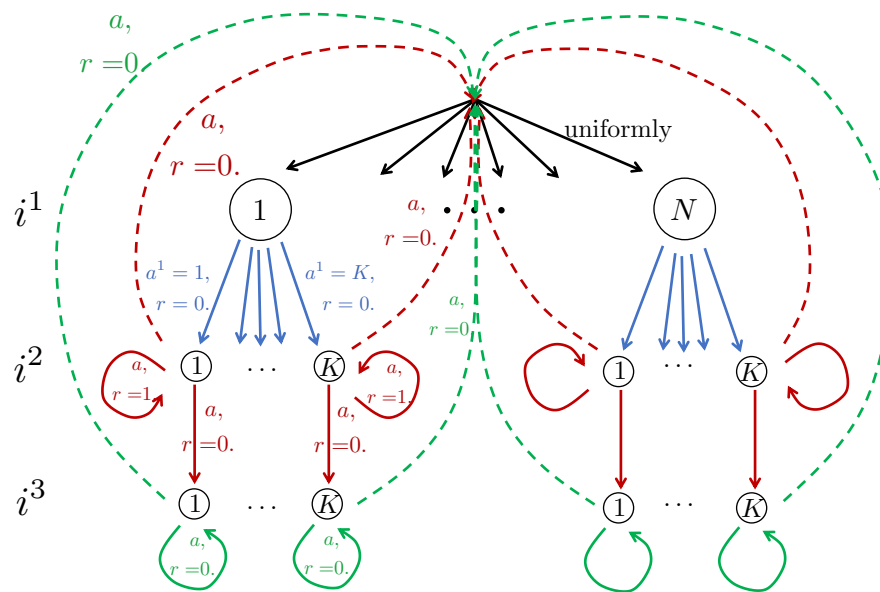
- state space $s \in S$
- actions A_s for $a \in S$ and $A \stackrel{\text{def}}{=} \bigcup_{s \in S} A_s$
- transition probabilities $p_{s,a} \in \Delta^S$
- rewards $r_{s,a} \in [-1,1]$

Level 1: N states, each has K actions that transit to different level 2 state

Level 2: each state s goes uniformly to level 1 with probability $1 - \gamma$, stays with probability γp_s , and goes to level 3 with probability $\gamma(1 - p_s)$

Level 3: each state goes uniformly to to level 1 with probability $1 - \gamma$ and stays with probability γ

Rewards: All 0 except at level 1



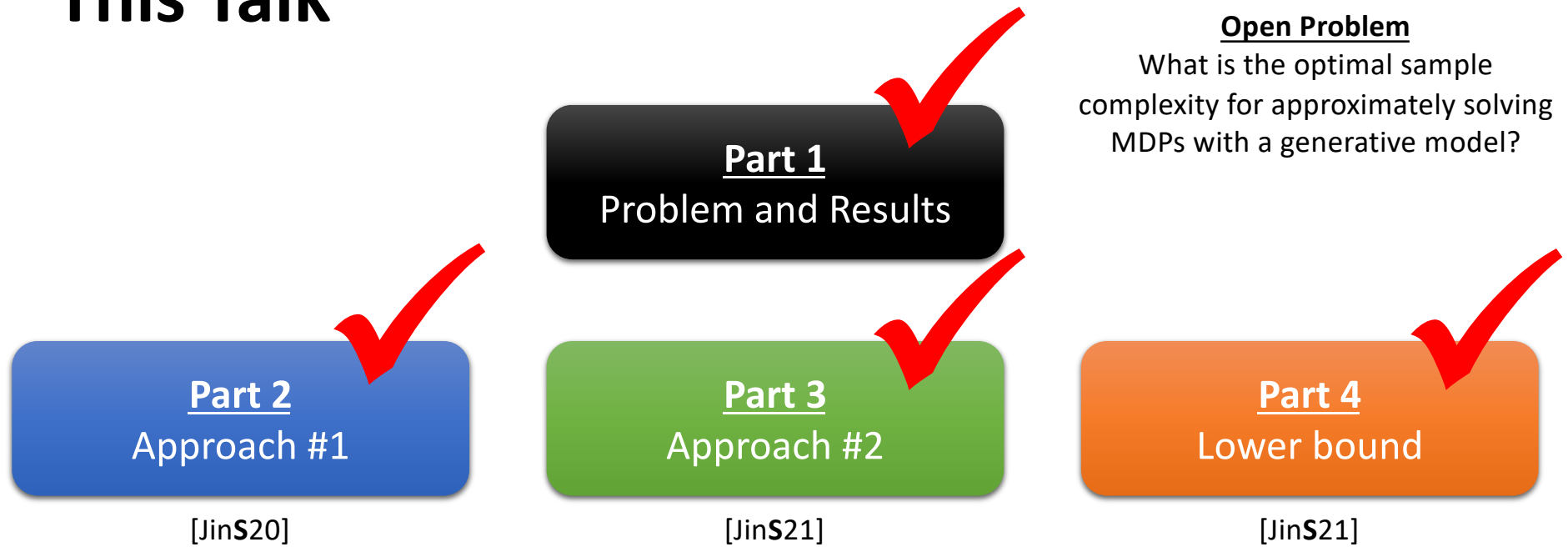
$$\begin{aligned}
 p(\text{red circle}) &= \gamma p_{(i^1, a^1)}, \\
 p(\text{red down arrow}) &= \gamma(1 - p_{(i^1, a^1)}), \\
 p(\text{red left arrow}) &= 1 - \gamma.
 \end{aligned}$$

$$\begin{aligned}
 p(\text{green circle}) &= \gamma, \\
 p(\text{green left arrow}) &= 1 - \gamma.
 \end{aligned}$$

Lower bound strategy

- Each level 1 state has on action to a level 2 state with a higher γ_s
- Lower bound how many samples need to find enough higher p_s

This Talk



Open Problem
What is the optimal sample complexity for approximately solving MDPs with a generative model?

Take-aways

- Algorithmic tools for solving MDPs!
- Open problem!

Thank you

Questions?

arXiv:2008.12776
arXiv:2106.07046



Yujia Jin

Contact Info:

- email: sidford@stanford.edu
- website: www.aaronsidford.com