# CS292F StatRL Lecture 11
# Exploration in Linear MDP
# & Introduction to offline RL

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

# Logistics

- Project midterm milestone due
  - Important as I need to allocate space for student presentation

- For those who haven't submitted HW1
  - You don't have to solve everything, just submit what you have
  - HW1 is long I am thinking of adjusting grading criteria

- HW2 is not as long
  - Don't wait

# Recap: Lecture 10

- Exploration in Linear MDPs

- Properties of Linear MDPs

- Algorithm: UCB-VI for Linear MDPs

- Regret analysis

# Recap: Impossibility results

- What are the assumptions to make?
  - **Q\*(s,a) approximately linear?**

  - **$Q^\pi$(s,a) is approximately linear for all π?**

  - Q\*(s,a) is exactly linear?

  - $Q^\pi$(s,a) is exactly linear for all π?

Weisz et al (ALT-2020):
http://proceedings.mlr.press/v132/weisz21a.html

$s_a \rightarrow \phi(s_a)$

$\exists \theta$
s.t $\theta^T \phi(s_a)$
$\approx Q^*(s_a)$

for each $\pi$
$\exists \theta^\pi$ $Q^\pi_{(s_a)} = (\theta^\pi)^T \phi(s_a)$

open problem

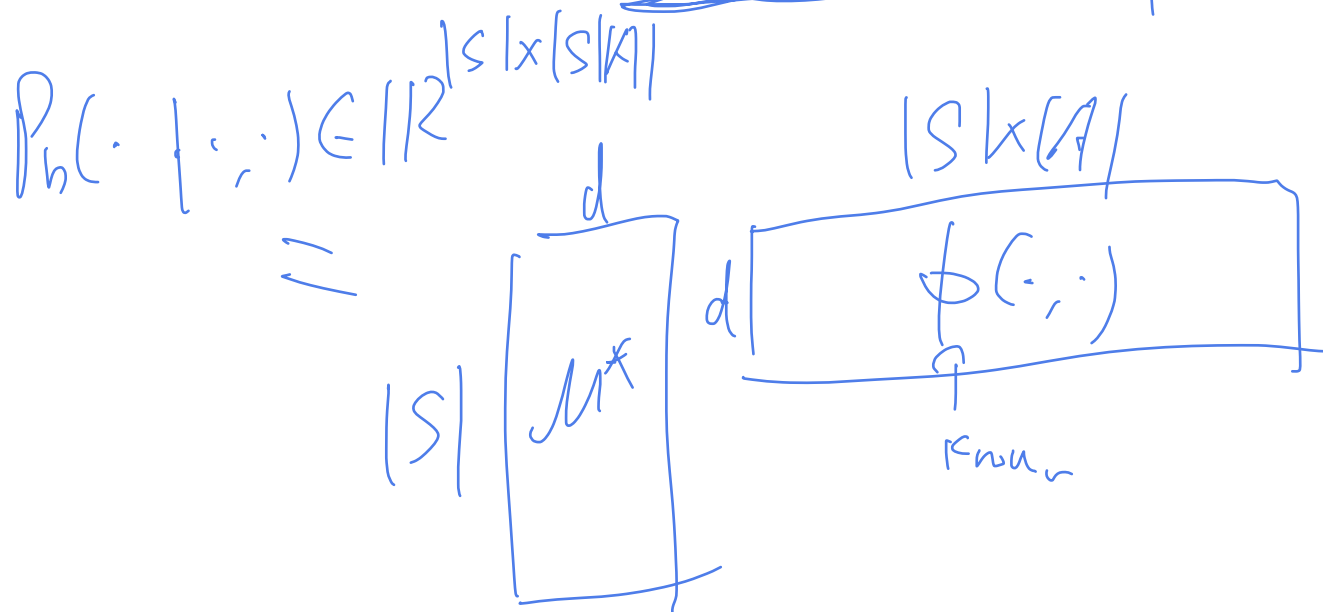Exponential sample complexity / regret lower bounds for the approximate case…

(Du, Kakade, Wang, Yang, 2019) Is a good representation sufficient for sample efficient reinforcement learning?

4

# Recap: Linear MDPs

- Exists feature map $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$
  - Such that:

$$r_h(s,a) = \theta_h^{\star} \cdot \phi(s,a), \quad P_h(\cdot|s,a) = \mu_h^{\star}\phi(s,a), \forall h$$

$|S|$ is exponentially large

$P_h(\cdot | \cdot, \cdot) \in \mathbb{R}^{|S| \times |S||A|}$

$=$

$|S| \begin{bmatrix} d \\ \mu^{\star} \end{bmatrix}$

$d \left| \begin{array}{c} |S| \times |A| \\ \phi(\cdot, \cdot) \\ \text{known} \end{array} \right.$

(Jin et al., 2020) Provably efficient reinforcement learning with linear function approximation

# Recap: UCB-VI for Linear MDPs

- In every round:

  1. Run Ridge regression for estimating the model

  $$\widehat{\mu}_h^n = \mathrm{argmin}_{\mu \in \mathbb{R}^{|\mathcal{S}| \times d}} \sum_{i=0}^{n-1} \left\| \mu\phi(s_h^i, a_h^i) - \delta(s_{h+1}^i) \right\|_2^2 + \lambda \|\mu\|_F^2.$$

  $$\widehat{\mu}_h^n = \sum_{i=0}^{n-1} \delta(s_{h+1}^i)\phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$$

  2. Construct the exploration bonuses

  $$b_h^n(s, a) = \beta\sqrt{\phi(s,a)^\top (\Lambda_h^n)^{-1}\phi(s,a)},$$

  3. Run optimistic value iterations, and update greedy policy

# Recap: Regret bound

- Choose
$$\beta = Hd \left( \sqrt{\ln \frac{H}{\delta}} + \sqrt{\ln(W+H)} + \sqrt{\ln B} + \sqrt{\ln d} + \sqrt{\ln N} \right)$$

$$\lambda = 1$$

$$b_h^n(s,a) = \beta \sqrt{\phi(s,a)^\top (\Lambda_h^n)^{-1} \phi(s,a)},$$

$$\| \phi(s,a) \|_{\Lambda_h^{n-1}}$$

$$\Lambda_h^n = \sum_i \phi(s_h^i, a_h^i)$$

$$\phi(s_h^i, a_h^i)$$

$$R^{d \times q}$$

$$+ \lambda I$$

- Regret
$$\tilde{O} \left( H^2 \sqrt{d^3 N} \right)$$

# of episodes

# Recap: Regret analysis

- Regret of episode t
  - Optimism / simulation lemma

- Sum them up to get total regret

  - Same information-gain bound from linear bandits

$$-\hat{P}(\cdot|s,a) + \hat{P}(\cdot|s,a) = \hat{\mu} \cdot \phi(s,a) - \mu^* \cdot \phi(s,a)$$

# Recap: It remains to prove

$$f: S \to R$$

- 1. Uniform convergence bound

$$\left(\hat{P}(\cdot|s,a) - P(\cdot|s,a)\right) \cdot f(\cdot)$$

$$\sup_{f \in \mathcal{F}} \left| \left(\hat{P}(\cdot|s,a) - P(\cdot|s,a)\right) \cdot f \right| \le \varepsilon$$

$$= \sum_{s'} \left(\hat{P}(s'|s,a) - P(s'|s,a)\right) f(s')$$

$$\sum_{s'} \hat{P}(s'|s,a) \cdot f(s')$$

- 2. "Optimism"

The same induction argument as in the UCB-VI for tabular MDP
(Read Lemma 7.10 in AJKS)

- 3. "Information gain" bound

The same argument as in the Linear Bandits case.
(Read Lemma 7.12 in AJKS)

se the fact that $\delta(s)^\top V = V(s)$. Thus the operator $P_h^n(\cdot|s,a) \cdot V$ simply requires storing all data

d via simple linear algebra and the computation complexity is simply $\text{poly}(d,n)$—no poly depend

# Recap: Bound for a fixed V

late the difference between $\widehat{\mu}_h^n$ and $\mu_h^\star$.

(Difference between $\widehat{\mu}_h$ and $\mu_h^\star$). *For all $n$ and $h$, we must have:*
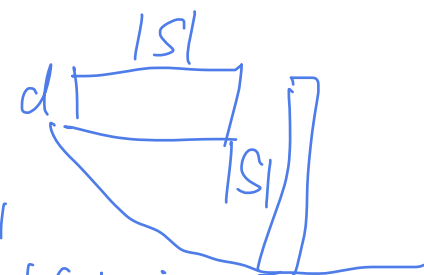
- Lemma 7.3 AJKS

$$\widehat{\mu}_h^n - \mu_h^\star = -\lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum_{i=1}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}.$$

- The quantity of interest is a inner product with this:

start from the closed-form solution of $\widehat{\mu}_h^n$:

$(s_{h+1}^i)\phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} = \sum_{i=0}^{n-1} (P(\cdot|s_h^i, a_h^i) + \epsilon_h^n)\phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$

$\mu_h^\star \phi(s_h^i, a_h^i) + \epsilon_h^i)\phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} = \sum_{i=0}^{n-1} \mu_h^\star \phi(s_h^i, a_h^i)\phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} + \sum_{i=0}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top$

$_h^\star \phi(s_h^i, a_h^i)\phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} + \sum_{i=0}^{n-1} \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}$

$_h^n - \lambda I)(\Lambda_h^n)^{-1} + \sum \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1} = \mu_h^\star - \lambda \mu_h^\star (\Lambda_h^n)^{-1} + \sum \epsilon_h^i \phi(s_h^i, a_h^i)^\top (\Lambda_h^n)^{-1}.$

# Challenge: we cannot use union bound because we have an infinite number of value functions

$$\sup_{f \in F} \left\| \sum \phi(S_n^i, A_n^i) \cdot \varepsilon_i^T f \right\|$$

$$\left( \frac{\#}{\sqrt{n}} \right)$$

- A covering number argument.



$$\forall f \in F$$
$$\exists \check{f} \in F_\varepsilon$$

$$s.t. \left\| f - \check{f} \right\|_\infty \leq \varepsilon$$

- Covering number: the number of balls with radius ε that is needed to cover all points in a set.

$$\sup_X \left| f(x) - \check{f}(x) \right|$$
$$\leq \varepsilon$$

$$N_\varepsilon = \# \text{ of balls needed to cover } F$$
$$s.t. \left| f - \check{f} \right| \leq \varepsilon \quad \forall f \; \exists \check{f} \in F_\varepsilon$$



11

# Family of value functions we consider

$$V_{h,\theta} = \max_a \{ b_h(s,a) + \phi(s,a)^\top \left( \mu^{*\top} V_{h+1} \in \mathbb{R}^d \right) \}$$

$$f_{w,\beta,\Lambda}(s) = \min \left\{ \max_a \left( \overline{w^\top \phi(s,a) + \beta \sqrt{\phi(s,a)^\top \Lambda^{-1} \phi(s,a)}} \right), H \right\}, \forall s \in \mathcal{S}.$$

$$b_h$$

$$\mathcal{F} = \{ f_{w,\beta,\Lambda} : \|w\|_2 \leq L, \beta \in [0, B], \sigma_{\min}(\Lambda) \geq \lambda \}.$$

$$V \in \mathcal{F}$$

What is a finite set to cover this class such that for every f in this set, there is a function in the finite set, such that they are ε-close in sup-norm?

$$\text{Lemma:} \quad N_\varepsilon \left( \{ x \in \mathbb{R}^d \mid \|x\|_2 \leq B \} \right) = O \left( \left( \frac{B}{\varepsilon} \right)^d \right)$$

# Covering number calculations

$$\bar{F}_\xi = W_{\frac{\xi}{3}} \times B_{\frac{\xi}{3}} \times \Lambda^{-1}_{\frac{\xi}{3}} \qquad \|\phi_{sa}\|_2 \le 1$$

$$f \in \bar{F}(w, \beta, \Lambda) \qquad \hat{f} \in \bar{F}_\xi \qquad (\hat{w}, \hat{\beta}, \hat{\Lambda})$$

$$\left| f(\xi) - \hat{f}(\xi) \right| \le \left| \max_a \left( w^T \phi(s_a) + \beta \sqrt{\phi_{sa}^T \Lambda^{-1} \phi_{sa}} \right) - \max_a \left( \hat{w}\phi(s_a) + \hat{\beta}\sqrt{\phi_{(s_a)}^T \hat{\Lambda}^{-1} \phi_{(s_a)}} \right) \right|$$

$$\le \max_a \left| (w - \hat{w})^T \phi(s_a) \right| + \max_a \left| (\beta - \hat{\beta}) \sqrt{\phi_{sa}^T \Lambda^{-1} \phi_{sa}} \right| + \max_a \left( \hat{\beta} \sqrt{\phi_{(s_a)}^T \Lambda^{-1} \phi_{sa}} \right)$$

$$\le \|w - \hat{w}\|_2 + \frac{\|\beta - \hat{\beta}\|}{\sqrt{\lambda}} + B\sqrt{\|\Lambda^{-1} - \hat{\Lambda}^{-1}\|_F} \qquad \sqrt{\phi_{sa}^T \Lambda^{-1} \phi_{r_i}}$$

$$\le \frac{\xi}{3} \qquad \le \frac{\xi}{3} \qquad \frac{\xi}{3}$$

$$N_\xi(\bar{F}) \le \xi \qquad N_\xi = |\bar{F}_\xi| = \|W_{\frac{\xi}{3}}\| \cdot \|B_{\frac{\xi}{3}}\| \cdot \|\Lambda^{-1}_{\frac{\xi}{3}}\| \qquad x^T A x = tr(x^T A x)$$
$$= tr(A x x^T)$$
$$\le \langle A, x x^T \rangle$$
$$\le \|A\|_F \|x x^T\|_F$$

$$\left( \frac{1}{\xi} \right)^{3d}$$

$$\bar{F}_\xi = \left\{ f_{W \beta \Lambda} \ s.t \ \middle| \ @ \ w \in W_{\frac{\xi}{3}}, \ B \in B_{\frac{\xi}{3}}, \ \Lambda^{-1} \in \Lambda^T_{\frac{\xi}{3}} \right\}$$

13

# From covering number to a uniform convergence bound

$$\sup_{f \in F} \left\| \sum \phi_{(s_i, a_i)} \cdot \varepsilon_i^T f \right\|_{(\Lambda_n^n)^{-1}} \leq \sup_{f \in F} \left\| \sum \phi_{s_i a_i} \varepsilon_i^T (f - \breve{f}_f + \breve{f}_f) \right\|_{\Lambda_n^{-1}}$$

$$\leq \sup_f \left\| \sum \phi_{\infty} \varepsilon_i^T (f - \breve{f}_f) \right\|_{\Lambda^{-1}} + \sup_f \left\| \sum \phi_{s_i a_i} \varepsilon_i^T \breve{f}_f \right\|_{\Lambda^{-1}}$$

$$\leq \underbrace{2\varepsilon \cdot \left\| \sum \phi_i \right\|_{\Lambda^{-1}}}_{\leq 2\varepsilon M} + \sup_{\breve{f}} \left\| \sum \phi_{s_i a_i} \varepsilon^T \breve{f} \right\|_{\Lambda^{-1}}$$

$$\left\| \sum \phi_i \right\|_{\Lambda^{-1}}$$

$$\leq \sum_i \sqrt{\phi_i^T \Lambda^{-1} \phi_i}$$

$$\leq \sqrt{N \sum_i \phi_i^T \Lambda^{-1} \phi_i}$$

$$\leq N$$

$$\varepsilon_i := \delta_i - \mathbb{E}\delta_i$$

$$\|\varepsilon_i\|_1 \leq 2$$

$$\|f - \breve{f}\|_\infty \leq \varepsilon$$

1. apply pointwise result for fixed $\breve{f}$
2. apply union bound

$$\left( d \sqrt{d \log N} + \log \frac{N\varepsilon}{\delta} \right)$$

$$\Lambda_i = \sum_{b=1}^{N} \phi_b \phi_b^T + \lambda I$$

14

# Final notes about linear MDPs

- A semi-parametric model
  - The number of parameters to describe the model can be exponentially large:  $\underline{d\ S}$  *describe* $M^*$
  - Efficient algorithm with regret independent to S

- Still very strong assumption on the feature map
  - Interesting open problems:
    - Representation learning    $\underline{\phi}$ *is unknown*
    - Nonlinear parametric models
    - Suboptimal rates when naively applying to the tabular case
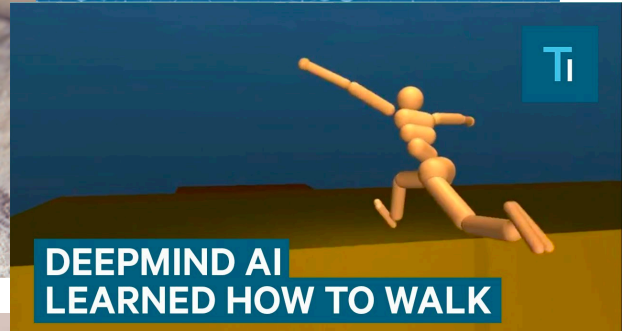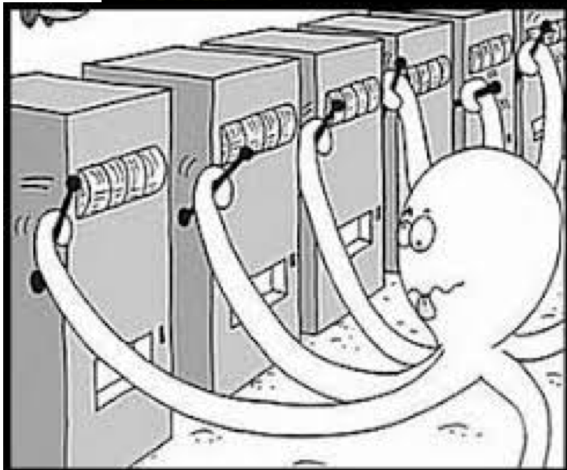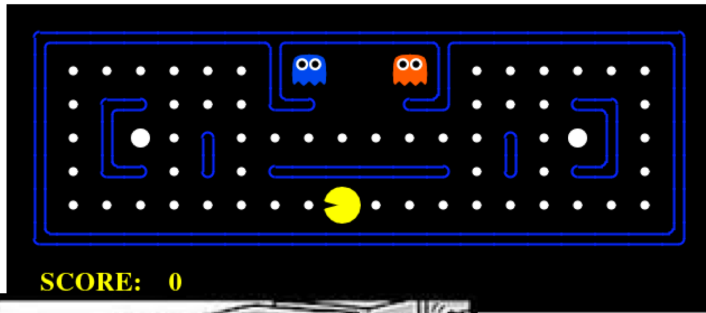
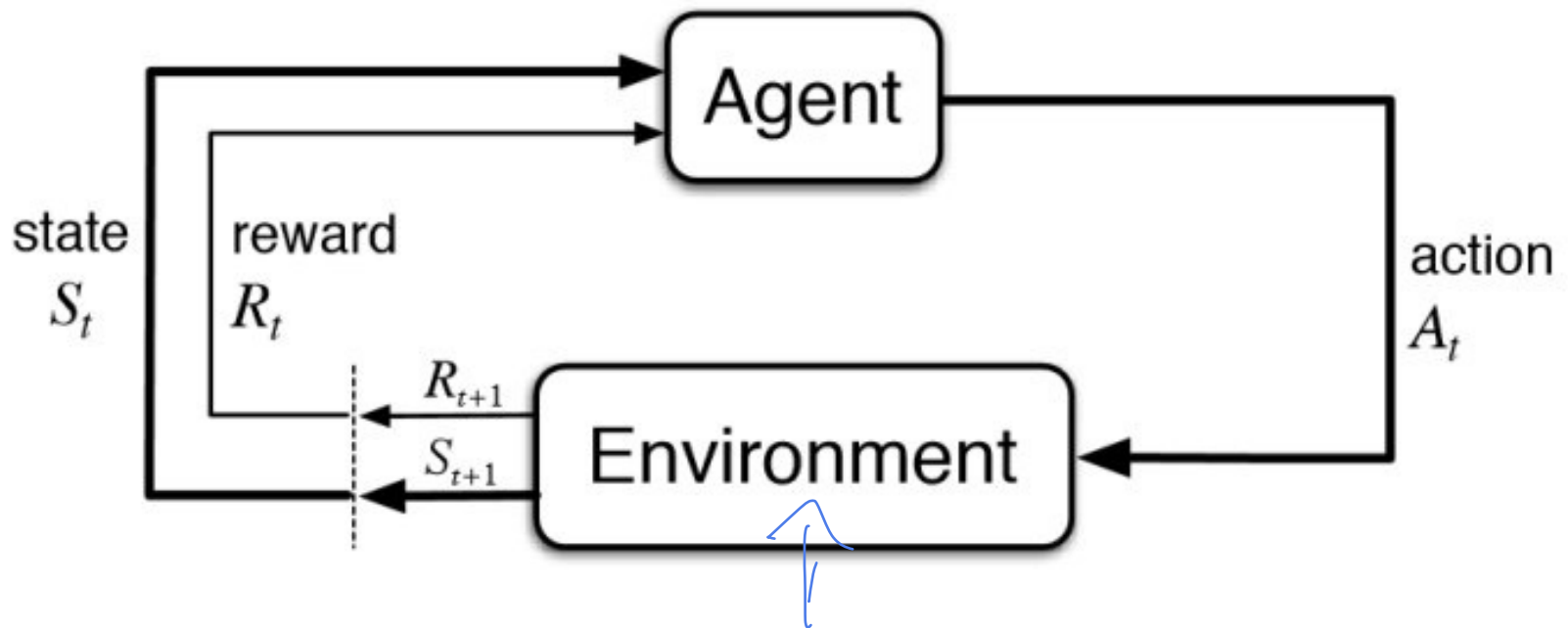$d = S$    *but* $\dfrac{d\sqrt{T}}{\sqrt{ST}}$

# Remainder of the lecture

- Introduction to offline reinforcement learning

- Off-policy evaluation in contextual bandits

# Recap: RL is among the hottest area of research in ML!

# An RL agent learns interactively through the feedbacks of an environment.



state $S_t$  reward $R_t$  action $A_t$  $R_{t+1}$  $S_{t+1}$

Agent — Environment

- Learning how the world works (dynamics) and how to maximize the long-term reward (control) at the same time.

# Applications of RL in the real life

- RL for robotics.
- RL for dialogue systems.
- RL for personalized medicine.
- RL for self-driving cars.
- RL for new material discovery.
- RL for sustainable energy.
- RL for feature-based dynamic pricing.
- RL for maximizing user satisfaction.
- RL for QoE optimization in networking
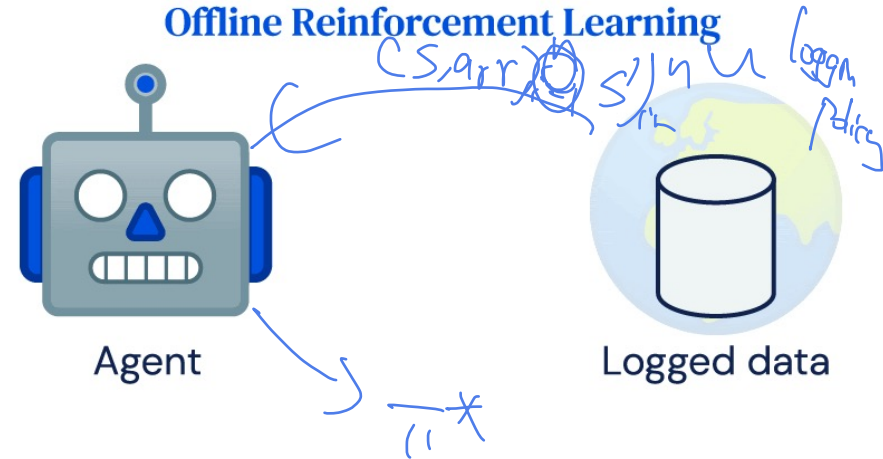- …

# Challenges of Reinforcement in the real life

- No access to a simulator
- Every data point is costly.
- Legal, safety issues associated with exploration.
- Large / complex state-space, action space.
- Long horizon
- Limited adaptivity (cannot run too many iterations)

# From an Applied ML Scientist point of view, the starting point of a project is often:

# Online RL vs Offline RL



**Online Reinforcement Learning**

$a = \pi(s)$

$s', r$

Agent → Environment

**Offline Reinforcement Learning**

$(s, a, r, s'), \eta \sim \mu$ logging policy

$\pi^*$

Agent — Logged data
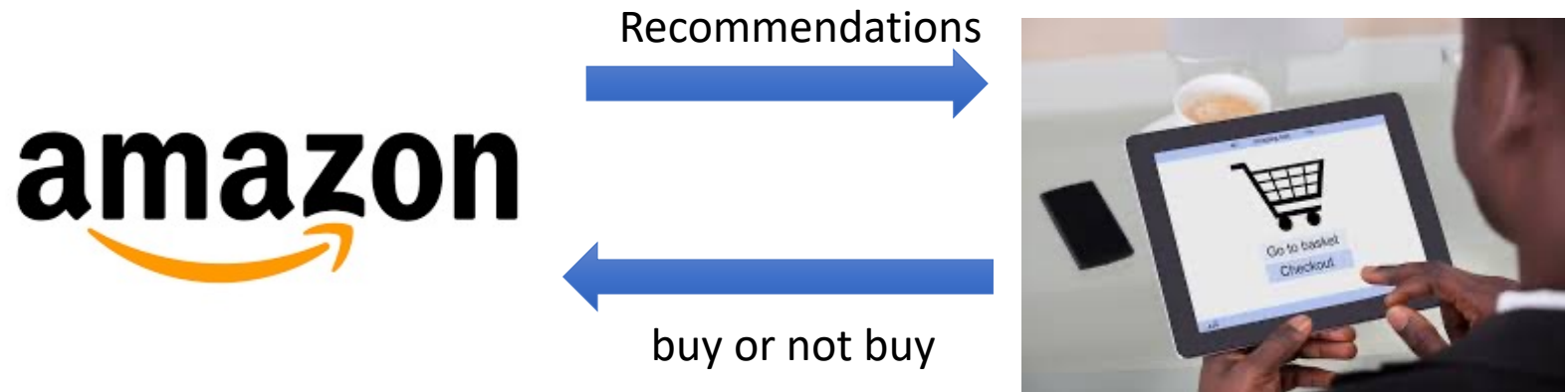
Exploration is often **expensive**, **unsafe**, **unethical** or **illegal** in practice, e.g., in self-driving cars, or in medical applications.

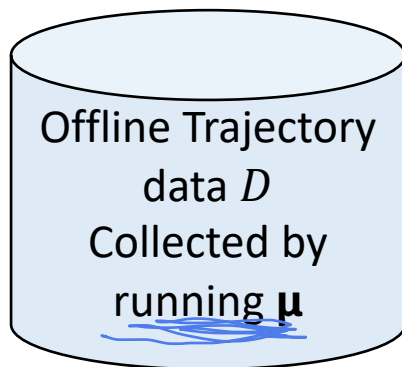Can we learn a policy from already **logged interaction data**?

# Off-Policy learning: an example

Recommendations

buy or not buy

- How to evaluate a new algorithm without actually running it live?
- How to learn a better system than the one that is deployed.

# Offline Reinforcement Learning, aka. Batch RL

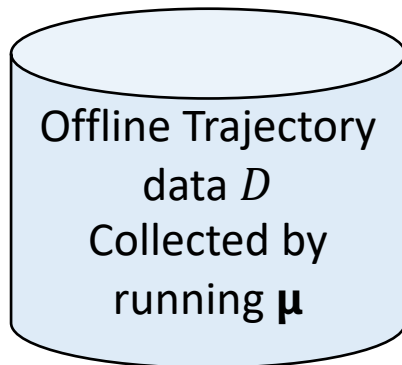- Task 1: Offline Policy Evaluation. (OPE)



Offline Trajectory data $D$ Collected by running $\mu$

Task: design OPE methods

Evaluate fixed Target Policy $\pi$

**Via Uniform OPE**

- Task 2: Offline Policy Learning. (OPL)

Offline Trajectory data $D$ Collected by running $\mu$

Task: design OPO methods

Find near optimal Policy $\hat{\pi}^*$

# Contextual bandits model

- Contexts: /state
  - $x_1, ..., x_n \sim \lambda$    drawn iid, possibly infinite domain    $R^d$

- Actions:
  - $a_i \sim \mu(a|x_i)$    Taken by a randomized "Logging" policy

- Reward:
  - $r_i \sim D(r|x_i, a_i)$    Revealed only for the action taken

- Value:
  - $v^{\mu} = \mathbb{E}_{x \sim \lambda} \mathbb{E}_{a \sim \mu(\cdot|x)} \mathbb{E}_D[r|x, a]$

- We collect data $(x_i, a_i, r_i)_{i=1}^n$ by the above processes.

# Off-policy Evaluation and Learning

**Off-policy evaluation**

Estimate the value of a fixed target policy $\pi$

$$v_\pi := \mathbb{E}_\pi[\text{Reward}]$$

**Off-policy learning**

find $\in \Pi$

that maximizes $v_\pi$

- Using data $(x_i, a_i, r_i)_{i=1}^n$

$$(x_i, a_i, r_i, \mu_i)_{i=1}^n \qquad \mu(a_i|x)$$

- often the policy $\mu$ or logged propensities $(\mu_i)_{i=1}^n$

# ATE estimation is a special case of off-policy evaluation

- a:  Action          $\Leftrightarrow$     T: Treatment  {0,1}
- r: Reward          $\Leftrightarrow$     Y: Response variable
- x: Contexts        $\Leftrightarrow$     X: covariates

# Direct Method / Regression-estimator

- Fit a regression model of the reward

$$\hat{r}(x,a) \approx \mathbb{E}(r|x,a) \quad \text{using the data}$$

- Then for any target policy

$$\hat{v}^{\pi}_{\text{DM}} = \frac{1}{n} \sum_{i=1}^{n} \sum_{a \in \mathcal{A}} \boxed{\hat{r}(x_i, a)} \pi(a|x_i)$$

Pros:

- Low-variance.
- Can evaluate on unseen contexts

Cons:

- Often high bias
- The model can be wrong/hard to learn

# Inverse propensity score / Importance sampling

(Horvitz & Thompson, 1952)

Importance weights

$$\hat{v}^{\pi}_{\text{IPS}} = \frac{1}{n} \sum_{i=1}^{n} \boxed{\frac{\pi(a_i|x_i)}{\mu(a_i|x_i)}} r_i$$

$$=: \rho_i$$

$$\mathbb{E}\left[\hat{v}^{\pi}_{\text{IPS}}\right] = \frac{1}{n}\sum_i \mathbb{E}\left[\frac{\pi(a_i|x_i)}{\mu(a_i|x_i)} \sigma_i\right]$$

$$= E_{x_i}\left[\sum_a \mu(a|x_i) \frac{\pi(a|x_i)}{\mu(a|x_i)} E[r_i]\right]$$

**Pros:**

- No assumption on rewards
- Unbiased
- Computationally efficient

**Cons:**

- High variance when the weight is large

29

# Analyzing the performance of importance sampling estimator

# Importance Sampling and Direct Method are surprisingly similar in some cases

- Consider the MAB case

# Next lecture: OPE for reinforcement learning

- Importance sampling

- Marginalized importance sampling