

# Provably Efficient Offline Reinforcement Learning for Partially Observable Markov Decision Processes

Hongyi Guo<sup>1</sup>, Qi Cai<sup>1</sup>, Yufeng Zhang<sup>1</sup>, Zhuoran Yang<sup>2</sup>, Zhaoran Wang<sup>1</sup>

<sup>1</sup>Northwestern University

<sup>2</sup>Yale University

July 21, 2022

# Background

POMDP

Introduction

Algorithm

Offline reinforcement learning (RL) for partially observable Markov decision processes (POMDPs).

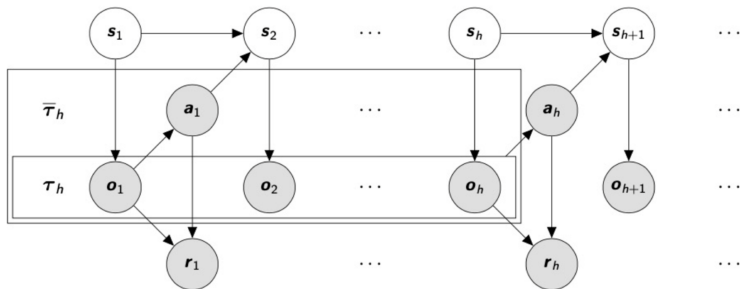


Figure: POMDP

# Overview

POMDP

Introduction

Algorithm

- Linear POMDP (emission kernel and transition kernel are linear in known feature mappings)
- Undercompleteness assumption (observation dist.  $\Rightarrow$  state dist.)

# Overview

POMDP

Introduction

Algorithm

- Linear POMDP (emission kernel and transition kernel are linear in known feature mappings)
- Undercompleteness assumption (observation dist.  $\Rightarrow$  state dist.)
- A pessimistic offline RL algorithm
- Finite sample guarantee  $\tilde{O}(1/\epsilon^2)$

# Intuition

POMDP

Introduction

Algorithm

Define random functions

$$X_{h,a}(o) = \mathbb{1}_{\text{do}(A_{h-1}=a)}^{\bar{\pi}} \{O_h = o\},$$
$$Y_{h,a,a'}(o, o') = \mathbb{1}_{\text{do}(A_{h-1:h}=(a,a'))}^{\bar{\pi}} \{O_{h:h+1} = (o, o')\}.$$

# Intuition

POMDP

Introduction

Algorithm

Define random functions

$$\begin{aligned}X_{h,a}(o) &= \mathbb{1}_{\text{do}(A_{h-1}=a)}^{\bar{\pi}} \{O_h = o\}, \\Y_{h,a,a'}(o, o') &= \mathbb{1}_{\text{do}(A_{h-1:h}=(a,a'))}^{\bar{\pi}} \{O_{h:h+1} = (o, o')\}.\end{aligned}$$

Regress  $Y$  on  $X$  and in the mean time we can solve the model parameter  $\theta$ :

$$Y_{h,a,a'} = \mathbb{F}_{h,a'}^{\theta} X_{h,a} + U_{h,a,a'}, \quad (2.1)$$

where  $\mathbb{F}$  is some linear operator and  $U_{h,a,a'}$  is the zero-mean perturbation term.

# Intuition

POMDP

Introduction

Algorithm

Define random functions

$$\begin{aligned}X_{h,a}(o) &= \mathbb{1}_{\text{do}(A_{h-1}=a)}^{\bar{\pi}} \{O_h = o\}, \\Y_{h,a,a'}(o, o') &= \mathbb{1}_{\text{do}(A_{h-1:h}=(a,a'))}^{\bar{\pi}} \{O_{h:h+1} = (o, o')\}.\end{aligned}$$

Regress  $Y$  on  $X$  and in the mean time we can solve the model parameter  $\theta$ :

$$Y_{h,a,a'} = \mathbb{F}_{h,a'}^{\theta} X_{h,a} + U_{h,a,a'}, \quad (2.1)$$

where  $\mathbb{F}$  is some linear operator and  $U_{h,a,a'}$  is the zero-mean perturbation term.

- $U$  is correlated with  $X$ !  $\Rightarrow$  IV regression.

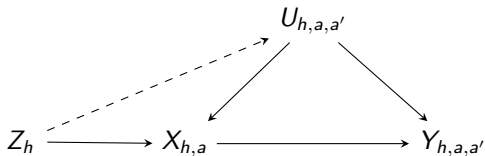
# IV Regression

POMDP

Introduction

Algorithm

An instrumental variable (IV) is correlated with  $X$  but uncorrelated with the perturbation  $U$ . We choose  $Z_h = O_{h-1}$  as the IV.



**Figure:** The relationship between  $X_{h,a}$ ,  $Y_{h,a,a'}$ ,  $U_{h,a,a'}$ , and  $Z_h$ . The arrows indicate the dependency between those variables. The dashed arrow indicates two uncorrelated variables. In this figure,  $U_{h,a,a'}$  affects both  $X_{h,a}$  and  $Y_{h,a,a'}$  directly,  $Z_h$  only affects  $X_{h,a}$  directly, and  $Z_h$  and  $U_{h,a,a'}$  are uncorrelated.



# Algorithm Overview

POMDP

Introduction

Algorithm

With the help of the IV, we have

$$\mathbb{E}[Y_{h,a,a'} | Z_h] = \mathbb{F}_{h,a'}^\theta \mathbb{E}[X_{h,a} | Z_h]. \quad (2.2)$$

Our algorithm has the follows steps:

- 1 Construct estimators of  $\mathbb{E}[Y_{h,a,a'} | Z_h]$  and  $\mathbb{E}[X_{h,a} | Z_h]$  from the dataset.
- 2 Construct confidence region  $\hat{\Theta}$  of the model parameter so that  $\mathbb{E}[Y_{h,a,a'} | Z_h]$  and  $\mathbb{F}_{h,a'}^\theta \mathbb{E}[X_{h,a} | Z_h]$  are close enough for  $\theta \in \hat{\Theta}$ .
- 3 Pessimism Planning

$$(\hat{\pi}, \hat{\theta}) = \operatorname{argmax}_{\pi \in \Pi} \operatorname{argmin}_{\theta \in \hat{\Theta}} \mathcal{J}(\theta, \pi).$$

*Thank you!*

Please check our poster @Hall E