

# 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 ✓
  - Not Sure
- F dagger does

# Refresh Your Understanding

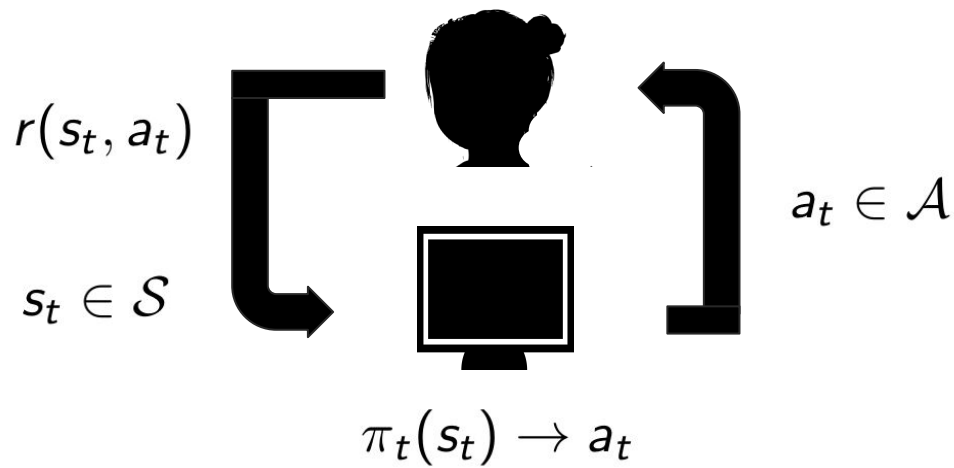
Importance sampling (select all that are true)

- 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  $\tau, \tau \sim \pi_b$

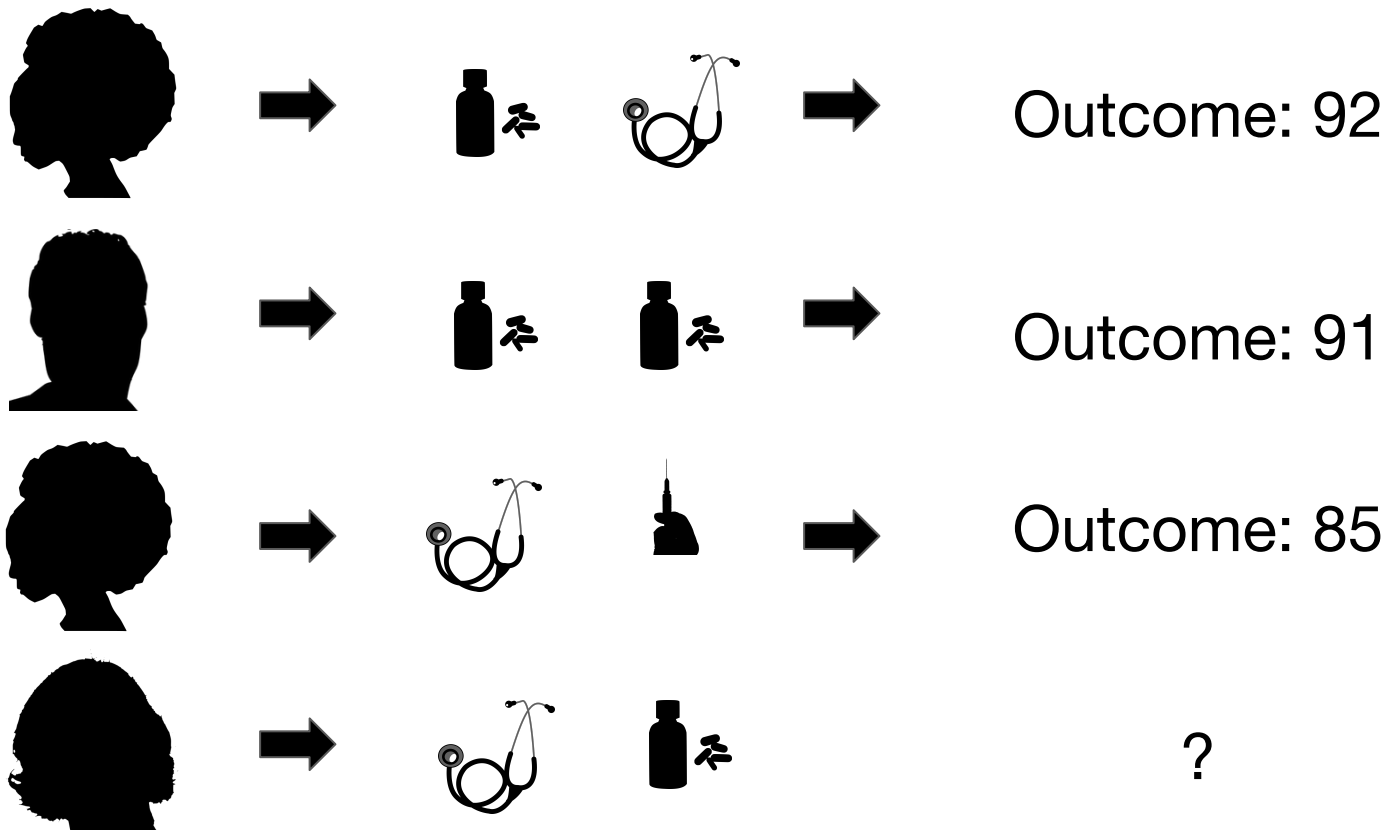
# Where We Are In The Course

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?





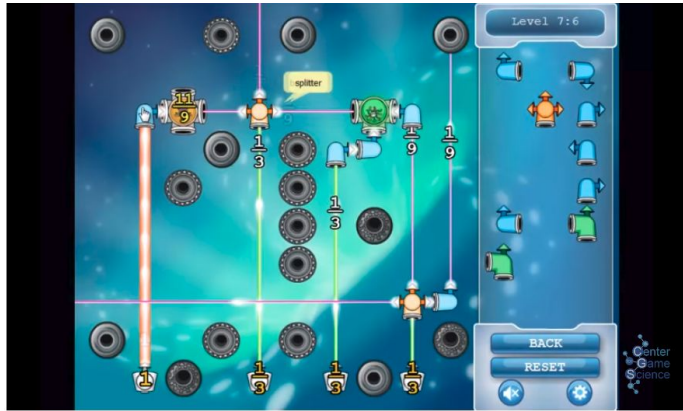
Level 1:8  
Fork

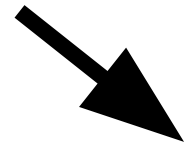
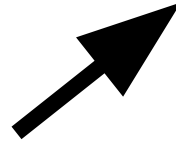
MENU

OPTIONS





Took > 30s



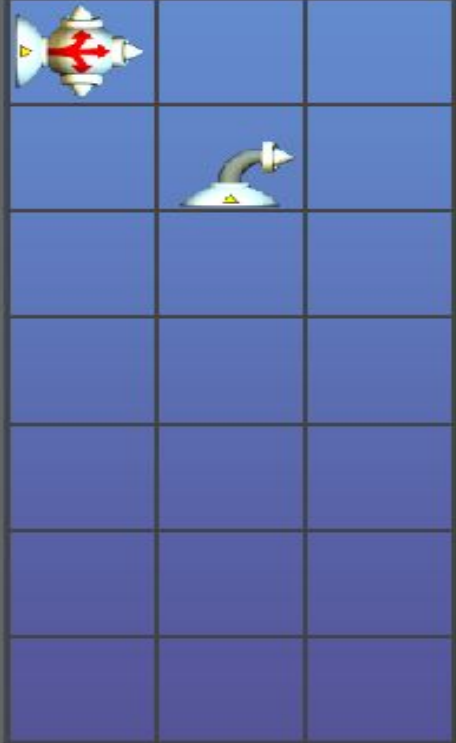
Took <= 30s



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

Goal: Learn a New Policy to  
Maximize Student Persistence

Level 1:8  
Fork



MENU

OPTIONS

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

Learn a Policy that Increases  
Student Persistence

(Mandel, Liu, Brunskill, Popovic 2014)

Level 1:8  
Fork



MENU

OPTIONS

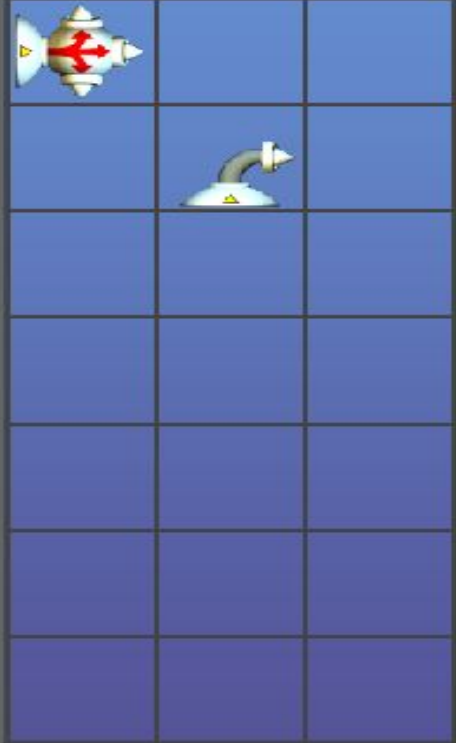


Given ~11k Learners' Trajectories  
With Random Action (Levels)

**Learned a Policy that Increased  
Student Persistence by +30%**

(Mandel, Liu, Brunskill, Popovic 2014)

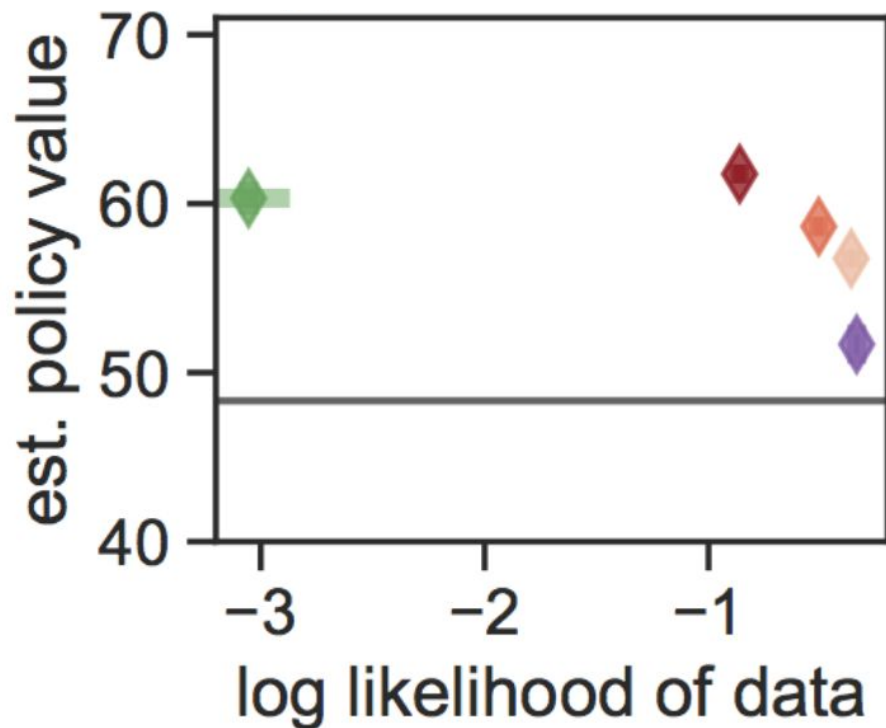
Level 1:8  
Fork



MENU

OPTIONS

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



- ◆ Value term only (ESS: 79±5)
- ◆ POPCORN  $\lambda=.316$  (ESS: 87±4)
- ◆ POPCORN  $\lambda=.031$  (ESS: 78±3)
- ◆ POPCORN  $\lambda=.003$  (ESS: 77±3)
- ◆ 2-stage (EM then PBVI) (ESS: 52±2)
- Behavior policy value

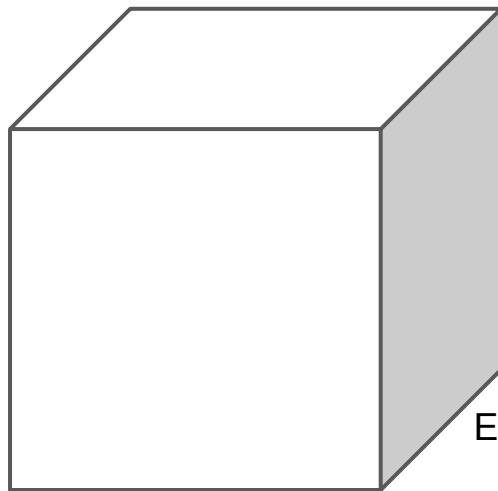
# Today

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

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



Assumptions

Evaluation  
Criteria

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

- Markov?
- Overlap?
- Sequential ignorability?

$\mathcal{D}$ : 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  $\mathcal{D}$

# 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  $\mathcal{D}$

Levine  
2020  
overview



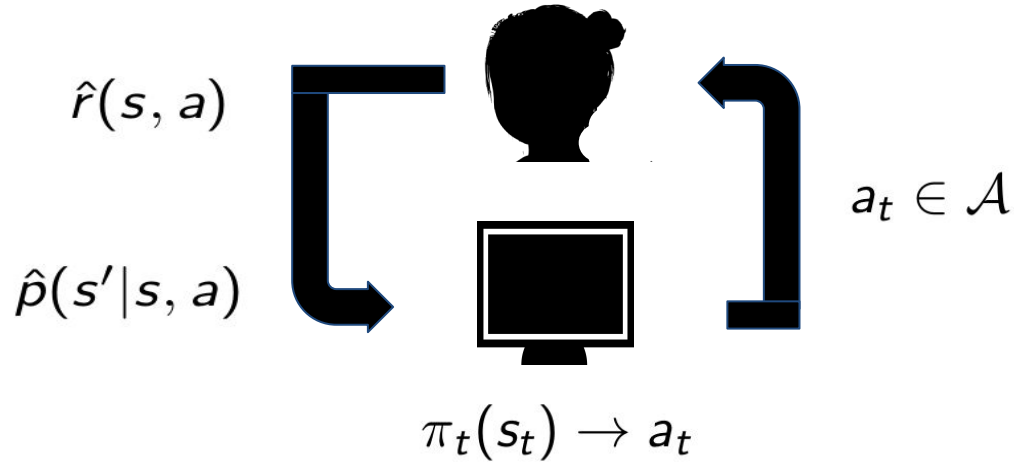
# Policy Optimization: Find Good Policy to Deploy

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

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

$\mathcal{D}$ : 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



$$|V^* - V^\pi|$$

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

$$|\hat{V}^* - V^\pi|$$

$$\pi(s) = \underset{a}{\operatorname{argmax}} Q^*(s, a)$$

# Model Free Value Function Approximation: Fitted Q Iteration

DQN

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

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

$\mathcal{D}$ : 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  $v^* - v^{\pi_{f_k}} \leq \epsilon \cdot V_{\max}$  when  $k \rightarrow \infty$  and<sup>11</sup>*

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

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

$\forall v$      $\swarrow$  density  $sa$   
 $\searrow$  target  $v$      $Q^* \in \mathcal{F}$

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

$\uparrow$  density  
 behavior 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>*

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

Munos  
2003  
2008-2010

Bellman backup

$\forall f \in \mathcal{F}, Tf \in \mathcal{G}$ .  
BF  
Completeness

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

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

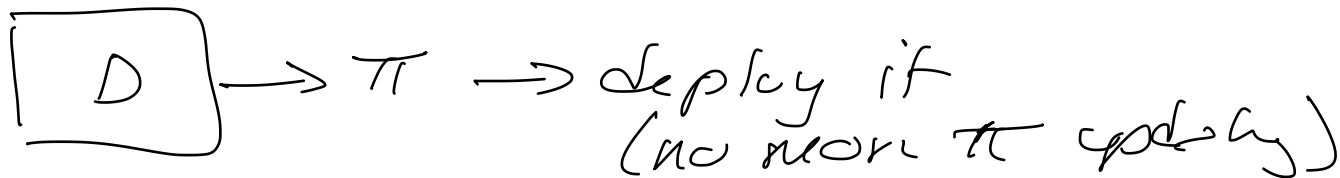
Overlap assumption: Concentratability coefficient

# 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
- 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



robust MDP 2003-2005

param uncertainty in MDPs

1990s

# Check Your Intuition Solutions

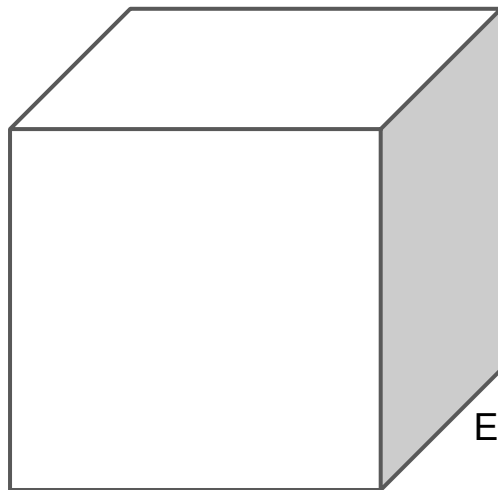
- Optimism under uncertainty can enable sublinear regret in online multi-armed bandits
- Pessimism under uncertainty can lead to linear regret in online multi-armed bandits
- With high probability the optimistic upper bound on the selected arm in UCB algorithms is an upper bound on the performance of 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

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



**Assumptions**

- Markov?
- **Overlap?**
- **Sequential ignorability?**

**Evaluation  
Criteria**

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

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

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

$S_0$ : Set of initial states

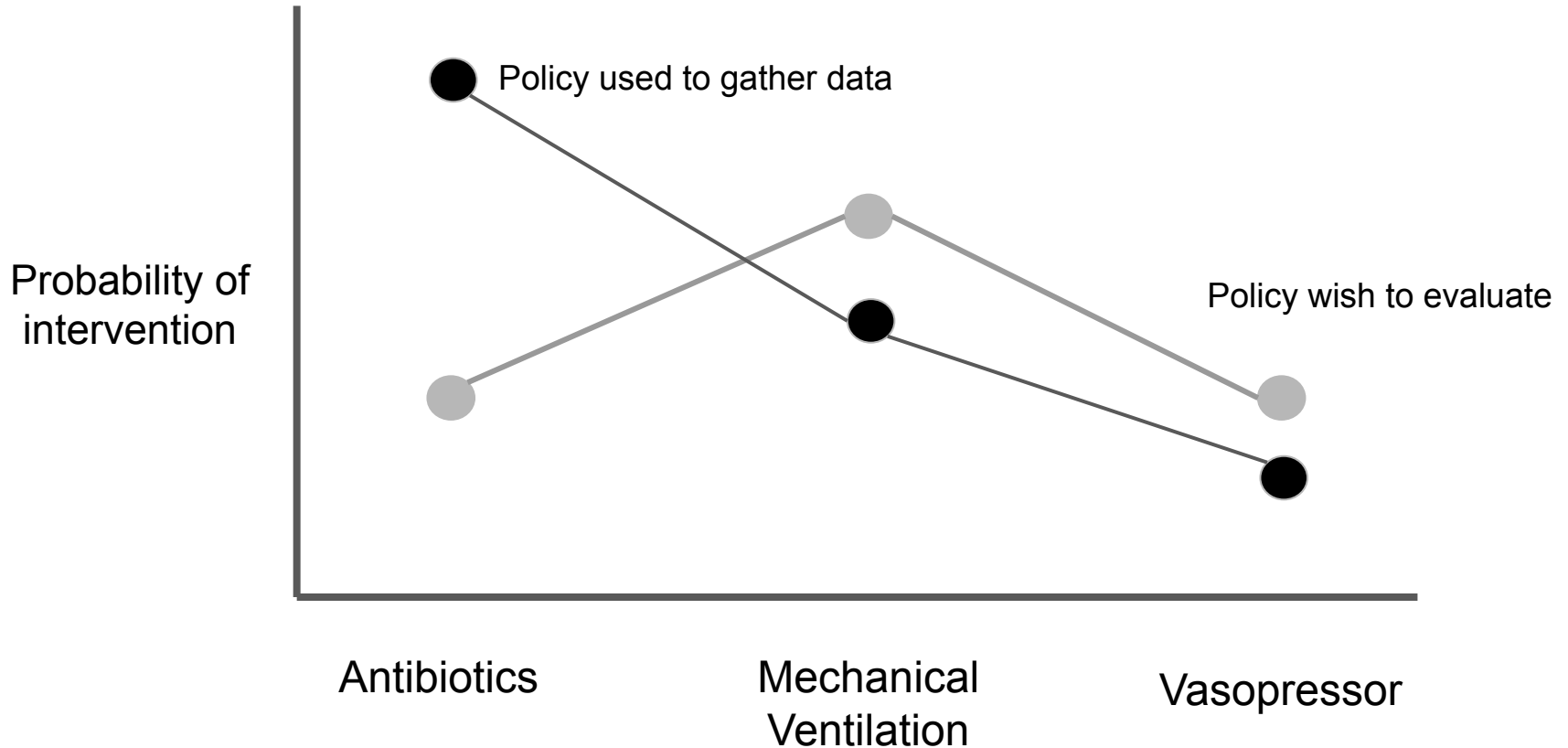
$\hat{V}^\pi(s, \mathcal{D})$ : Estimate  $V(s)$  w/dataset  $\mathcal{D}$

# Standard Assumptions for Off Policy / Counterfactual Estimation & Optimization

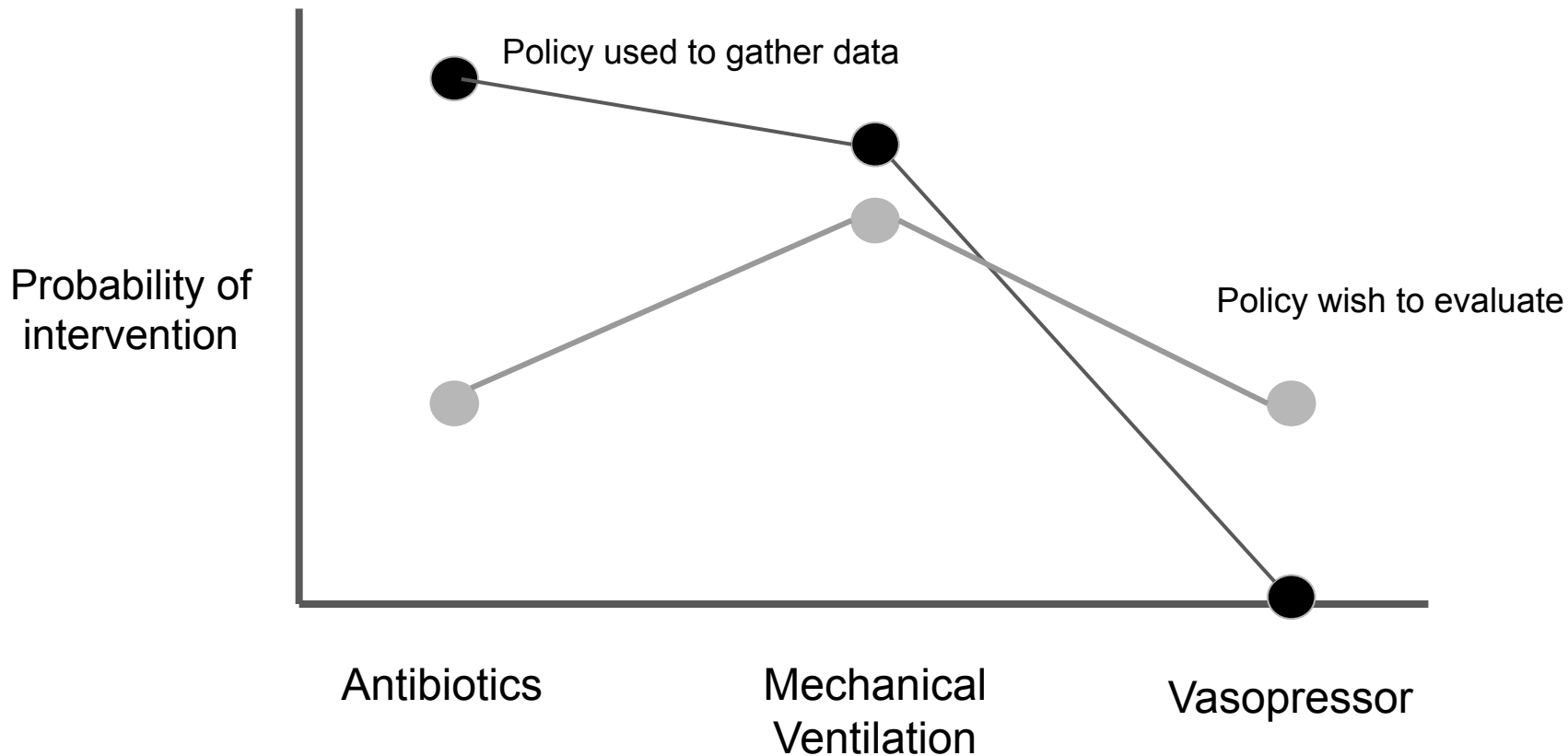
- Overlap
  - Have to take all actions that target policy would take
  - In infinite data / finite data
- No confounding

$\mathcal{D}$ : 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

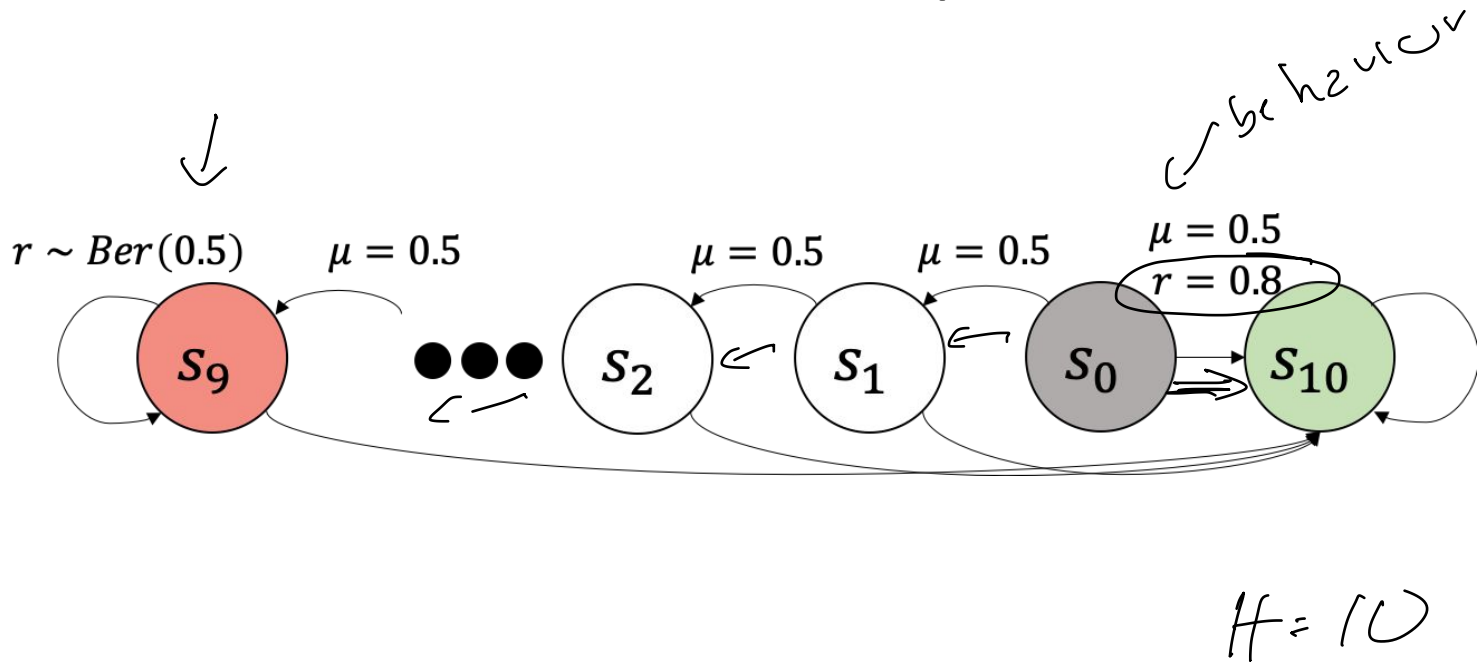
- Typically assume overlap
  - Off policy estimation: for policy of interest
  - Off policy optimization: for all policies including optimal one (see concentrability assumption in batch RL)
- Unlikely to be true in many settings
- Many real datasets don't include complete random exploration
- Assuming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

## 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

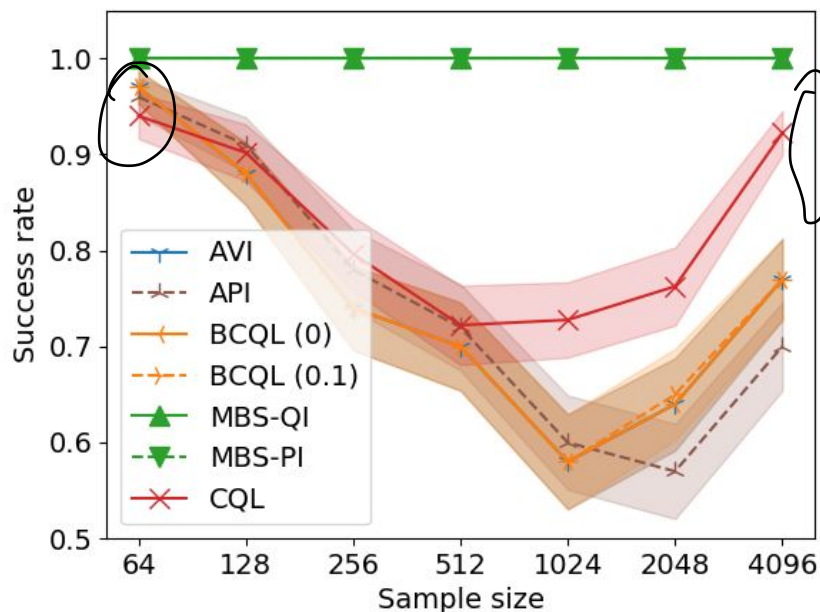
*Common challenge that's attracted substantial interest in last few years but...*

# Illustrative Examples

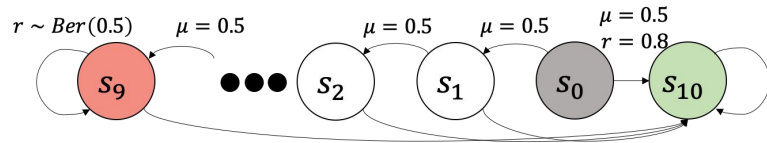




# Recent Conservative Batch Reinforcement Learning Are Insufficient



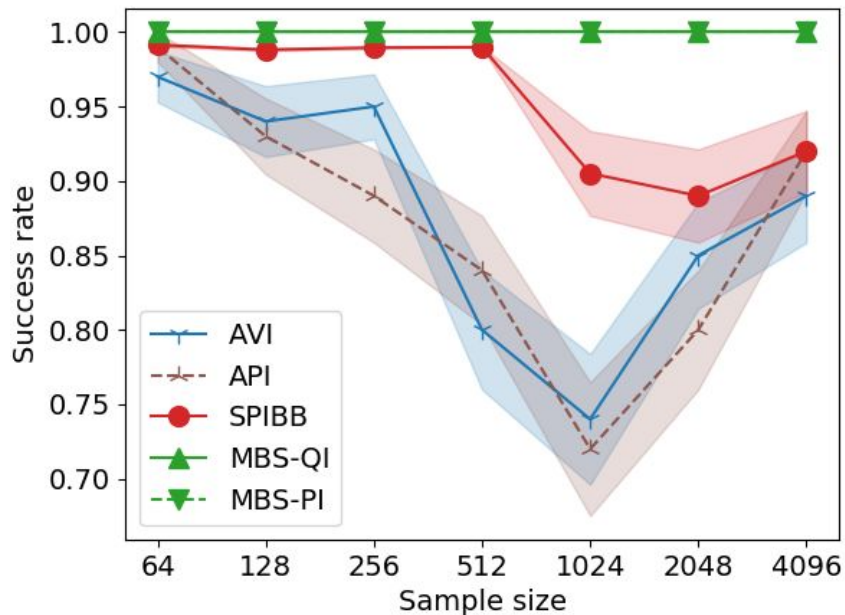
Success rate: #(getting the optimal policy)/#(trials)



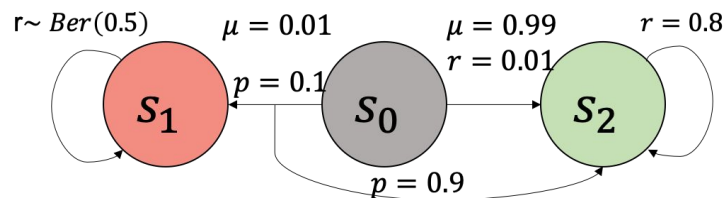
Reasons why baselines fail:

- 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



Success rate: #(getting the optimal policy)/#(trials)



Reasons why baselines fail:

- SPIBB adds conservatism based on estimates of  $\pi_b$  &  $V$  of  $\pi_b$ .
- In this example, the actions which is 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

density in behavior data

- Filtration function:

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

# 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:

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

# 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:

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

$\Rightarrow = 0$  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:

$$\begin{aligned} \mathcal{T}f(s, a) &= r(s, a) + \gamma \mathbb{E}_{s'} \left[ \max_{a'} \zeta(s', a') f(s', a') \right] \\ \mathcal{T}^\pi f(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P, a' \sim \pi} [\zeta(s', a') f(s', a')] \end{aligned}$$

# 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:

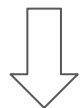
$$\begin{aligned} \mathcal{T}f(s, a) &= r(s, a) + \gamma \mathbb{E}_{s'} \left[ \max_{a'} \zeta(s', a') f(s', a') \right] \\ \mathcal{T}^\pi f(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P, a' \sim \pi} [\zeta(s', a') f(s', a')] \end{aligned}$$



# 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

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

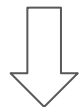


$$V^* - V^{\pi_{\mathcal{A}}} \leq \epsilon$$

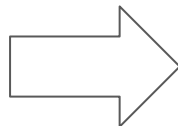
# Best in Well Supported Policy Class\*

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

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



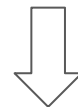
$$V^* - V^{\pi_{\mathcal{A}}} \leq \epsilon$$



**Define**

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

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



$$\max_{\pi' \in \Pi_{all}} V^{\pi'} - V^{\pi_{\mathcal{A}}} \leq \epsilon$$

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

$$\eta_h^\pi(s) := \Pr[s_h = s | \pi],$$

$$\eta_h^\pi(s, a) = \eta_h^\pi(s) \pi(a | s)$$

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

# Theoretical Result

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

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

Error bounds <sup>1:</sup>

• PI:

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

*v = ln func*

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

**Note: Results are for  
function approximation,  
finite sample setting**

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>:

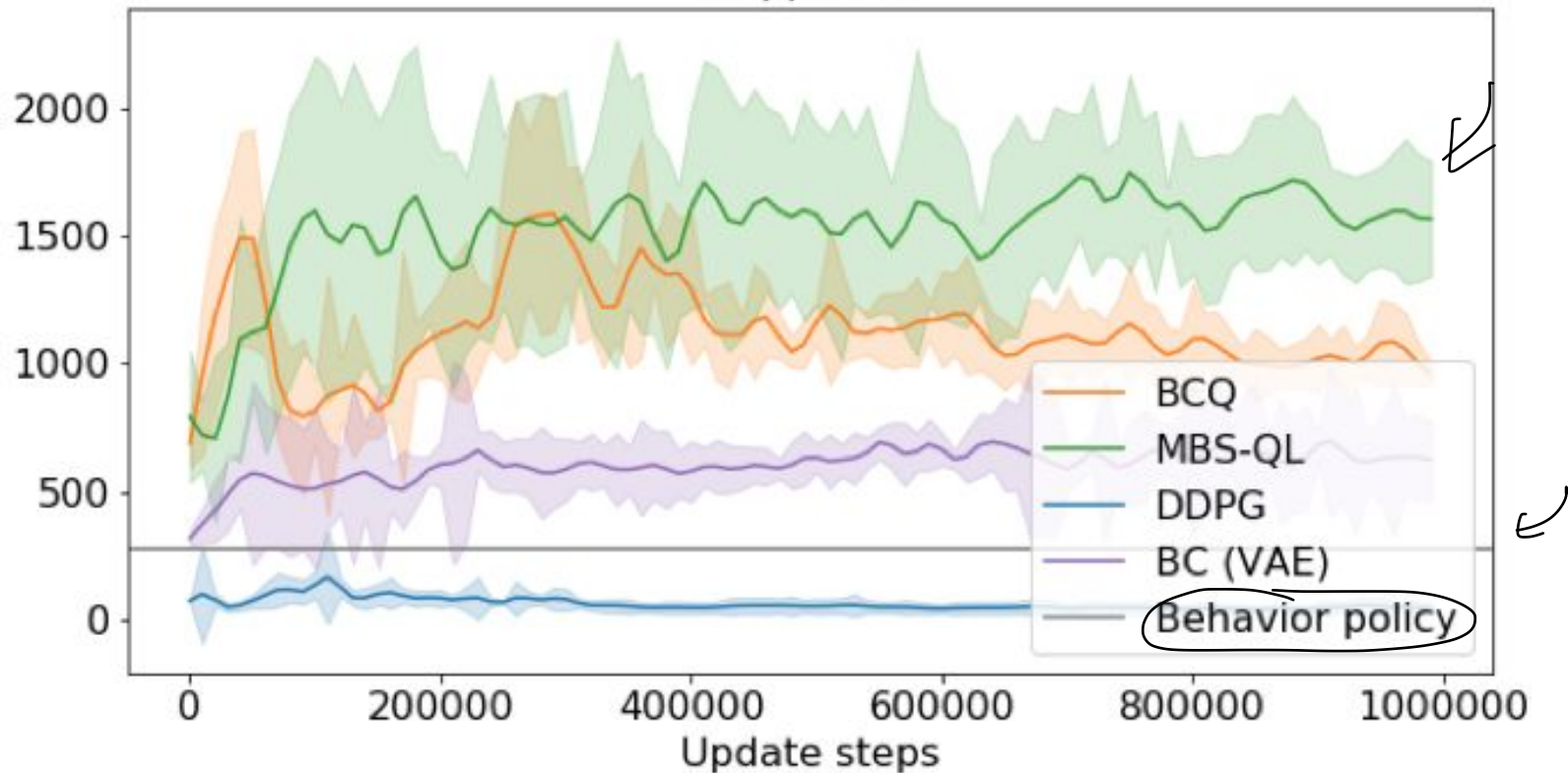
$$O\left(\frac{V_{\max}}{(1-\gamma)^2 b} \sqrt{\frac{\ln(|\mathcal{F}|/\delta)}{n}}\right) + \frac{V_{\max} \epsilon_\zeta}{1-\gamma}$$

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.

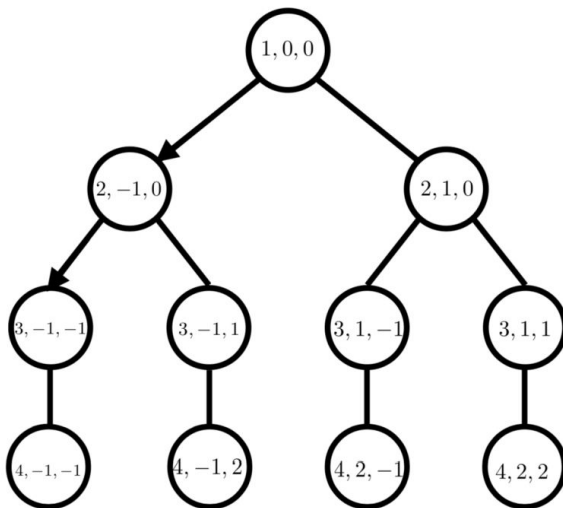
# Can Do Get Substantially Better Solutions, With Same Data

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



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

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

Tu & Recht COLT 2019

Sun, Jiang, Krishnamurthy,  
Agarwal, Langford 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  $\mathcal{D}$

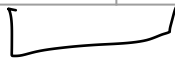


# Early Comparison with Concurrent Work

	<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

# 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



# 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>

*Science* November 2019

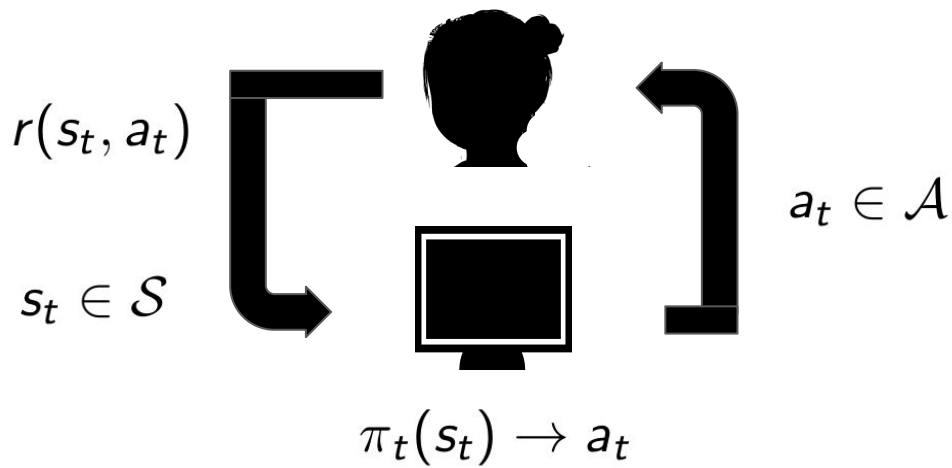


# Optimizing while Ensuring Solution Won't, in the Future, Exhibit Undesirable Behavior

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



# 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

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

Default policy

Potential policy

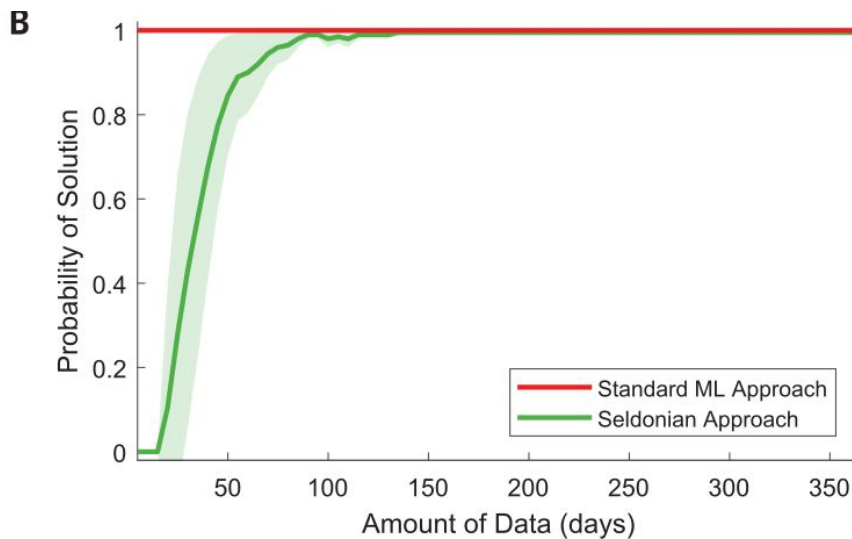
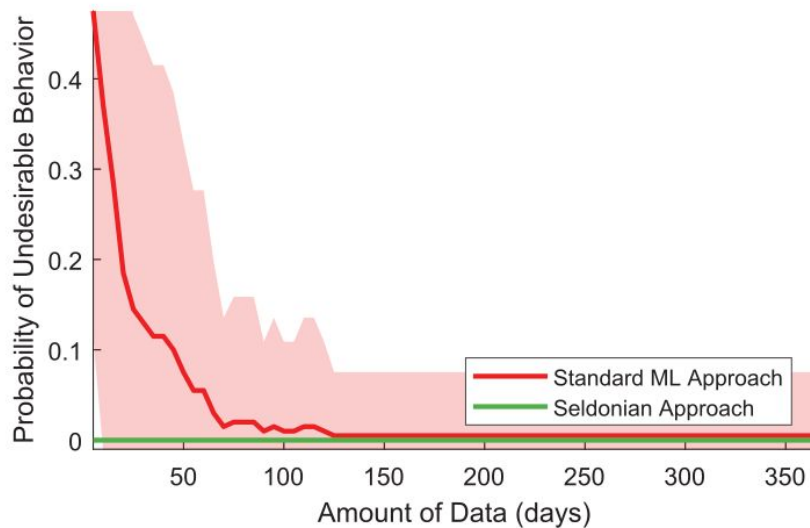
- $r'(H)$  is a function of the trajectory  $H$

# Diabetes Insulin Management

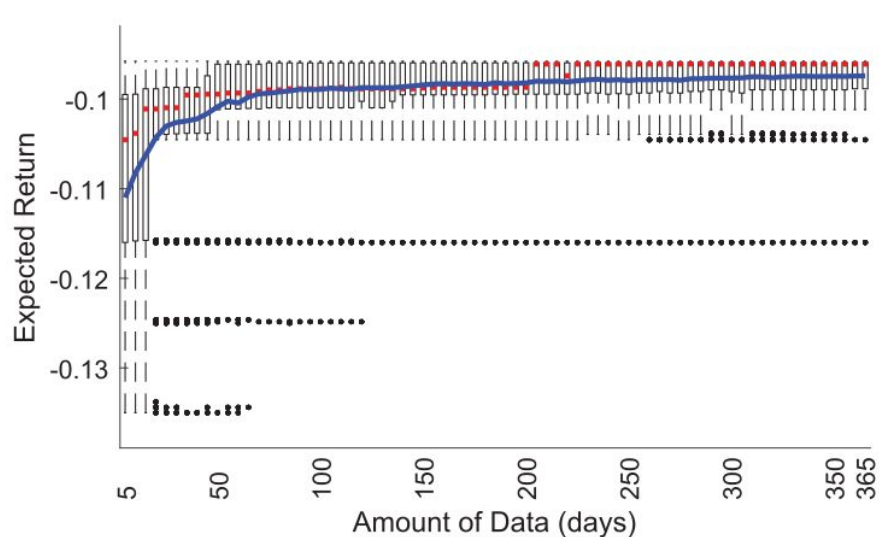


- Blood glucose control
- Action: insulin dosage
- Search over policies
- Constraint:  
hypoglycemia
- Very accurate simulator:  
approved by FDA to  
replace early stage  
animal trials

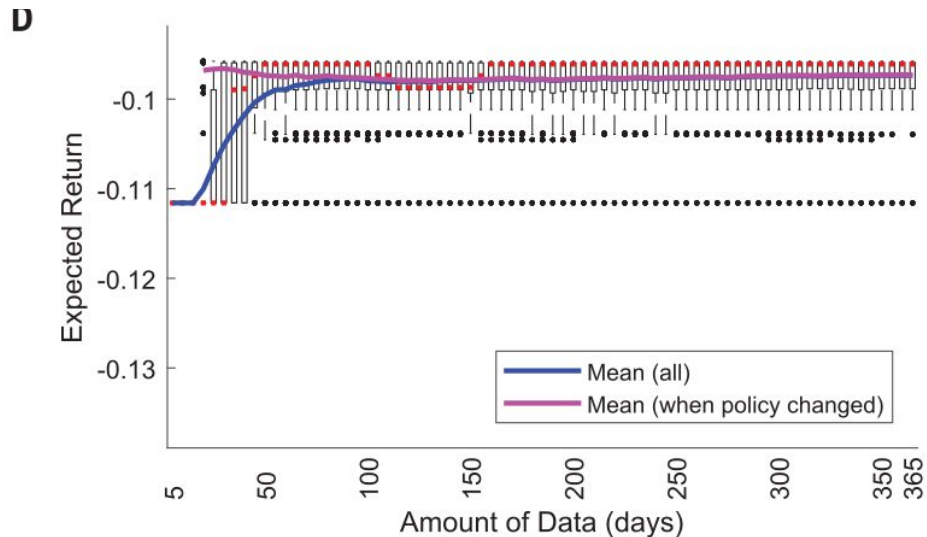
# Personalized Insulin Dosage: Safe Batch Policy Improvement



# 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

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

⇒ 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