

# Outline

1. Introduction and Setting
2. Offline batch evaluation using models
3. Offline batch evaluation using Q functions
4. Offline batch evaluation using importance sampling

# Per Decision Importance Sampling (PDIS)

- Leverage temporal structure of the domain (**similar to policy gradient**)

$$IS(D) = \frac{1}{n} \sum_{i=1}^n \left( \prod_{t=1}^L \frac{\pi_e(a_t | s_t)}{\pi_b(a_t | s_t)} \right) \left( \sum_{t=1}^L \gamma^t R_t^i \right)$$

$$\text{PDIS} = \sum_{t=1}^L \gamma^t \underbrace{\frac{1}{n} \sum_{i=1}^n \left( \prod_{\tau=1}^t \frac{\pi_e(a_\tau | s_\tau)}{\pi_b(a_\tau | s_\tau)} \right) R_t^i}_{\text{}}$$

On the board I also covered weighted importance sampling (WIS)

# Importance Sampling Variance

- Importance sampling, like Monte Carlo estimation, is generally high variance
- Importance sampling is particularly high variance for estimating the return of a policy in a sequential decision process

$$= \sum_{i=1, \tau_i \sim \pi_b}^N R_{\tau_i} \prod_{t=1}^{H_i} \frac{p(a_{it} | \pi, s_{it})}{p(a_{it} | \pi_b, s_{it})}$$

$$\checkmark = 1 \quad \checkmark^S \leq 10^H$$

- Variance can generally scale exponentially with the horizon
  - a. Concentration inequalities like Hoeffding scale with the largest range of the variable
  - b. The largest range of the variable depends on the product of importance weights
  - c. **Check your understanding: for a  $H$  step horizon with a maximum reward in a single trajectory of 1, and if  $p(a|s, \pi_b) = .1$  and  $p(a|s, \pi) = 1$  for each time step, what is the maximum importance-weighted return for a single trajectory?**

$$10^H \left( \frac{1}{1/10} \right)^H$$

$$R_{\tau_i} \prod_{t=1}^{H_i} \frac{p(a_{it} | \pi, s_{it})}{p(a_{it} | \pi_b, s_{it})}$$

# Recent Directions

- Leveraging Markov structure to break curse of horizon.
  - Marginalized importance sampling (state-action distribution)
  - Dai, Nachum, Chow, Li (dualdice, coindice) 2019/2020
  - Liu, Li, Tang, Zhou Neurips 2018
- Doubly robust estimation
- Blended estimators

# Control variates

- Given:  $X, Y, \mathbb{E}[Y]$
- Estimate  $\mu = \mathbb{E}[X]$
- $\hat{\mu} = X - Y + \mathbb{E}[Y]$
- Unbiased:  
$$\mathbb{E}[\hat{\mu}] = \mathbb{E}[X - \mathbb{E}[Y] + \mathbb{E}[Y]] = \mathbb{E}[X] + \underbrace{\mathbb{E}[Y] - \mathbb{E}[Y]}_{=0} = \mathbb{E}[X] = \mu$$
- Variance:

$$\begin{aligned}\text{Var}(\hat{\mu}) &= \text{Var}(X - Y + \mathbb{E}[Y]) = \text{Var}(X - Y) \\ &= \text{Var}(X) + \text{Var}(Y) - 2 \text{cov}(X, Y)\end{aligned}$$

# Control variates

- Given:  $X, Y, \mathbb{E}[Y]$
- Estimate:  $\mu = \mathbb{E}[X]$
- $\hat{\mu} = X - \underline{Y} + \underline{\mathbb{E}[Y]}$
- Unbiased:  
$$\mathbb{E}[\hat{\mu}] = \mathbb{E}[X - Y + \mathbb{E}[Y]] = \mathbb{E}[X] - \mathbb{E}[Y] + \mathbb{E}[Y] = \mathbb{E}[X] = \mu$$
- Variance:

$$\begin{aligned}\text{Var}(\hat{\mu}) &= \text{Var}(X - Y + \mathbb{E}[Y]) = \text{Var}(X - Y) \\ &= \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y)\end{aligned}$$

- Lower variance if  $2\text{Cov}(X, Y) > \text{Var}(Y)$
- We call  $Y$  a control variate
- We saw this idea before: **baseline term in policy gradient estimation**

## Off-policy policy evaluation (revisited)

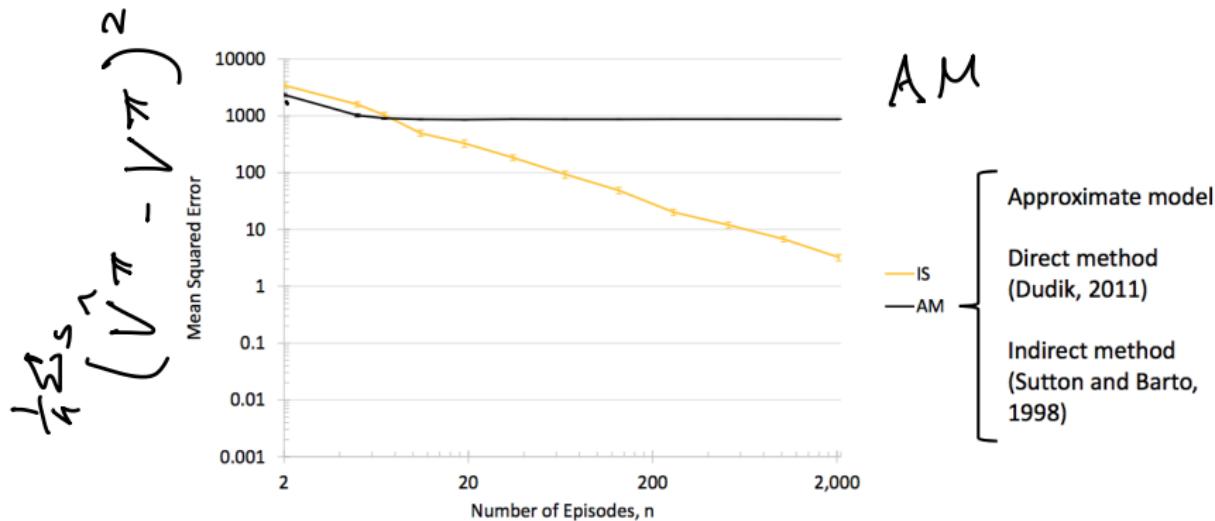
- Idea: add a control variate to importance sampling estimators
  - $X$  is the importance sampling estimator
  - $Y$  is a control variate build from an approximate model of the MDP
- Called the **doubly robust estimator** (Jiang and Li, 2015)
  - Robust to (1) poor approximate model, and (2) **error in estimates of  $\pi_b$** 
    - If the model is poor, the estimates are still unbiased
    - If the sampling policy is unknown, but the model is good, MSE will still be low
- Non-recursive and weighted forms, as well as control variate view provided by Thomas and Brunskill (ICML 2016)

## Off-policy policy evaluation (revisited)

$$DR(\pi_e \mid D) = \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^{\infty} \gamma^t \underbrace{w_t^i}_{\downarrow} \underbrace{(R_t^i - \hat{q}^{\pi_e}(S_t^i, A_t^i))}_{\text{Bellman Error}} + \gamma^t \rho_{t-1}^i \hat{v}^{\pi_e}(S_t^i),$$

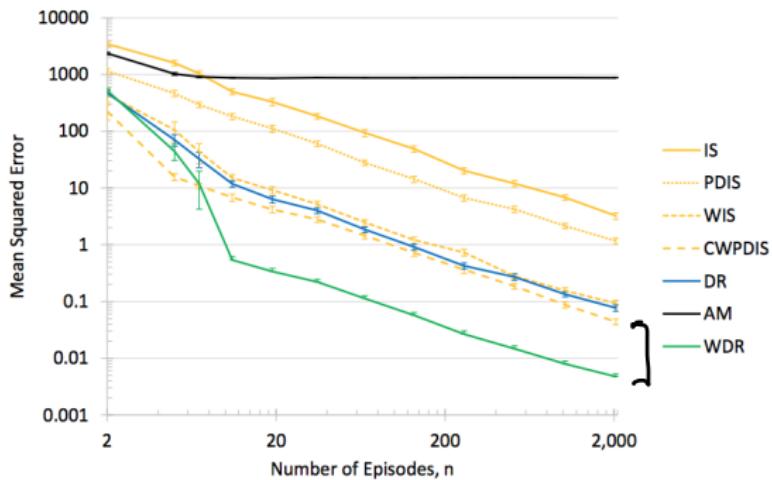
where  $w_t^i = \prod_{\tau_1}^t \frac{\pi_e(a_\tau \mid s_\tau)}{\pi_b(a_\tau \mid s_\tau)}$

# Empirical Results (Gridworld)



AM

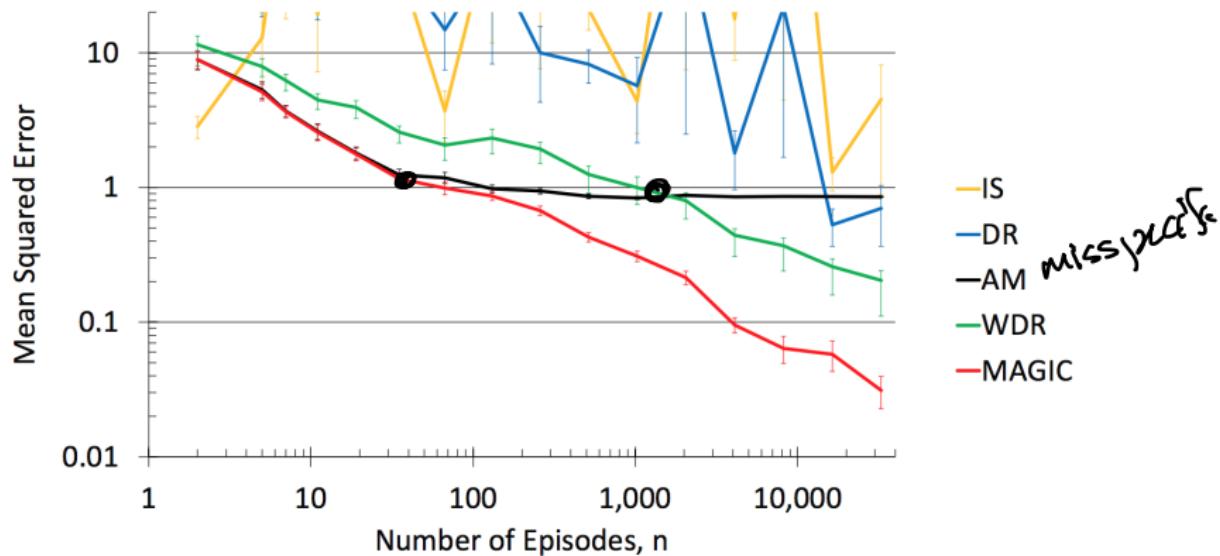
# Empirical Results (Gridworld)



## Off-policy policy evaluation (revisited): Blending

- Importance sampling is unbiased but high variance
- Model based estimate is biased but low variance
- Doubly robust is one way to combine the two
- Can also trade between importance sampling and model based estimate within a trajectory
- MAGIC estimator (Thomas and Brunskill ICML 2016)
- Can be particularly useful when part of the world is non-Markovian in the given model, and other parts of the world are Markov

# Can Need an Order of Magnitude Less Data To Get Good Estimates



# What You Should Know

- Be able to define and apply importance sampling for off policy policy evaluation
- Define some limitations of IS (variance)
- Define why we might want to do batch offline RL policy evaluation and potential applications
- Be aware of the main potential limitations of model and model free methods

# Batch / Offline RL Policy Learning

Emma Brunskill

March 2 2023

CS234

Thanks to Phil Thomas for some figures

# Refresh Your Understanding

Importance sampling (select all that are true)

- Requires the behavior policy to visit all the state--action pairs that would be visited under the evaluation policy in order to get an unbiased estimator ✓
- Is likely to be high variance ✓
- Not Sure

Behavior cloning from demonstrations:

- Reduces batch/offline learning to supervised learning ✓
- May learn a low performing policy if the demonstrations come from a non-expert ✓
- May learn a low performing policy if the demonstrations from an expert ✓
- Could be used to warm start an online reinforcement learning algorithm ✓
- Requires a human to label what they would do at the states visited by the policy learned F dage ✓  
dels ↴
- Not Sure

# Refresh Your Understanding

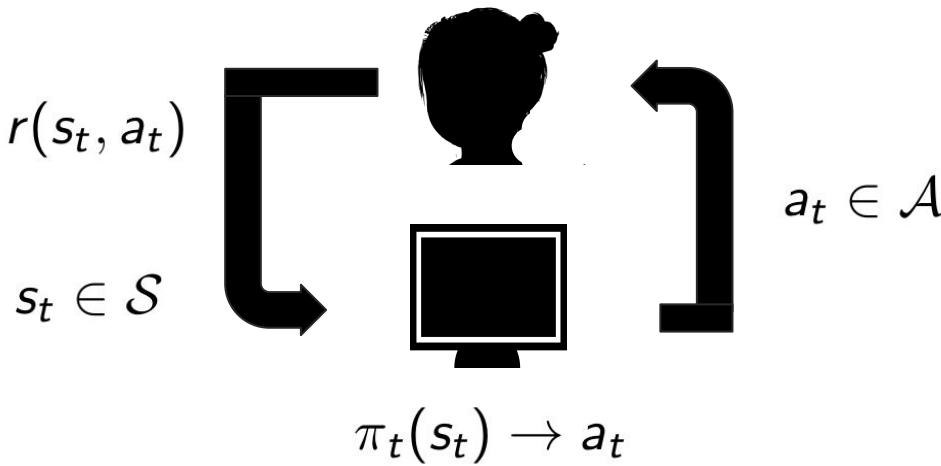
Importance sampling (select all that are true)

- Requires the behavior policy to visit all the state--action pairs that would be visited under the evaluation policy in order to get an unbiased estimator (true)
- Is likely to be high variance (true)
- Not Sure

Behavior cloning from demonstrations:

- Reduces batch/offline learning to supervised learning
- May learn a low performing policy if the demonstrations come from a non-expert
- May learn a low performing policy if the demonstrations from an expert
- Could be used to warm start an online reinforcement learning algorithm
- Requires a human to label what they would do at the states visited by the policy learned
- Not Sure

# Today: Counterfactual / Batch RL



$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

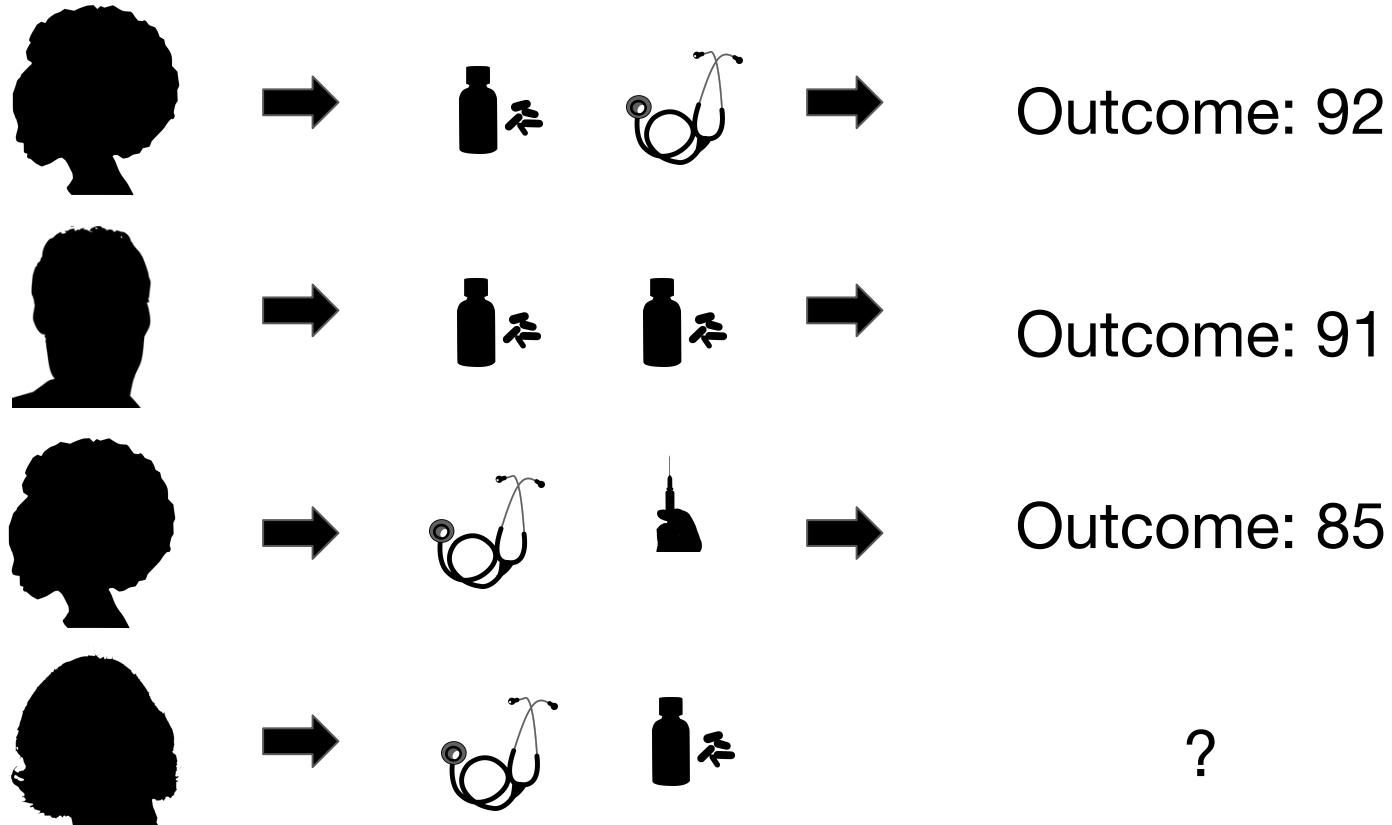
# Where We Are In The Course

1. Learning from offline data
  - a. Imitation learning
  - b. Batch/offline policy evaluation
  - c. **Batch/offline policy learning**
2. Next week
  - a. Guest lecture
  - b. Quiz

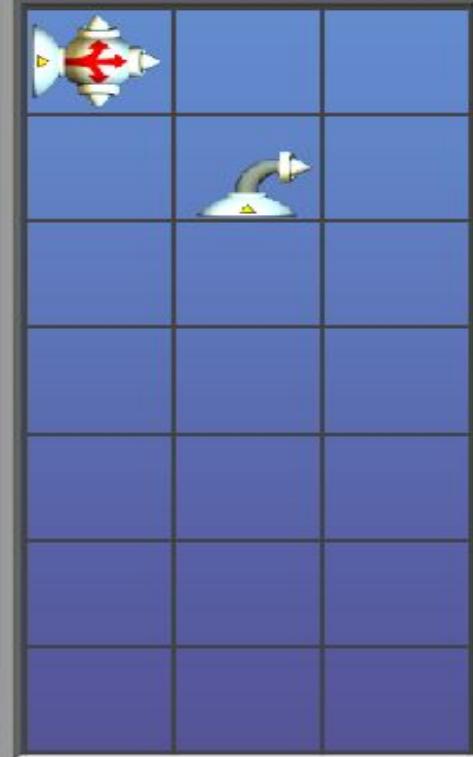
# Today

1. Imitation vs batch/offline RL policy learning
2. Fitted Q Iteration / Offline Q Learning
3. Pessimism
4. Case Study

# Is the Hope for Batch RL over Imitation Learning?



Level 1:8  
Fork

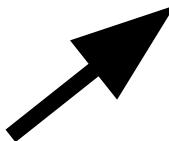


MENU

OPTIONS

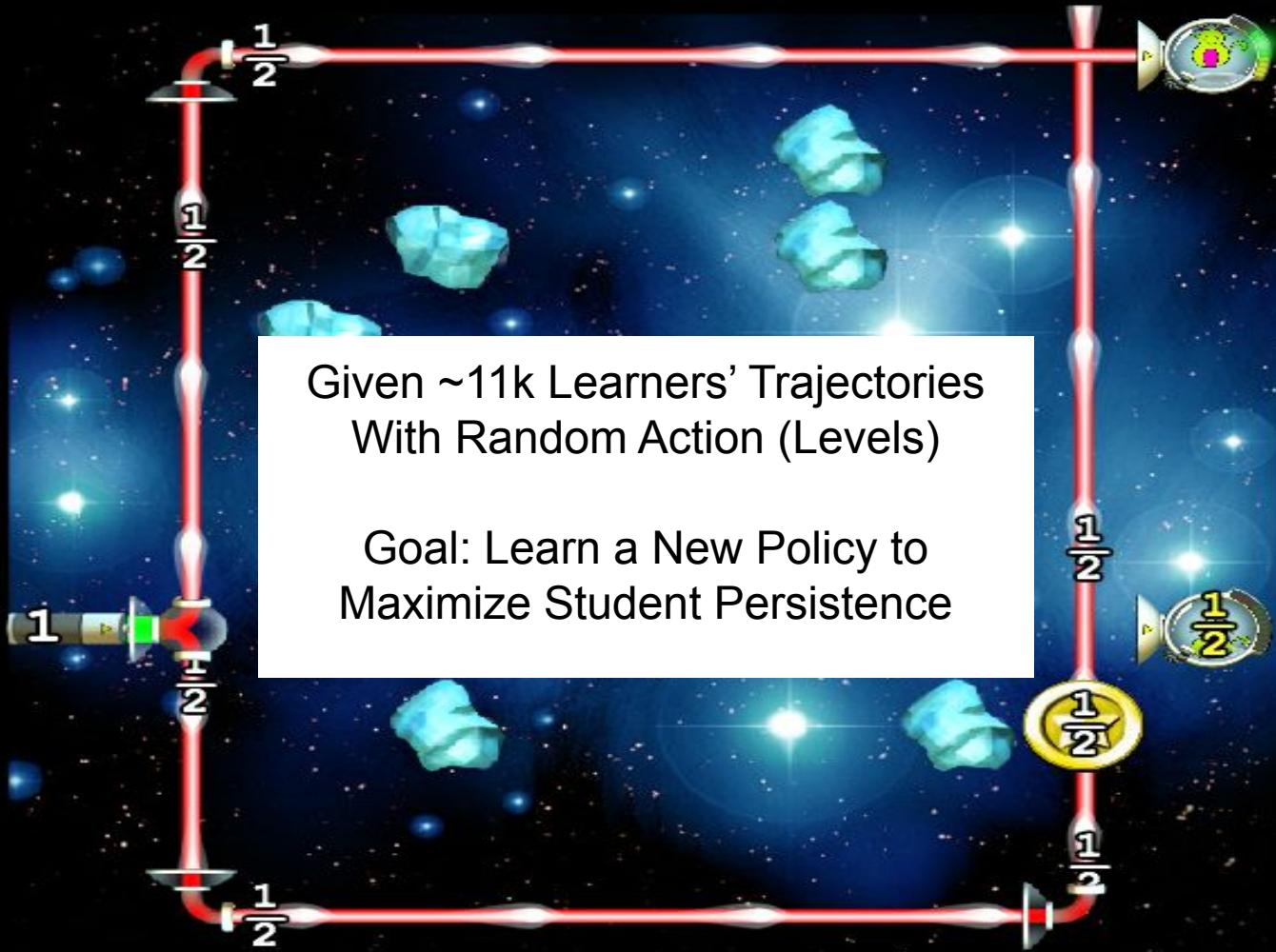


Took > 30s



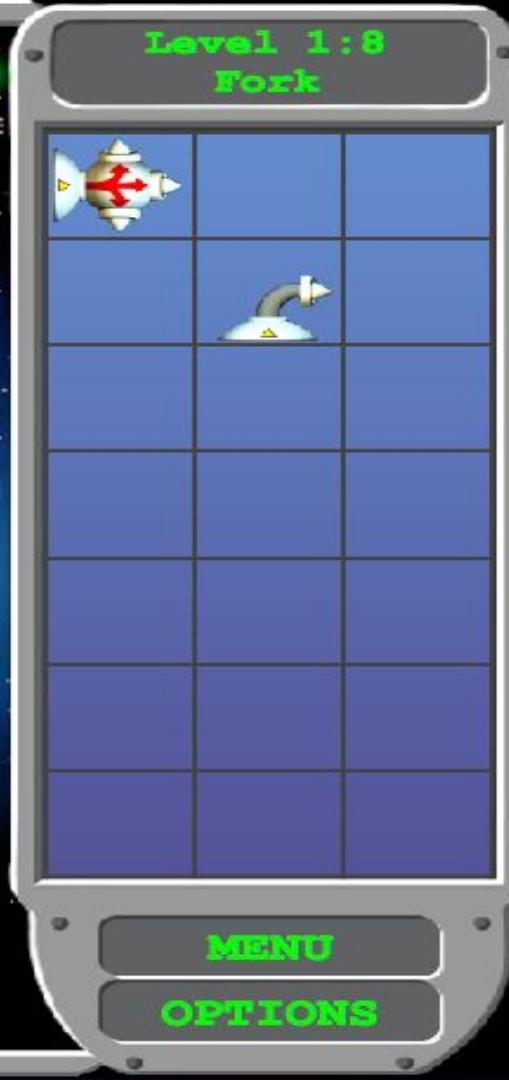
Took <= 30s





Given ~11k Learners' Trajectories  
With Random Action (Levels)

Goal: Learn a New Policy to  
Maximize Student Persistence







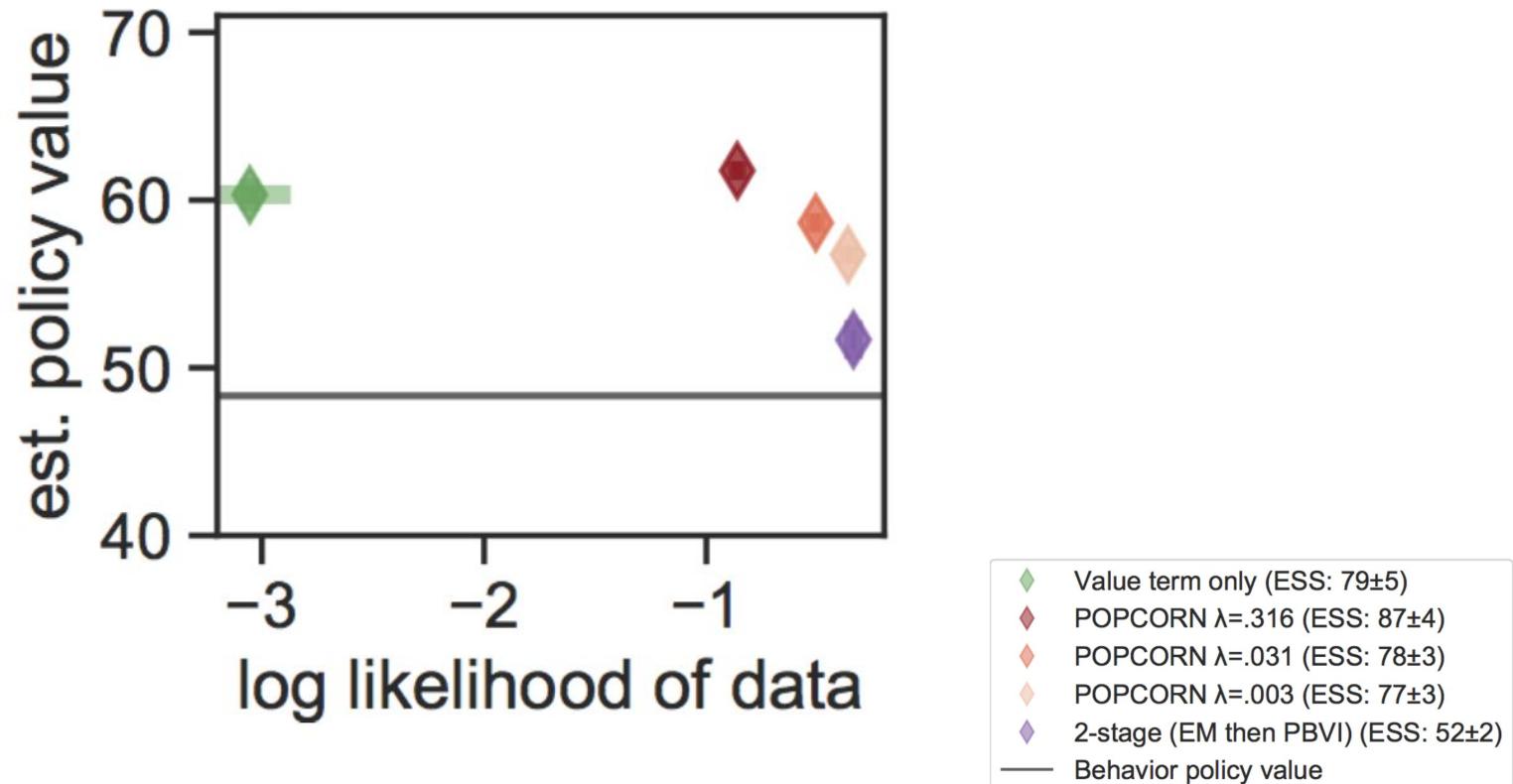
Given ~11k Learners' Trajectories  
With Random Action (Levels)

Learned a Policy that Increased  
Student Persistence by +30%

(Mandel, Liu, Brunskill, Popovic 2014)



# Encouraging Recent Work on Observational Health Data (MIMIC) Hypotension



# Today

1. Imitation vs batch/offline RL policy learning
2. Fitted Q Iteration / Offline Q Learning
3. Pessimism
4. Case study

# Offline / Batch Reinforcement Learning

Tasks

$$\int_{s \in S_0} \hat{V}^\pi(s, \mathcal{D}) ds$$

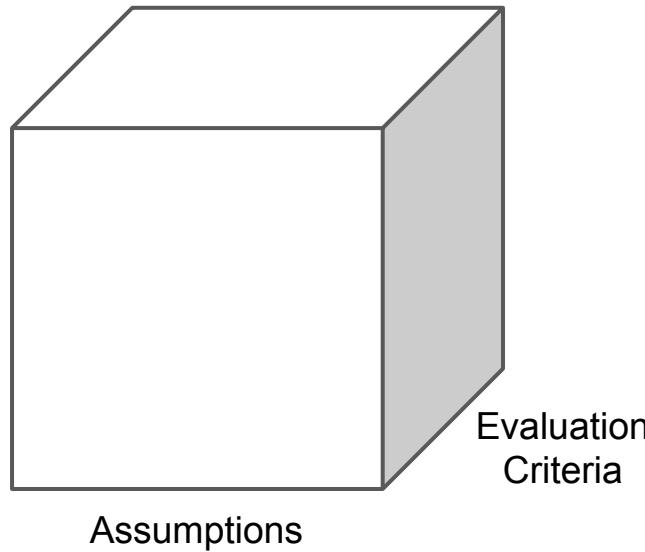
$$\arg \max_{\pi \in \mathcal{H}_i} \int_{s \in S_0} \hat{V}^\pi(s, \mathcal{D}) ds$$

$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

$\pi$ : Policy mapping  $s \rightarrow a$

$S_0$ : Set of initial states

$\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$



- Markov?
- Overlap?
- Sequential ignorability?

- Empirical accuracy
- Consistency
- Robustness
- Asymptotic efficiency
- Finite sample bounds
- Computational cost

# Batch Policy Optimization: Find a Good Policy That Will Perform Well in the Future

$$\underbrace{\arg \max_{\pi \in \mathcal{H}_i} \max_{\mathcal{H}_i \in \{\mathcal{H}_1, \mathcal{H}_2, \dots\}}}_{\text{Policy Optimization}} \quad \underbrace{\int_{s \in S_0} \hat{V}^\pi(s, \mathcal{D}) ds}_{\text{Policy Evaluation}}$$

$$\mathcal{H} = \mathcal{M}, \mathcal{V}, \Pi ?$$

- Today will not be a comprehensive overview, but instead highlight some of the challenges involved & some approaches with desirable statistical properties convergence, sample efficiency & bounds

$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$   
 $\pi$ : Policy mapping  $s \rightarrow a$   
 $S_0$ : Set of initial states  
 $\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$

*Lecine  
2020  
overview*

# Policy Optimization: Find Good Policy to Deploy

$$\arg \max_{\pi \in \mathcal{H}_i} \max_{\mathcal{H}_i \in \{\mathcal{H}_1, \mathcal{H}_2, \dots\}} \int_{s \in S_0} \hat{V}^\pi(s, \mathcal{D}) ds$$

$\mathcal{H} = \mathcal{M}, \mathcal{V}, \Pi ?$

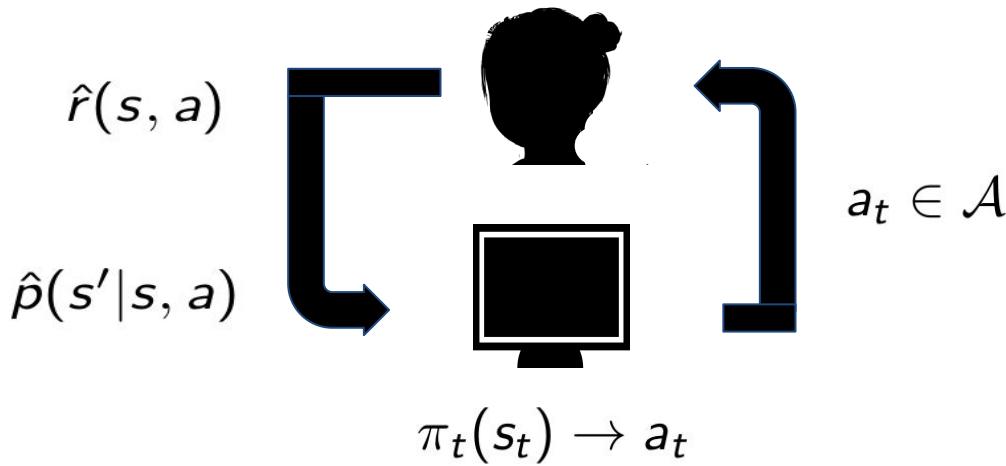
$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

$\pi$ : Policy mapping  $s \rightarrow a$

$S_0$ : Set of initial states

$\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$

# Learn Dynamics and Reward Models from Data, Plan



$$|\hat{V}^* - V^\pi|$$
  
$$\hat{V}^*(s) = \max_a \hat{r}(s, a) + \gamma \sum_{s'} \hat{p}(s'|s, a) \hat{V}^*(s')$$

$$\underline{\pi(s)} = \arg \max_a \hat{Q}^1(s, a)$$

# Model Free Value Function Approximation: Fitted Q Iteration

DQN

$$\mathcal{D} = (s_i, a_i, r_i, s_{i+1}) \quad \forall i$$

$$(\mathcal{T}f)(s, a) := R(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)}[V_f(s')]$$

$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

$\pi$ : Policy mapping  $s \rightarrow a$

$S_0$ : Set of initial states

$\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$

# Value Function Estimation, Fitted Q Iteration

**Theorem 2** (Sample complexity of FQI). *Given a dataset  $D = \{(s, a, r, s')\}$  with sample size  $|D| = n$  and  $\mathcal{F}$  that satisfies completeness (Assumption 3 when  $\mathcal{G} = \mathcal{F}$ ), w.p.  $\geq 1 - \delta$ , the output policy of FQI after  $k$  iterations,  $\pi_{f_k}$ , satisfies  $\underline{v^*} - \underline{v^{\pi_{f_k}}} \leq \epsilon \cdot V_{\max}$  when  $k \rightarrow \infty$  and<sup>11</sup>*

$$n = O\left(\frac{C \ln \frac{|\mathcal{F}|}{\delta}}{\epsilon^2 (1 - \gamma)^4}\right).$$

$$\forall f \in \mathcal{F}, \mathcal{T}f \in \mathcal{G}.$$

$$\forall v \quad \xrightarrow[\text{fix } f]{\text{doubling}} \quad \stackrel{\text{sa}}{\vee} \quad Q^* \in \mathcal{F}$$

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \frac{\nu(s, a)}{\mu(s, a)} \leq C.$$

*↑  
doubling  
policy*

# Value Function Estimation, Fitted Q Iteration

**Theorem 2** (Sample complexity of FQI). *Given a dataset  $D = \{(s, a, r, s')\}$  with sample size  $|D| = n$  and  $\mathcal{F}$  that satisfies completeness (Assumption 3 when  $\mathcal{G} = \mathcal{F}$ ), w.p.  $\geq 1 - \delta$ , the output policy of FQI after  $k$  iterations,  $\pi_{f_k}$ , satisfies  $v^* - v^{\pi_{f_k}} \leq \epsilon \cdot V_{\max}$  when  $k \rightarrow \infty$  and<sup>11</sup>*

$$n = O\left(\frac{C \ln(\frac{|\mathcal{F}|}{\delta})}{\epsilon^2(1-\gamma)^4}\right).$$

Munos  
2003  
2008, 2010

Bellman  
Backup

$\forall f \in \mathcal{F}, T_f \in \mathcal{G}$ .

$B_f$   
Completeness

$Q^* \in \mathcal{F}$   
Realizability

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \frac{\nu(s, a)}{\mu(s, a)} \leq C.$$

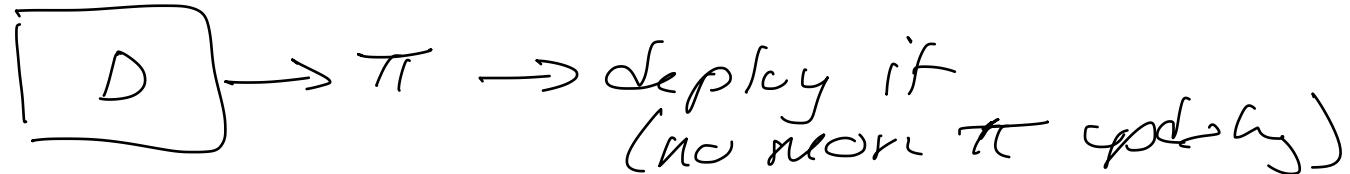
Overlap assumption: Concentrability coefficient

# Today

1. Imitation vs batch/offline RL policy learning
2. Fitted Q Iteration / Offline Q Learning
- 3. Pessimism**
4. Case Study

# Check Your Intuition

- Optimism under uncertainty can enable sublinear regret in online multi-armed bandits  $\checkmark$
- Pessimism under uncertainty can lead to linear regret in online multi-armed bandits  $\times$
- With high probability the optimistic upper bound on the selected arm in UCB algorithms is an upper bound on the performance of any arm  $\checkmark$
- In offline / batch RL selecting the optimistic best arm is likely to be best  $\times$
- In offline / batch RL selecting the arm with the highest mean is likely to be best  $\times$
- Not sure



robust MDP 2003-2005  
param uncertainty in MDPs 1990s

# Check Your Intuition Solutions

- Optimism under uncertainty can enable sublinear regret in online multi-armed bandits
- Pessimism under uncertainty can lead to linear regret in online multi-armed bandits
- With high probability the optimistic upper bound on the selected arm in UCB algorithms is an upper bound on the performance of any arm
- In offline / batch RL selecting the optimistic best arm is likely to be best
- In offline / batch RL selecting the arm with the highest mean is likely to be best
- Not sure

# Offline / Batch Reinforcement Learning

Tasks

$$\int_{s \in S_0} \hat{V}^\pi(s, \mathcal{D}) ds$$

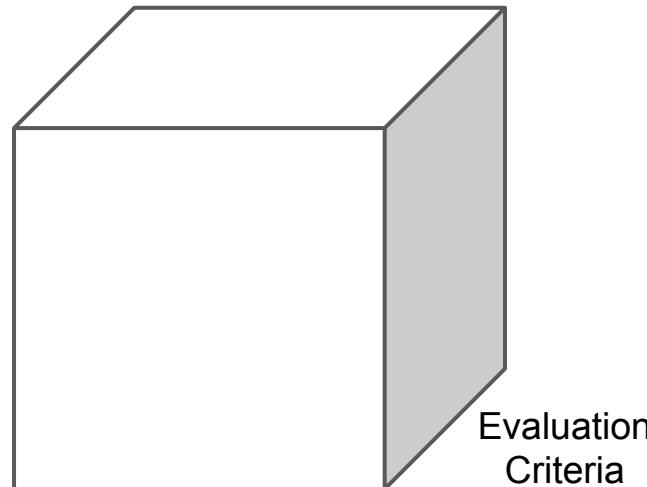
$$\arg \max_{\pi \in \mathcal{H}_i} \int_{s \in S_0} \hat{V}^\pi(s, \mathcal{D}) ds$$

$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

$\pi$ : Policy mapping  $s \rightarrow a$

$S_0$ : Set of initial states

$\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$



## Assumptions

- Markov?
- Overlap?
- Sequential ignorability?

- Empirical accuracy
- Consistency
- Robustness
- Asymptotic efficiency
- Finite sample bounds
- Computational cost
- Constraints?

# Standard Assumptions for Off Policy / Counterfactual Estimation & Optimization

- Overlap
  - Have to take all actions that target policy would take
  - In infinite data / finite data
- No confounding

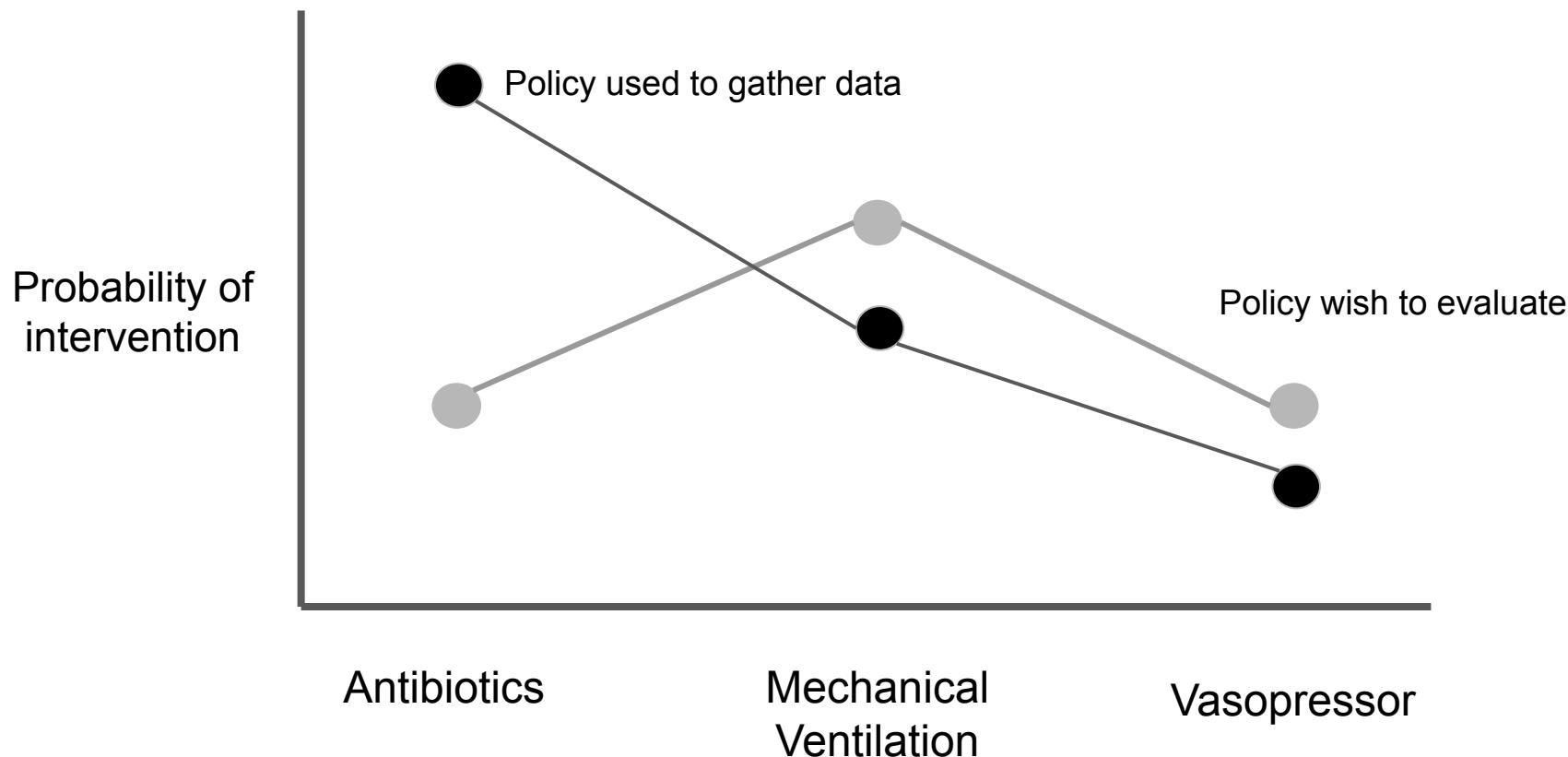
$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

$\pi$ : Policy mapping  $s \rightarrow a$

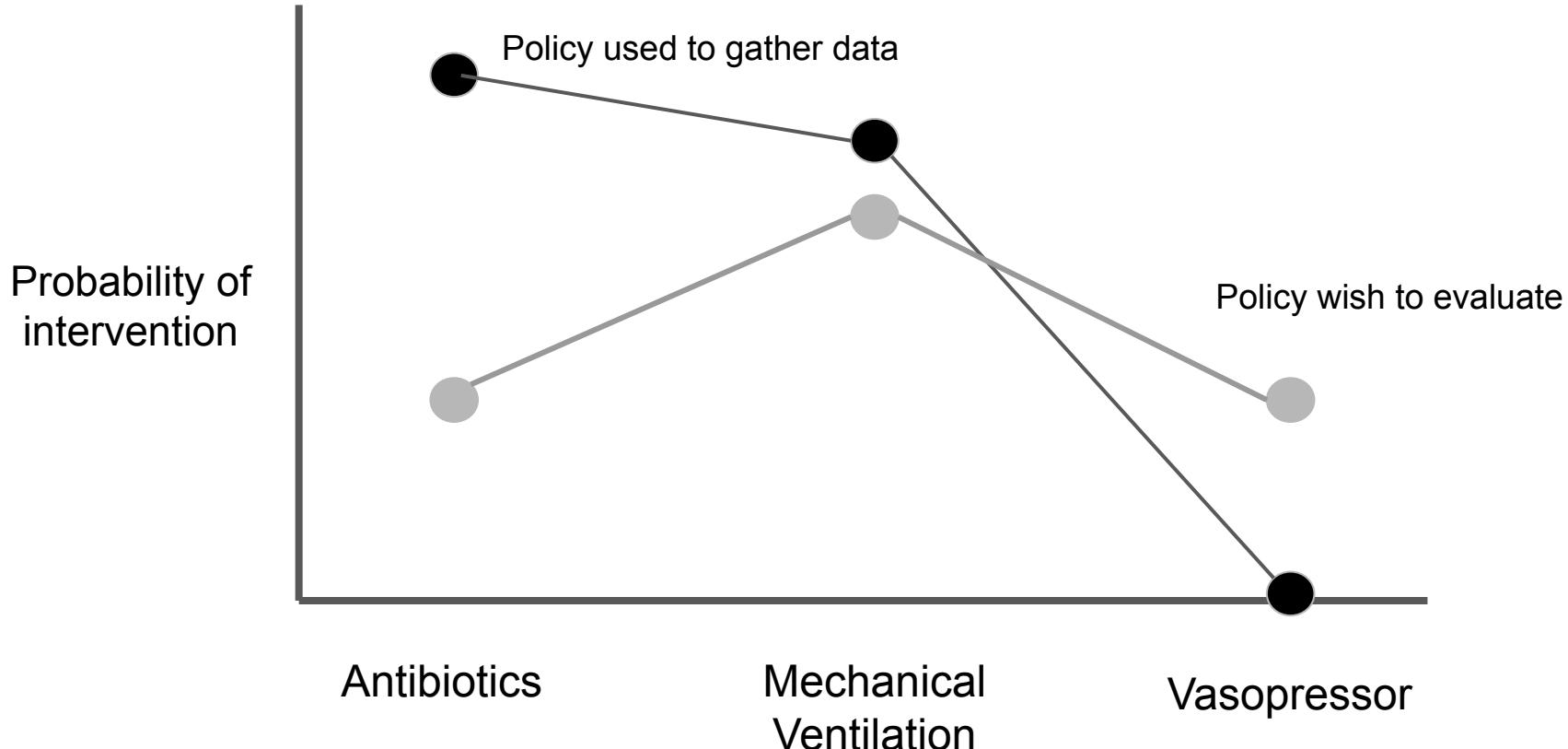
$S_0$ : Set of initial states

$\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$

# Overlap Requirement: Data Must Support Policy Wish to Evaluate



# No Overlap for Vasopressor $\Rightarrow$ Can't Do Off Policy Estimation for Desired Policy



# Limitations of Prior Work

- Typically assume overlap
  - Off policy estimation: for policy of interest
  - Off policy optimization: for all policies including optimal one  
(see concentrability assumption in batch RL)
- Unlikely to be true in many settings
- Many real datasets don't include complete random exploration

# Limitations of Prior Work

- Typically assume overlap
  - Off policy estimation: for policy of interest
  - Off policy optimization: for all policies including optimal one (see concentrability assumption in batch RL)
- Unlikely to be true in many settings
- Many real datasets don't include complete random exploration
- Assuming overlap when it's not there can be a problem:
  - We can end up with a policy with estimated high performance, but actually does poorly when deployed

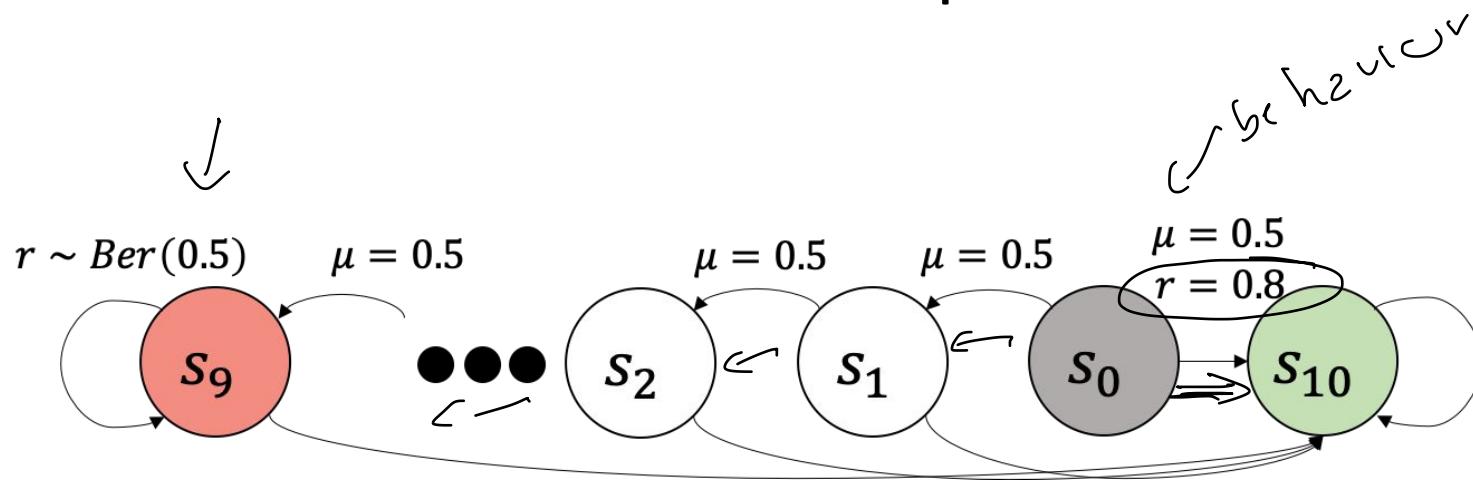
Pessimistic Batch RL 2019-now

## Doing the Best with What We've Got: Off Policy Optimization Without Full Data Coverage

- Idea: restrict off policy optimization to those with overlap in data
- Computationally tractable algorithm
- Simple idea: assume **pessimistic outcomes** for areas of state--action space with insufficient overlap/support

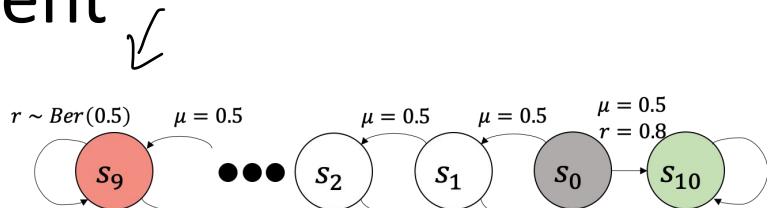
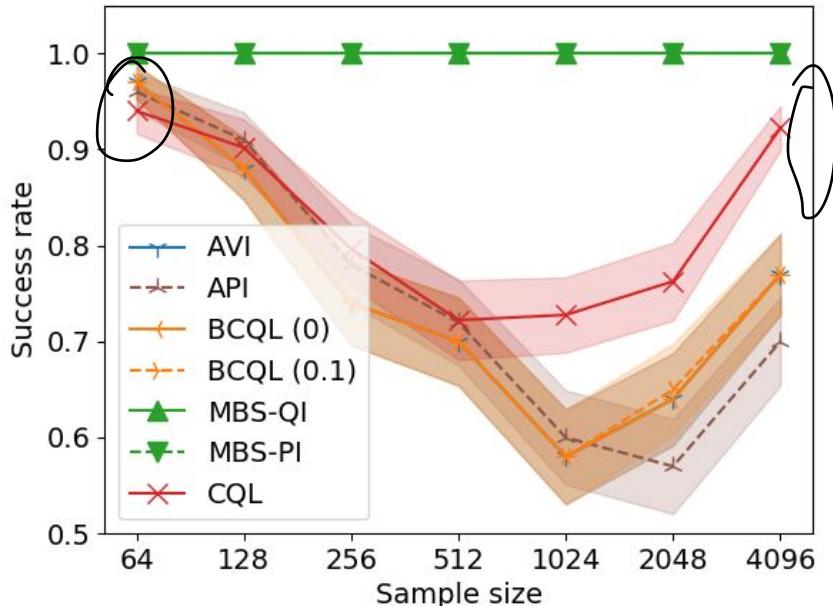
*Common challenge that's attracted substantial interest in last few years but...*

# Illustrative Examples



$$H = 10$$

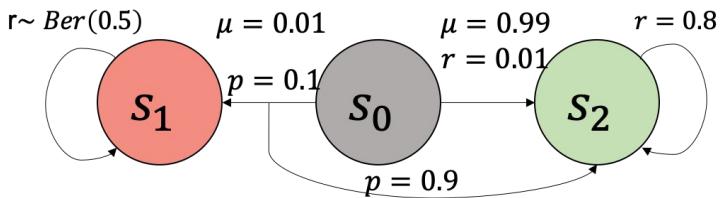
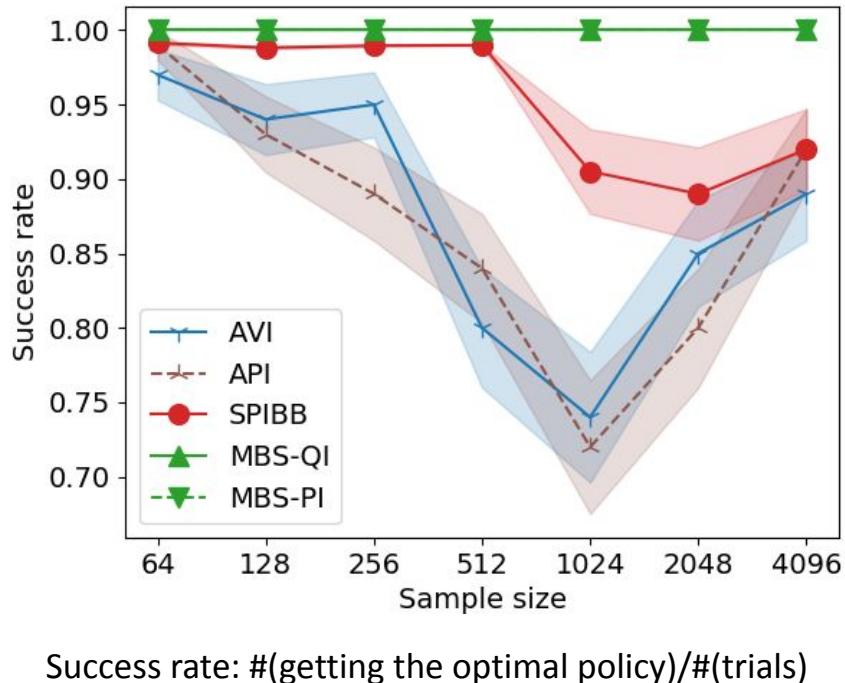
# Recent Conservative Batch Reinforcement Learning Are Insufficient



Reasons why baselines fail:

- Many baselines focus on penalty/constraints that are based on  $\text{dist}(\pi(a|s), \pi_b(a|s))$ .
- In this example a sequence of large action conditional probabilities leads to a rare state.
- Due to finite samples, estimates of the reward of this rare state can be overestimated.

# Recent Conservative Batch Reinforcement Learning Are Insufficient



Reasons why baselines fail:

- SPIBB adds conservatism based on estimates of  $\pi_b$  &  $V$  of  $\pi_b$ .
- In this example, the actions which are rare under  $\pi_b$  also have a stochastic transition and reward, thus the  $\pi_b$ 's  $V$  is overestimated.

Idea: Use pessimistic value for state-action space with insufficient data

- Filtration function:

$$\zeta(s, a; \hat{\mu}, b) = 1(\hat{\mu}(s, a) > b)$$

↓  
dansing  
s-a bahuton  
dar

# Idea: Use pessimistic value for state-action space with insufficient data

- Filtration function:

$$\zeta(s, a; \hat{\mu}, b) = 1(\hat{\mu}(s, a) > b)$$

**b can account for statistical uncertainty due to finite samples**

# Idea: Use pessimistic value for state-action space with insufficient data

- Filtration function:

$$\zeta(s, a; \hat{\mu}, b) = 1(\hat{\mu}(s, a) > b)$$

- Bellman operator and Bellman evaluation operator:

$$\mathcal{T}f(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} \left[ \max_{a'} \underbrace{\zeta(s', a') f(s', a')}_{\text{V}} \right]$$

# Idea: Use pessimistic value for state-action space with insufficient data

- Filtration function:

$$\zeta(s, a; \hat{\mu}, b) = 1(\hat{\mu}(s, a) > b)$$

- Bellman operator and Bellman evaluation operator:

$$Tf(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} \left[ \max_{a'} \zeta(s', a') f(s', a') \right]$$

$\Rightarrow = 0$  for  $(s', a')$  with insufficient data.

We assume  $r(s, a) \geq 0$

Therefore pessimistic estimate for such tuples

# Idea: Use pessimistic value for state-action space with insufficient data

- Filtration function:

$$\zeta(s, a; \hat{\mu}, b) = 1(\hat{\mu}(s, a) > b)$$

- Bellman operator and Bellman evaluation operator:

$$\mathcal{T}f(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} \left[ \max_{a'} \zeta(s', a') f(s', a') \right]$$

$$\mathcal{T}^\pi f(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P, a' \sim \pi} [\zeta(s', a') f(s', a')]$$

# Marginalized Behavior Supported (MBI) Policy Optimization

- Filtration function:

$$\zeta(s, a; \hat{\mu}, b) = 1(\hat{\mu}(s, a) > b)$$

- Bellman operator and Bellman evaluation operator:

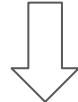
$$\mathcal{T}f(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} \left[ \max_{a'} \zeta(s', a') f(s', a') \right]$$

$$\mathcal{T}^\pi f(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P, a' \sim \pi} [\zeta(s', a') f(s', a')]$$

# Majority of Past Model-Free Batch RL Theory for Function Approximation Setting

**Assume** for any  $\nu(s,a)$  distribution possible  
under some policy in this MDP

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \frac{\nu(s, a)}{\mu(s, a)} \leq C.$$

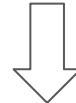


$$V^* - V^{\pi_A} \leq \epsilon$$

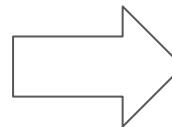
# Best in Well Supported Policy Class\*

**Assume** for any  $v(s,a)$  distribution possible under some policy in this MDP

$$\forall (s,a) \in \mathcal{S} \times \mathcal{A}, \frac{v(s,a)}{\mu(s,a)} \leq C.$$



$$V^* - V^{\pi_A} \leq \epsilon$$



**Define**

$$\Pi_{all} : \underline{\pi} \text{ s.t.}$$

$$\mathbb{E}_{s,a \sim \eta^\pi} [\mathbb{1} (\zeta(s,a) = 0)] \leq \epsilon_\zeta$$



$$\max_{\pi' \in \Pi_{all}} V^{\pi'} - V^{\pi_A} \leq \epsilon$$

\*Note: Policy set  $\Pi_{all}$  is not constructed, but implicitly our algorithm only considers elements in it

**Assumption 1** (Bounded densities). *For any non-stationary policy  $\pi$  and  $h \geq 0$ ,  $\eta_h^\pi(s, a) \leq U$ .*

**Assumption 2** (Density estimation error). *With probability at least  $1 - \delta$ ,  $\|\hat{\mu} - \mu\|_{TV} \leq \epsilon_\mu$ .*

**Assumption 3** (Completeness under  $\tilde{\mathcal{T}}^\pi$ ).  $\forall \pi \in \Pi$ ,  $\max_{f \in \mathcal{F}} \min_{g \in \mathcal{F}} \|g - \tilde{\mathcal{T}}^\pi f\|_{2,\mu}^2 \leq \epsilon_{\mathcal{F}}$ .

**Assumption 4** ( $\Pi$  Completeness).  $\forall f \in \mathcal{F}$ ,  $\min_{\pi \in \Pi} \|\mathbb{E}_\pi [\zeta \circ f(s, a)] - \max_a \zeta \circ f(s, a)\|_{1,\mu} \leq \epsilon_\Pi$ .

$$\boxed{\begin{aligned}\eta_h^\pi(s) &:= \Pr[s_h = s | \pi], \\ \eta_h^\pi(s, a) &= \eta_h^\pi(s) \pi(a | s)\end{aligned}}$$

$$\zeta(s, a; \hat{\mu}, b) = \mathbb{1}(\hat{\mu}(s, a) \geq b)$$

# Theoretical Result

We bound the error w.r.t. the best policy in the following policy set:

$$\{\text{all policies such that } \Pr(\zeta(s, a) = 0 | \pi) \leq \epsilon_\zeta\}$$

Error bounds<sup>1</sup>:

- PI:

$$O\left(\frac{V_{\max}}{(1-\gamma)^2 b} \sqrt{\frac{\ln(|\mathcal{F}| |\Pi| / \delta)}{n}}\right) + \frac{V_{\max} \epsilon_\zeta}{1-\gamma}$$

*✓  $\zeta$   $\mathcal{F}$   $n$*

- VI<sup>2</sup>:

$$O\left(\frac{V_{\max}}{(1-\gamma)^2 b} \sqrt{\frac{\ln(|\mathcal{F}| / \delta)}{n}}\right) + \frac{V_{\max} \epsilon_\zeta}{1-\gamma}$$

1: We omit some constant terms that is same as standard ADP analysis with function approximation.

2: For VI results there is another important constant term, see our paper for detailed result and discussion.

$$\zeta(s, a; \hat{\mu}, b) = \mathbb{1}(\hat{\mu}(s, a) \geq b)$$

# Theoretical Result

We bound the error w.r.t. the best policy in the following policy set:

$$\{\text{all policies such that } \Pr(\zeta(s, a) = 0 | \pi) \leq \epsilon_\zeta\}$$

Error bounds<sup>1</sup>:

- PI:

$$O\left(\frac{V_{\max}}{(1-\gamma)^3 b} \sqrt{\frac{\ln(|\mathcal{F}| |\Pi| / \delta)}{n}}\right) + \frac{V_{\max} \epsilon_\zeta}{1-\gamma}$$

- VI<sup>2</sup>:

$$O\left(\frac{V_{\max}}{(1-\gamma)^2 b} \sqrt{\frac{\ln(|\mathcal{F}| / \delta)}{n}}\right) + \frac{V_{\max} \epsilon_\zeta}{1-\gamma}$$

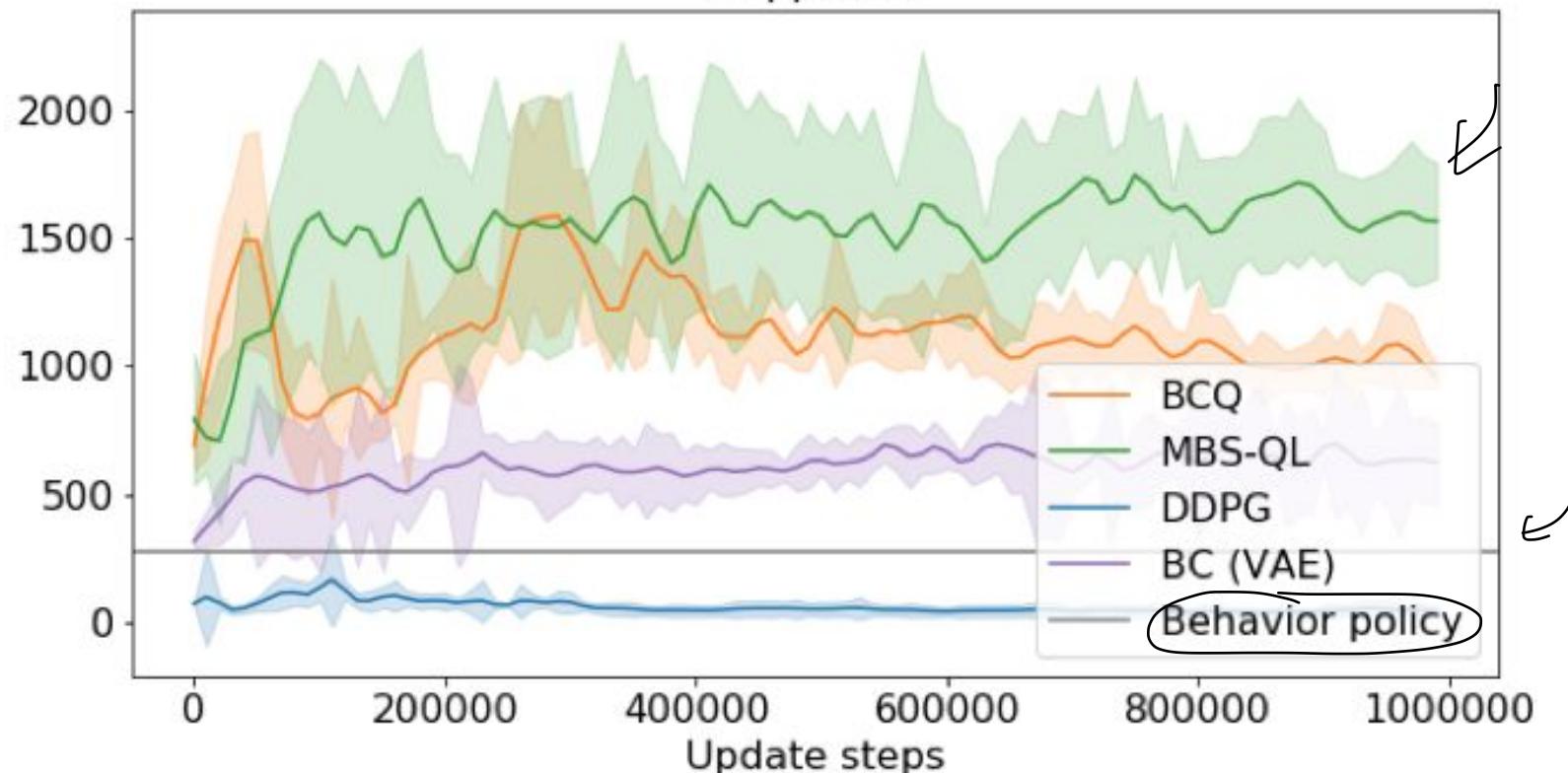
Note: Results are for function approximation, finite sample setting

1: We omit some constant terms that is same as standard ADP analysis with function approximation.

2: For VI results there is another important constant term, see our paper for detailed result and discussion.

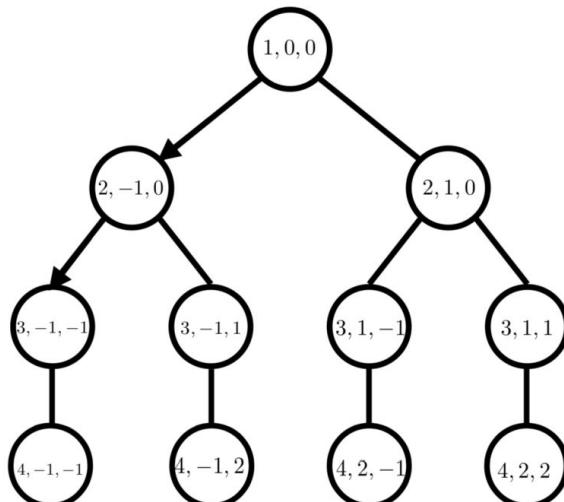
# Can Do Get Substantially Better Solutions, With Same Data

Hopper-v3



# This Was Model Free. Might Models Be Even Better?

- Model based approaches can be provably more efficient than model free value function for *online* evaluation or control



Sun, Jiang, Krishnamurthy,  
Agarwal, Langford COLT 2019

$$x_{t+1} = A_\star x_t + B_\star u_t + w_t ,$$

$$V^K(x) := \lim_{T \rightarrow \infty} \mathbb{E} \left[ \sum_{t=0}^{T-1} (x_t^\top Q x_t + u_t^\top R u_t - \lambda_K) \mid x_0 = x \right]$$

Tu & Recht COLT 2019

# Concurrent Work on Conservative Model-Based Offline Batch Reinforcement Learning

- Ex. Yu, Thomas, Yu, Ermon, Zou, Levine, Finn & Ma (NeurIPS 2020) and Kidambi, Rajeswaran, Netrapalli & Joachims (NeurIPS 2020)
- Learn a model and penalize model uncertainty during planning
- Empirically very promising on D4RL tasks
- Their work has more limited theoretical analysis

$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

$\pi$ : Policy mapping  $s \rightarrow a$

$S_0$ : Set of initial states

$\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$

# Early Comparison with Concurrent Work

	<b>MBS-BCQ</b>	<b>MBS-BEAR</b>	BCQ	BEAR	MOPO	CQL
Hopper-medium	75.9	32.3	54.5	52.1	26.5	58.0
HalfCheetah-medium	38.4	39.7	40.7	41.7	40.2	44.4
Walker2d-medium	64.4	75.4	53.1	59.1	14.0	79.2

- Preliminary draft results: on some D4RL recent model-based pessimistic approaches or CQL do better
- In general suspect recent model-based approaches will dominate our MBS empirically but our theoretical results are stronger
- Interesting to see further theoretical work on model based approaches

# Pessimistic Model-Free Batch/Offline Policy Learning

- Restrict off policy optimization to those with overlap in data
- Computationally tractable algorithm
- **Simple idea: assume pessimistic outcomes for areas of state--action space with insufficient overlap/support**
- Theoretical results bound distance to best supported policy
  - Considers finite sample & function approximation
- Model free value function method

⇒ ***Pessimism under uncertainty has received a lot of attention in last 1-2 years for offline RL***

# Today

1. Imitation vs batch/offline RL policy learning
2. Fitted Q Iteration / Offline Q Learning
3. Pessimism
4. **Case Study**

RESEARCH

---

COMPUTER SCIENCE

# Preventing undesirable behavior of intelligent machines

Philip S. Thomas<sup>1\*</sup>, Bruno Castro da Silva<sup>2</sup>, Andrew G. Barto<sup>1</sup>, Stephen Giguere<sup>1</sup>,  
Yuriy Brun<sup>1</sup>, Emma Brunskill<sup>3</sup>

bioRxiv preprint doi: <https://doi.org/10.1101/2019.10.22.853400>; this version posted October 22, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [aCC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

Science November 2019

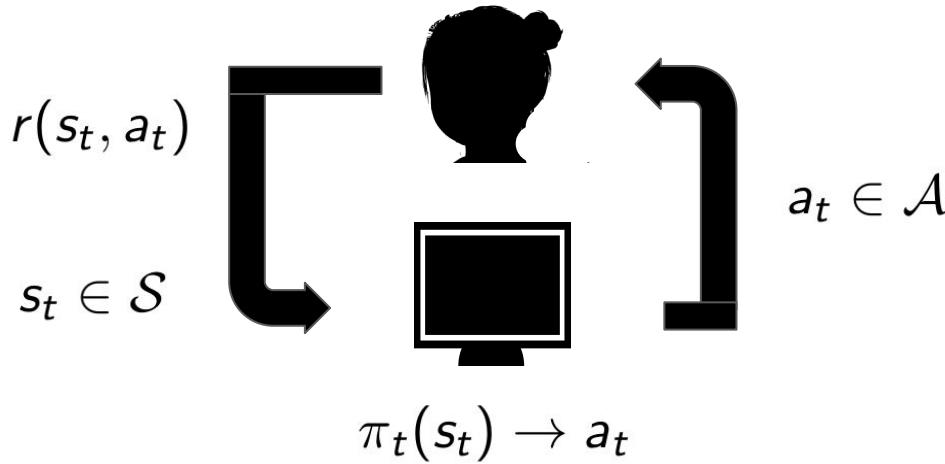


# Optimizing while Ensuring Solution Won't, in the Future, Exhibit Undesirable Behavior

$$\begin{aligned} & \arg \max_{a \in \mathcal{A}} f(a) \\ \text{s.t. } & \text{s.t. } \forall i \in \{1, \dots, n\}, \Pr\left(g_i(a(D)) \leq 0\right) \geq 1 - \delta_i \end{aligned}$$

↓  
Constraints

# Counterfactual RL with Constraints on Future Performance of Policy



$\mathcal{D}$ : Dataset of  $n$  traj.s  $\tau$ ,  $\tau \sim \pi_b$

# Related Work in Decision Making

$$\arg \max_{a \in \mathcal{A}} f(a)$$

$$\text{s.t. } \forall i \in \{1, \dots, n\}, \Pr(g_i(a(D)) \leq 0) \geq 1 - \delta_i$$

- Chance constraints, data driven robust optimization have similar aims
- Most of this work has focused on ensuring computational efficiency for  $f$  and/or constraints  $g$  with certain structure (e.g. convex)
- Also need to be able to capture broader set of aims & constraints

# Batch RL with Safety Constraints

$$g(\theta) = \mathbf{E}[r'(H)|\theta_0] - \mathbf{E}[r'(H)|\theta]$$

The equation  $g(\theta) = \mathbf{E}[r'(H)|\theta_0] - \mathbf{E}[r'(H)|\theta]$  is displayed. Two arrows point upwards from the terms  $\mathbf{E}[r'(H)|\theta_0]$  and  $\mathbf{E}[r'(H)|\theta]$  to the labels "Default policy" and "Potential policy" respectively.

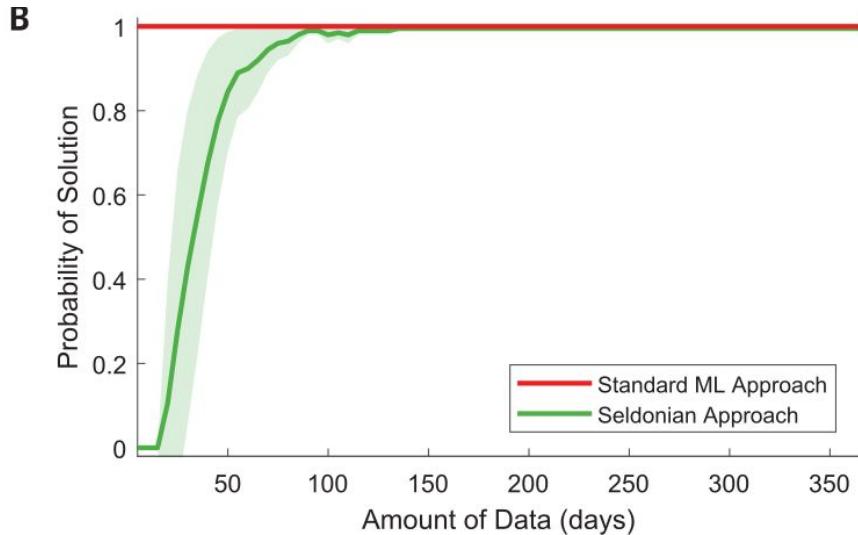
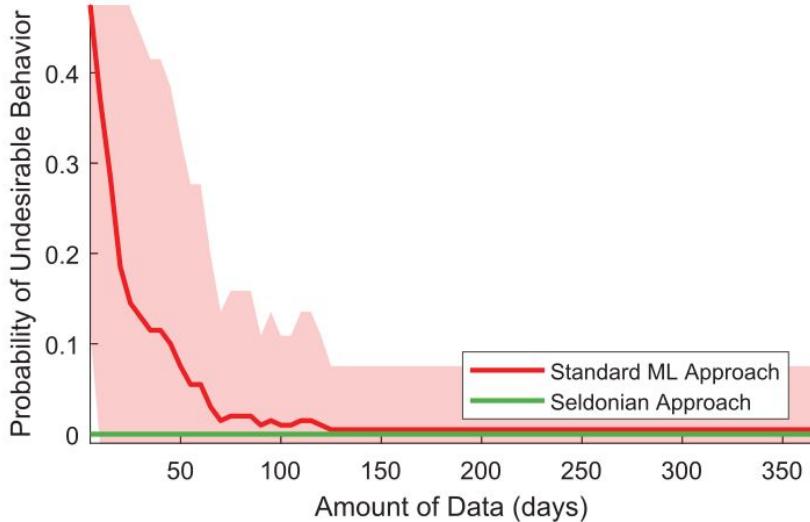
- $r'(H)$  is a function of the trajectory  $H$

# Diabetes Insulin Management

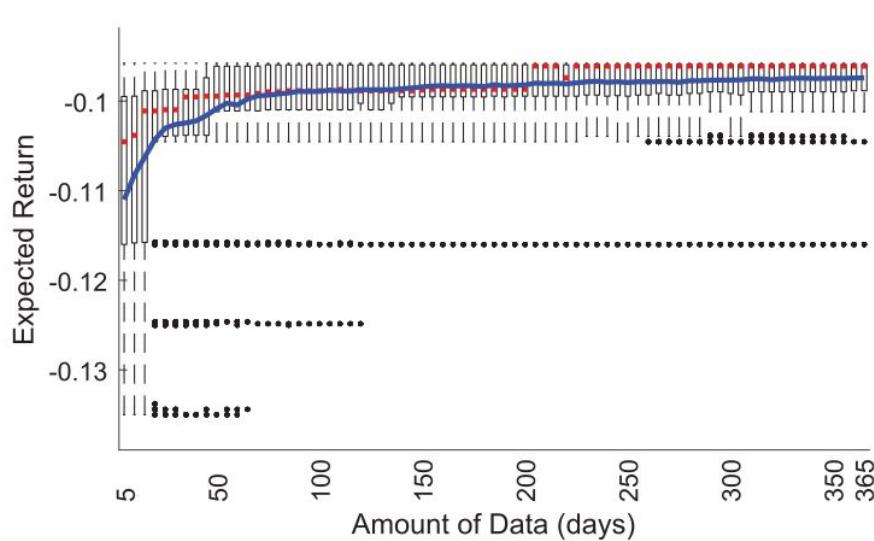


- Blood glucose control
- Action: insulin dosage
- Search over policies
- Constraint:  
hypoglycemia
- Very accurate simulator:  
approved by FDA to  
replace early stage  
animal trials

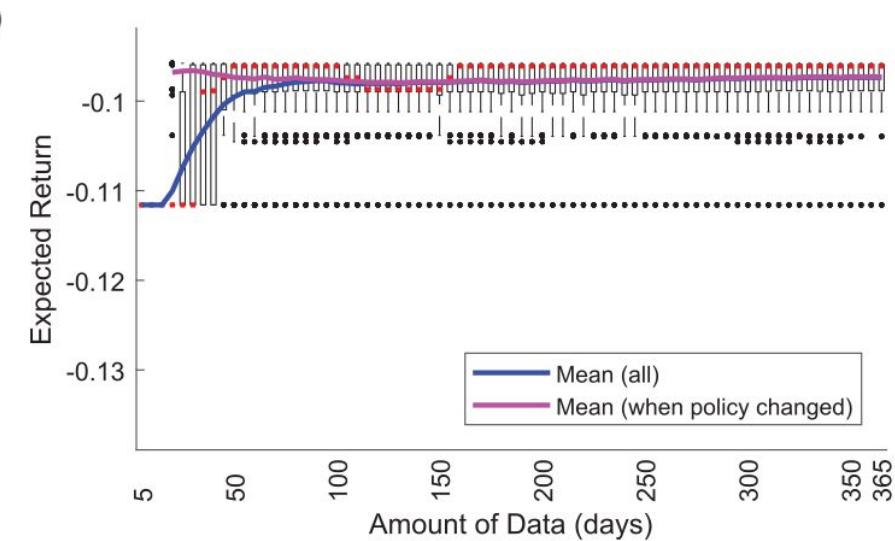
# Personalized Insulin Dosage: Safe Batch Policy Improvement



# Personalized Insulin Dosage: Quickly Can Have Confidence in Safe Better Policy



Standard RL



Our Safe Batch RL

# Optimizing while Ensuring Solution Won't, in the Future, Exhibit Undesirable Behavior

$$\begin{aligned} & \arg \max_{a \in \mathcal{A}} f(a) \\ \text{s.t. } & \text{s.t. } \forall i \in \{1, \dots, n\}, \Pr\left(g_i(a(D)) \leq 0\right) \geq 1 - \delta_i \end{aligned}$$

↓  
Constraints

⇒ Illustrated we can do this, for very general constraints, for several problems but many open questions around computational efficiency, other constraints ...

# What You Should Know

- Offline RL can do better than imitation learning / behavior cloning (Why?)
- Pessimism under uncertainty can be useful, particularly for high stakes applications
- Be able to give example application areas where offline RL might be useful

# Where We Are In The Course

1. Learning from offline data
  - a. Imitation learning
  - b. Batch/offline policy evaluation
  - c. Batch/offline policy learning
2. Next week
  - a. Guest lecture: Maria Dimakopoulou
  - b. Quiz