

Risk sensitive Reinforcement Learning with Low-rank MDPs (In progress)

March 30, 2023

Yulai Zhao

Backgrounds

- What is risk sensitive RL?
- First recall the definition of Conditional Value-at-Risk (CVaR), defined as:
$$CVaR_{\tau}(x) = E_{x \sim P}[x | x < P^{-1}(\tau)]$$
- Often, τ is set to be 0.05, or 0.01
- It is often used as an empirical **RISK** measure, for example, if a portfolio has a high expected return, is it good?
- Not necessarily! We need to figure out what the return will be with the WORST 5% probability.

Backgrounds

- What is risk sensitive RL?
- Some find we can write $CVaR_\tau(x) = E_{x \sim P}[x | x < P^{-1}(\tau)]$ in another way:

$CVaR_\tau(R) = \max_b [b - \frac{1}{\tau} E_{R \sim P} \max(b - R, 0)]$, we can understand b as a threshold thing, where R is the normal reward from a process.

- Imagine the randomness comes from policy π in the MDP, and R is the received reward, then in a risk-sensitive RL, we aim at maximizing the CVaR measure

$$\begin{aligned} \max_{\pi} CVaR_\tau(R) &= \max_{\pi} \max_b [b - \frac{1}{\tau} E_{R \sim P} \max(b - R, 0)] \\ &= \max_b [b - \frac{1}{\tau} \min_{\pi} E_R \max(b - R, 0)] \end{aligned}$$

Note that in a RL language, $R = \sum_{i=1} r(s_i, a_i)$

Setup

- MDP

- a tuple $M = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$: A set of states \mathcal{S} , a set of actions \mathcal{A} , a transition probability $\mathcal{P}: \mathcal{S} \times \mathcal{A} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, a known and deterministic reward function $r: \mathcal{S} \times \mathcal{A} \times \mathcal{A} \rightarrow [0, 1]$, a discounted factor $\gamma \in [0, 1)$.
- We can easily extend $r(s,a)$ to unknown and stochastic case.
- Start from the initial state distribution d_0

- For each threshold b , we have

$$V^\pi(s, b) = E [\max(0, b - \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t)) | s_0 = s, \pi, b]$$
$$J^*(b) = \min_{\pi} J(b) = \min_{\pi} E_{s \sim d_0} V^\pi(s, b)$$

The final goal is to maximize CVaR, i.e., $\max_b (b - \frac{1}{\tau} J^*(b))$

So this means we can only have a finite set of b and do this bi-level optimization?

Setup

-
- For each threshold b , we have Q function

$$Q^\pi(s, b, a) = E [\max(0, b - \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t)) | s_0 = s, a_0 = a, \pi, b] = E_{s' \sim P(\cdot | s, a)} V^\pi(s', b - r(s, a))$$

The final goal is to maximize CVaR, i.e., $\max_b (b - \frac{1}{\tau} J^*(b))$

- Assume the MDP admits a low-rank decomposition

$$\forall s, s', a, \quad P(s' | s, a) = \mu^*(s')^\top \phi^*(s, a)$$

Some Review

- We wish to learn some useful techniques from [Representation Learning for Online and Offline RL in Low-rank MDPs] by Masatoshi Uehara , Xuezhou Zhang, and Wen Sun
- Now we move on to introduce this paper and briefly sketch its proofs

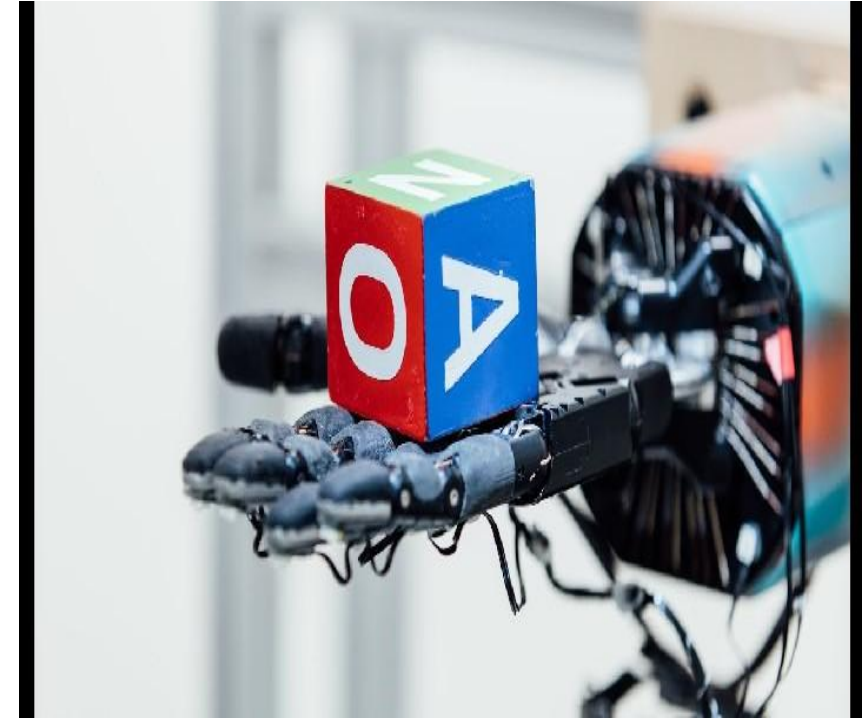
Empirical RL for large-scale problems



[AlphaGo, Silver et.al, 15]



[OpenAI Five, 18]



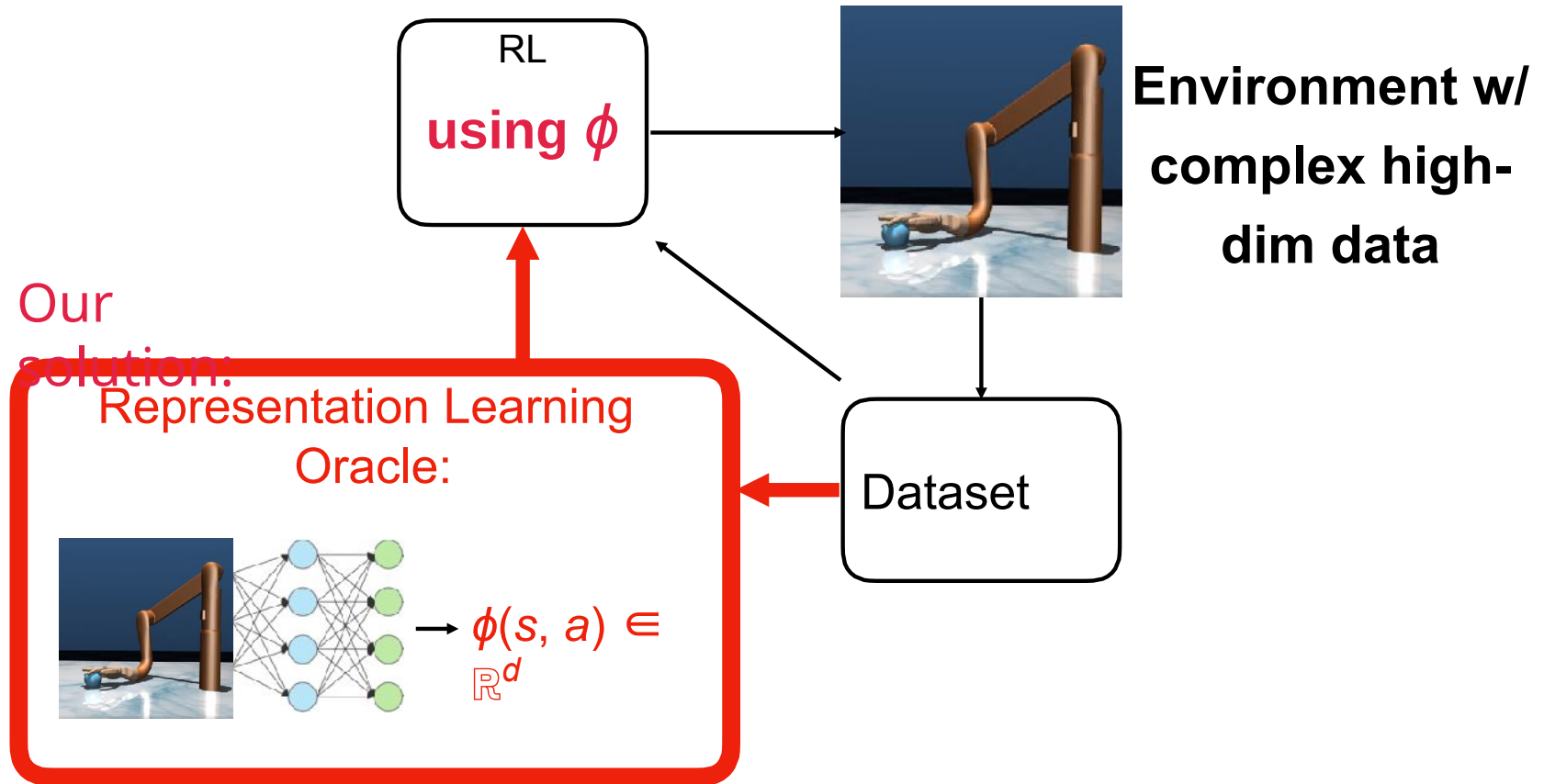
[OpenAI, 19]

Rich (nonlinear) function approximation + RL can work well w/ enough samples

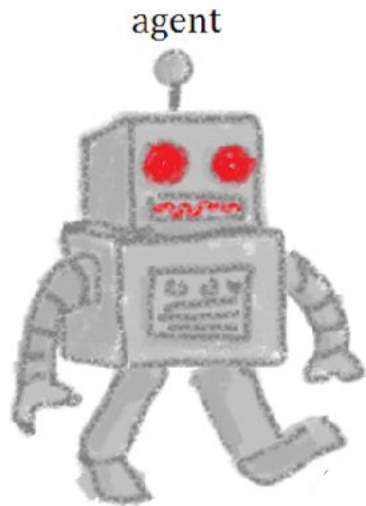
Can we design provably efficient algorithms for

Rich Function Approx + RL

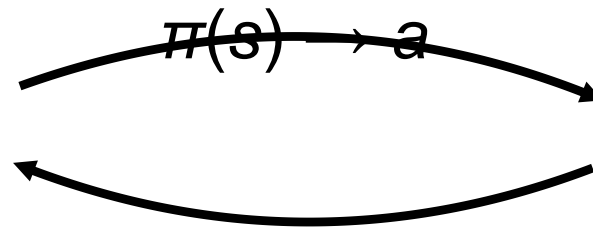
?



Episodic Infinite Horizon Discounted MDPs



Policy: state to action



Reward & Next State

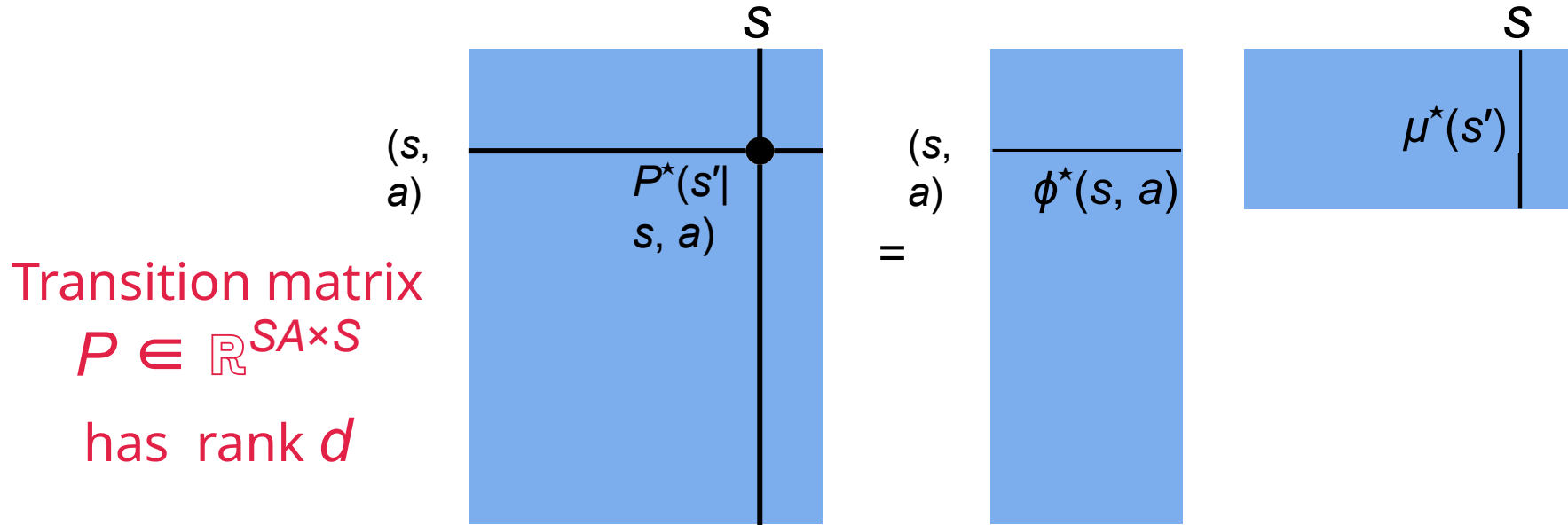
$$r(s, a), s' \sim P(\cdot | s, a)$$

Objective:

$$\max_{\pi} J(\pi; P, r), \text{ where } J(\pi; P, r) := \mathbb{E} [r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots | a \sim \pi, P]$$

Assume fixed initial state s

Low-rank MDP



$$\exists \mu^*, \phi^* \quad : \quad \forall s, a, s', P^*(s'|s, a) = \mu^*(s')^\top \phi^*(s, a)$$

Low-rank MDP **Linear MDPs** (Jin et al, Yang & Wang)

Linear MDP = low-rank + known ϕ^*

The formulation

1. Realizable hypothesis classes Ψ, Φ , and

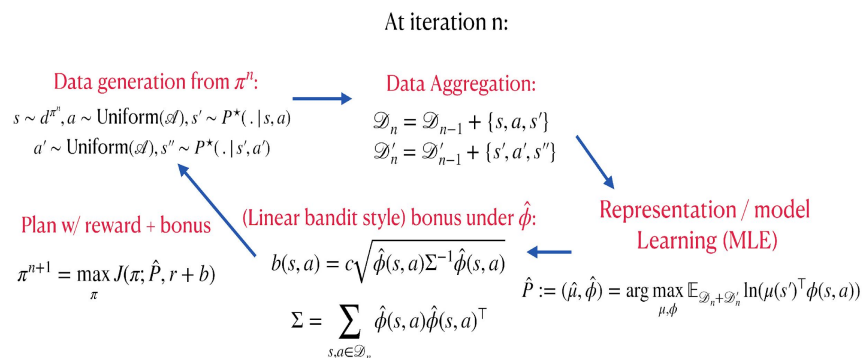
$$\mu^* \in \Psi, \phi^* \in \Phi$$

2. Computation oracle for learning representations: Maximum Likelihood Estimation (MLE): the sum is taken over a collection of tuples: $D = (s_i, a_i, s'_i), i = 1, \dots, n$

$$(\hat{\mu}, \hat{\phi}) = \arg \max_{\mu, \phi} \sum_i^n \ln(\mu(s'_i)^\top \phi(s_i, a_i))$$

Our algorithm: Rep-UCB

(UCB-driven Representation Learning for online RL)



Algorithm 1 UCB-driven representation learning, exploration, and exploitation (REP-UCB)

- 1: **Input:** Regularizer λ_n , parameter α_n , Models $\mathcal{M} = \{(\mu, \phi) : \mu \in \Psi, \phi \in \Phi\}$, Iteration N
- 2: Initialize $\pi_0(\cdot | s)$ to be uniform; set $\mathcal{D}_0 = \emptyset, \mathcal{D}'_0 = \emptyset$
- 3: **for** episode $n = 1, \dots, N$ **do**
- 4: Collect a tuple $(s, a, s', a', \tilde{s})$ with

$$s \sim d_{P^*}^{\pi_{n-1}}, a \sim U(\mathcal{A}), s' \sim P^*(\cdot | s, a), a' \sim U(\mathcal{A}), \tilde{s} \sim P^*(\cdot | s', a')$$

- 5: Update datasets by adding triples (s, a, s') and (s', a', \tilde{s}) :

$$\mathcal{D}_n = \mathcal{D}_{n-1} + \{(s, a, s')\}, \quad \mathcal{D}'_n = \mathcal{D}'_{n-1} + \{(s', a', \tilde{s})\}$$

- 6: Learn representation via ERM (i.e., MLE):

$$\hat{P}_n := (\hat{\mu}_n, \hat{\phi}_n) = \arg \max_{(\mu, \phi) \in \mathcal{M}} \mathbb{E}_{\mathcal{D}_n + \mathcal{D}'_n} \left[\ln \mu^\top(s') \phi(s, a) \right]$$

- 7: Update empirical covariance matrix $\hat{\Sigma}_n = \sum_{s, a \in \mathcal{D}_n} \hat{\phi}_n(s, a) \hat{\phi}_n(s, a)^\top + \lambda_n I$
- 8: Set the exploration bonus:

$$\hat{b}_n(s, a) := \min \left(\alpha_n \sqrt{\hat{\phi}_n(s, a)^\top \hat{\Sigma}_n^{-1} \hat{\phi}_n(s, a)}, 2 \right) \quad (1)$$

- 9: Update policy $\pi_n = \arg \max_{\pi} V_{\hat{P}_n, r + \hat{b}_n}^{\pi}$
 - 10: **end for**
 - 11: **Return** π_1, \dots, π_N
-

PAC-Bound of Rep-UCB in low-rank MDP

Assume trajectory-reward is normalized in $[0,1]$. W/ high probability, it finds an ϵ near optimal policy, with # of samples:

$$\tilde{O} \left(\frac{d^4 A^2}{\epsilon^2 (1-\gamma)^5} \cdot \ln(|\Gamma||\Phi|) \right)$$

For reference, prior SOTA
FLAMBE has the following bound:

$$\tilde{O} \left(\frac{d^7 A^9}{\epsilon^{10} (1-\gamma)^{22}} \cdot \ln(|\Gamma||\Phi|) \right)$$

Extension to offline RL

- We note that this algorithm can be extended to offline setting where the set D is pre-collected.
- Similar to other offline theories, a dataset coverage coefficient is needed

Coverage condition of the offline data

A comparator policy π is covered by offline data if the **relative condition number** is bounded:

$$C_{\pi^*} := \max_x \frac{x^\top (\mathbb{E}_{s,a \sim d^\pi} \phi^*(s,a) \phi^*(s,a)^\top) x}{x^\top (\mathbb{E}_{s,a \sim d^{\pi_b}} \phi^*(s,a) \phi^*(s,a)^\top) x} < \infty$$

Note coverage is wrt true representation only!

Goal is to learn **robustly**, i.e., as long as there is a high quality policy that is covered by d^{π_b} , we want to compete against it!

- However, the proof techniques are basically the same with the online algorithm.

Summary

1. Improved online Representation Learning algorithm for low-rank MDP:

Oracle-efficient + tight sample complexity

2. New offline RL algorithm for low-rank MDP:

Partial coverage + Oracle-efficient

Rep-UCB / LCB: <https://arxiv.org/pdf/2110.04652.pdf>

Proof Sketch

- Now we go over the proof sketch for this algorithm

Algorithm 1 UCB-driven representation learning, exploration, and exploitation (REP-UCB)

- 1: **Input:** Regularizer λ_n , parameter α_n , Models $\mathcal{M} = \{(\mu, \phi) : \mu \in \Psi, \phi \in \Phi\}$, Iteration N
- 2: Initialize $\pi_0(\cdot | s)$ to be uniform; set $\mathcal{D}_0 = \emptyset, \mathcal{D}'_0 = \emptyset$
- 3: **for** episode $n = 1, \dots, N$ **do**
- 4: Collect a tuple $(s, a, s', a', \tilde{s})$ with

$$s \sim d_{P^*}^{\pi_{n-1}}, a \sim U(\mathcal{A}), s' \sim P^*(\cdot | s, a), a' \sim U(\mathcal{A}), \tilde{s} \sim P^*(\cdot | s', a')$$

- 5: Update datasets by adding triples (s, a, s') and (s', a', \tilde{s}) :

$$\mathcal{D}_n = \mathcal{D}_{n-1} + \{(s, a, s')\}, \quad \mathcal{D}'_n = \mathcal{D}'_{n-1} + \{(s', a', \tilde{s})\}$$

- 6: Learn representation via ERM (i.e., MLE):

$$\hat{P}_n := (\hat{\mu}_n, \hat{\phi}_n) = \arg \max_{(\mu, \phi) \in \mathcal{M}} \mathbb{E}_{\mathcal{D}_n + \mathcal{D}'_n} \left[\ln \mu^\top(s') \phi(s, a) \right]$$

- 7: Update empirical covariance matrix $\hat{\Sigma}_n = \sum_{s, a \in \mathcal{D}_n} \hat{\phi}_n(s, a) \hat{\phi}_n(s, a)^\top + \lambda_n I$
- 8: Set the exploration bonus:

$$\hat{b}_n(s, a) := \min \left(\alpha_n \sqrt{\hat{\phi}_n(s, a)^\top \hat{\Sigma}_n^{-1} \hat{\phi}_n(s, a)}, 2 \right) \quad (1)$$

- 9: Update policy $\pi_n = \arg \max_{\pi} V_{\hat{P}_n, r + \hat{b}_n}^{\pi}$
 - 10: **end for**
 - 11: **Return** π_1, \dots, π_N
-

Define ρ_n and ρ'_n as the marginal distributions for s in $\mathcal{D}_n, \mathcal{D}'_n$, respectively.

Therefore, we have: for (s, a, s') in \mathcal{D}_n , $(s, a, s') \sim \rho_n(s)U(a)P^*(s' | s, a)$

Useful auxiliary lemmas

• Capture model error

Lemma 18 (MLE guarantee). For a fixed episode n , with probability $1 - \delta$,

$$\mathbb{E}_{s \sim \{0.5\rho_n + 0.5\rho'_n\}, a \sim U(\mathcal{A})} [\|\hat{P}_n(\cdot | s, a) - P^*(\cdot | s, a)\|_1^2] \lesssim \zeta, \quad \zeta := \frac{\ln(|\mathcal{M}|/\delta)}{n}.$$

As a straightforward corollary, with probability $1 - \delta$,

$$\forall n \in \mathbb{N}^+, \mathbb{E}_{s \sim \{0.5\rho_n + 0.5\rho'_n\}, a \sim U(\mathcal{A})} [\|\hat{P}_n(\cdot | s, a) - P^*(\cdot | s, a)\|_1^2] \lesssim 0.5\zeta_n, \quad \zeta_n := \frac{\ln(|\mathcal{M}|n/\delta)}{n}.$$

Algorithm 1 UCB-driven representation learning, exploration, and exploitation (REP-UCB)

- 1: **Input:** Regularizer λ_n , parameter α_n , Models $\mathcal{M} = \{(\mu, \phi) : \mu \in \Psi, \phi \in \Phi\}$, Iteration N
- 2: Initialize $\pi_0(\cdot | s)$ to be uniform; set $\mathcal{D}_0 = \emptyset, \mathcal{D}'_0 = \emptyset$
- 3: **for** episode $n = 1, \dots, N$ **do**
- 4: Collect a tuple $(s, a, s', a', \tilde{s})$ with

$$s \sim d_{P^*}^{\pi_{n-1}}, a \sim U(\mathcal{A}), s' \sim P^*(\cdot | s, a), a' \sim U(\mathcal{A}), \tilde{s} \sim P^*(\cdot | s', a')$$

- 5: Update datasets by adding triples (s, a, s') and (s', a', \tilde{s}) :

$$\mathcal{D}_n = \mathcal{D}_{n-1} + \{(s, a, s')\}, \quad \mathcal{D}'_n = \mathcal{D}'_{n-1} + \{(s', a', \tilde{s})\}$$

- 6: Learn representation via ERM (i.e., MLE):

$$\hat{P}_n := (\hat{\mu}_n, \hat{\phi}_n) = \arg \max_{(\mu, \phi) \in \mathcal{M}} \mathbb{E}_{\mathcal{D}_n + \mathcal{D}'_n} [\ln \mu^\top(s') \phi(s, a)]$$

- 7: Update empirical covariance matrix $\hat{\Sigma}_n = \sum_{s, a \in \mathcal{D}_n} \hat{\phi}_n(s, a) \hat{\phi}_n(s, a)^\top + \lambda_n I$
- 8: Set the exploration bonus:

$$\hat{b}_n(s, a) := \min \left(\alpha_n \sqrt{\hat{\phi}_n(s, a)^\top \hat{\Sigma}_n^{-1} \hat{\phi}_n(s, a)}, 2 \right) \quad (1)$$

- 9: Update policy $\pi_n = \arg \max_{\pi} V_{\hat{P}_n, r + \hat{b}_n}^\pi$
 - 10: **end for**
 - 11: **Return** π_1, \dots, π_N
-

Useful auxiliary lemmas

• Concentration for covariance

Lemma 11 (Concentration of the bonus term). *Set $\lambda_n = \Theta(d \ln(n|\Phi|/\delta))$ for any n . Define*

$$\Sigma_{\rho_n, \phi} = n \mathbb{E}_{s \sim \rho_n, a \sim U(\mathcal{A})} [\phi(s, a) \phi^\top(s, a)] + \lambda_n I, \quad \hat{\Sigma}_{n, \phi} = \sum_{i=0}^{n-1} \phi(s^{(i)}, a^{(i)}) \phi^\top(s^{(i)}, a^{(i)}) + \lambda_n I.$$

With probability $1 - \delta$, we have

$$\forall n \in \mathbb{N}^+, \forall \phi \in \Phi, c_1 \|\phi(s, a)\|_{\Sigma_{\rho_n \times U(\mathcal{A}), \phi}^{-1}} \leq \|\phi(s, a)\|_{\hat{\Sigma}_{n, \phi}^{-1}} \leq c_2 \|\phi(s, a)\|_{\Sigma_{\rho_n \times U(\mathcal{A}), \phi}^{-1}}.$$

Algorithm 1 UCB-driven representation learning, exploration, and exploitation (REP-UCB)

- 1: **Input:** Regularizer λ_n , parameter α_n , Models $\mathcal{M} = \{(\mu, \phi) : \mu \in \Psi, \phi \in \Phi\}$, Iteration N
- 2: Initialize $\pi_0(\cdot | s)$ to be uniform; set $\mathcal{D}_0 = \emptyset, \mathcal{D}'_0 = \emptyset$
- 3: **for** episode $n = 1, \dots, N$ **do**
- 4: Collect a tuple (s, a, s', a', \bar{s}) with

$$s \sim d_{P^*}^{n-1}, a \sim U(\mathcal{A}), s' \sim P^*(\cdot | s, a), a' \sim U(\mathcal{A}), \bar{s} \sim P^*(\cdot | s', a')$$

- 5: Update datasets by adding triples (s, a, s') and (s', a', \bar{s}) :

$$\mathcal{D}_n = \mathcal{D}_{n-1} + \{(s, a, s')\}, \quad \mathcal{D}'_n = \mathcal{D}'_{n-1} + \{(s', a', \bar{s})\}$$

Learn representation via ERM (i.e., MLE):

$$\hat{P}_n := (\hat{\mu}_n, \hat{\phi}_n) = \arg \max_{(\mu, \phi) \in \mathcal{M}} \mathbb{E}_{\mathcal{D}_n + \mathcal{D}'_n} [\ln \mu^\top(s') \phi(s, a)]$$

Update empirical covariance matrix $\hat{\Sigma}_n = \sum_{s, a \in \mathcal{D}_n} \hat{\phi}_n(s, a) \hat{\phi}_n(s, a)^\top + \lambda_n I$
Set the exploration bonus:

$$\hat{b}_n(s, a) := \min \left(\alpha_n \sqrt{\hat{\phi}_n(s, a)^\top \hat{\Sigma}_n^{-1} \hat{\phi}_n(s, a)}, 2 \right) \quad (1)$$

- 9: Update policy $\pi_n = \arg \max_{\pi} V_{\hat{P}_n, r + \hat{b}_n}^{\pi}$
 - 10: **end for**
 - 11: **Return** π_1, \dots, π_N
-

Performance difference lemma is always useful

- Substitute the two lemmas into the performance difference lemma then we got

$$\begin{aligned} V_{\hat{P}_n, r + \hat{b}_n}^\pi - V_{P^*, r}^\pi &= (1 - \gamma)^{-1} \mathbb{E}_{s, a \sim d_{\hat{P}_n}^\pi} \left[\hat{b}_n(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_n(\cdot | s, a)} V_{P^*}^\pi(s') - \gamma \mathbb{E}_{s' \sim P^*(\cdot | s, a)} V_{P^*}^\pi(s') \right] \\ &\geq (1 - \gamma)^{-1} \mathbb{E}_{s, a \sim d_{\hat{P}_n}^\pi} \left[\hat{b}_n(s, a) - \|\hat{P}_n(\cdot | s, a) - P^*(\cdot | s, a)\|_1 \right], \end{aligned}$$

$$\hat{b}_n(s, a) := \min \left(\alpha_n \sqrt{\hat{\phi}_n(s, a)^\top \hat{\Sigma}_n^{-1} \hat{\phi}_n(s, a)}, 2 \right)$$

Remark:

- The bonus term is defined using elliptical potential function under $\hat{\phi}$
- We need to bound $\|\hat{P} - P^*\|$ with the elliptical potential function, this is the following lemma

Relate model error with potential function

Lemma 12 (One-step back inequality for the learned model). *Take any $g \in \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ such that $\|g\|_\infty \leq B$. We condition on the event where the MLE guarantee (17):*

$$\mathbb{E}_{s \sim \rho_n, a \sim U(\mathcal{A})} [f_n(s, a)] \lesssim \zeta_n,$$

holds. Then, for any policy π , we have

$$\begin{aligned} & |\mathbb{E}_{(s,a) \sim d_{\hat{P}_n}^\pi} \{g(s, a)\}| \\ & \leq \mathbb{E}_{(\tilde{s}, \tilde{a}) \sim d_{\hat{P}_n}^\pi} \|\hat{\phi}_n(\tilde{s}, \tilde{a})\|_{\Sigma_{\rho_n \times U(\mathcal{A}), \hat{\phi}_n}^{-1}} \sqrt{\{n|\mathcal{A}| \mathbb{E}_{s \sim \rho'_n, a \sim U(\mathcal{A})} [g^2(s, a)]\} + B^2 \lambda_n d + nB^2 \zeta_n} \\ & + \sqrt{(1 - \gamma)|\mathcal{A}| \mathbb{E}_{s \sim \rho_n, a \sim U(\mathcal{A})} [g^2(s, a)]}. \end{aligned}$$

Recall $\Sigma_{\rho_n \times U(\mathcal{A}), \hat{\phi}_n} = n \mathbb{E}_{s \sim \rho_n, a \sim U(\mathcal{A})} [\hat{\phi}_n(s, a) \hat{\phi}_n^\top(s, a)] + \lambda_n I$.

Remark:

- Instantiate this lemma with $g = \|\hat{P}(\cdot |s, a) - P^*(\cdot |s, a)\|_1$ and $B = 2$
- This lemma is quite hard to prove, it requires both formulations of D_n, D_n'
- I think it's an issue to extend this result to risk sensitive setting

Back to regret: Optimism

$$\begin{aligned}
 V_{\hat{P}_n, r + \hat{b}_n}^\pi - V_{P^*, r}^\pi &= (1 - \gamma)^{-1} \mathbb{E}_{s, a \sim d_{\hat{P}_n}^\pi} \left[\hat{b}_n(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_n(\cdot | s, a)} V_{P^*}^\pi(s') - \gamma \mathbb{E}_{s' \sim P^*(\cdot | s, a)} V_{P^*}^\pi(s') \right] \\
 &\geq (1 - \gamma)^{-1} \mathbb{E}_{s, a \sim d_{\hat{P}_n}^\pi} \left[\hat{b}_n(s, a) - \|\hat{P}_n(\cdot | s, a) - P^*(\cdot | s, a)\|_1 \right],
 \end{aligned}$$

Use the previous lemma to
Relate model error with
potential function

$$\hat{b}_n(s, a) := \min \left(\alpha_n \sqrt{\hat{\phi}_n(s, a)^\top \hat{\Sigma}_n^{-1} \hat{\phi}_n(s, a)}, 2 \right)$$

$$\min \left(\alpha_n \|\hat{\phi}_n(s, a)\|_{\Sigma_{\rho_n \times U(\mathcal{A}), \hat{\phi}_n}^{-1}} + \sqrt{(1 - \gamma)|\mathcal{A}|\zeta_n}, 2 \right)$$

Finally we have:

Lemma 5 (Almost Optimism at the Initial State Distribution). *Set the parameters as in Theorem 4. With probability $1 - \delta$,*

$$\forall n \in [1, \dots, N], \forall \pi \in \Pi, V_{\hat{P}_n, r + \hat{b}_n}^\pi - V_{P^*, r}^\pi \geq -c_1 \sqrt{\frac{|\mathcal{A}| \ln(|\mathcal{M}|n/\delta)(1 - \gamma)^{-1}}{n}}.$$

Remark: From this result, we can see that the bonus term b is designed to provide near-optimism.

Regret

Optimism

$$\begin{aligned}
 & V_{P^*,r}^{\pi^*} - V_{P^*,r}^{\pi_n} \\
 & \leq V_{\hat{P}_n,r+\hat{b}_n}^{\pi^*} - V_{P^*,r}^{\pi_n} + \sqrt{|\mathcal{A}|\zeta_n(1-\gamma)^{-1}} \quad (\text{Lemma 8}) \\
 & \leq V_{\hat{P}_n,r+\hat{b}_n}^{\pi_n} - V_{P^*,r}^{\pi_n} + \sqrt{|\mathcal{A}|\zeta_n(1-\gamma)^{-1}} \quad (\pi_n = \arg \max_{\pi} V_{\hat{P}_n,r+\hat{b}_n}^{\pi}) \\
 & = (1-\gamma)^{-1} \mathbb{E}_{(s,a) \sim d_{P^*}^{\pi_n}} [\hat{b}_n(s,a) + \gamma \mathbb{E}_{\hat{P}_n(s'|s,a)} [V_{\hat{P}_n,r+\hat{b}_n}^{\pi_n}(s')] - \gamma \mathbb{E}_{P^*(s'|s,a)} [V_{\hat{P}_n,r+\hat{b}_n}^{\pi_n}(s')]] + \sqrt{|\mathcal{A}|\zeta_n(1-\gamma)^{-1}}.
 \end{aligned}$$

Relate model error with potential function under ϕ^*

Algorithm 1 UCB-driven representation learning, exploration, and exploitation (REP-UCB)

- 1: **Input:** Regularizer λ_n , parameter α_n , Models $\mathcal{M} = \{(\mu, \phi) : \mu \in \Psi, \phi \in \Phi\}$, Iteration N
- 2: Initialize $\pi_0(\cdot | s)$ to be uniform; set $\mathcal{D}_0 = \emptyset, \mathcal{D}'_0 = \emptyset$
- 3: **for** episode $n = 1, \dots, N$ **do**
- 4: Collect a tuple $(s, a, s', a', \tilde{s})$ with

$$s \sim d_{P^*}^{\pi_{n-1}}, a \sim U(\mathcal{A}), s' \sim P^*(\cdot | s, a), a' \sim U(\mathcal{A}), \tilde{s} \sim P^*(\cdot | s', a')$$
- 5: Update datasets by adding triples (s, a, s') and (s', a', \tilde{s}) :

$$\mathcal{D}_n = \mathcal{D}_{n-1} + \{(s, a, s')\}, \quad \mathcal{D}'_n = \mathcal{D}'_{n-1} + \{(s', a', \tilde{s})\}$$
- 6: Learn representation via ERM (i.e., MLE):

$$\hat{P}_n := (\hat{\mu}_n, \hat{\phi}_n) = \arg \max_{(\mu, \phi) \in \mathcal{M}} \mathbb{E}_{\mathcal{D}_n + \mathcal{D}'_n} [\ln \mu^\top(s') \phi(s, a)]$$
- 7: Update empirical covariance matrix $\hat{\Sigma}_n = \sum_{s,a \in \mathcal{D}_n} \hat{\phi}_n(s, a) \hat{\phi}_n(s, a)^\top + \lambda_n I$
- 8: Set the exploration bonus:

$$\hat{b}_n(s, a) := \min \left(\alpha_n \sqrt{\hat{\phi}_n(s, a)^\top \hat{\Sigma}_n^{-1} \hat{\phi}_n(s, a)}, 2 \right) \quad (1)$$
- 9: Update policy $\pi_n = \arg \max_{\pi} V_{\hat{P}_n,r+\hat{b}_n}^{\pi}$
- 10: **end for**
- 11: **Return** π_1, \dots, π_N

Policy update

Finally we get some terms like: $\|\phi^*\|_{\Sigma^{-1}}$, which are controlled by some data-dependent quantities. Q.E.D.

Back to risk sensitive RL

- For each threshold b , we have

$$V^\pi(s, b) = E [\max(0, b - \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t)) | s_0 = s, \pi, b]$$
$$J^*(b) = \min_{\pi} J(b) = \min_{\pi} E_{s \sim d_0} V^\pi(s, b)$$

The final goal is to maximize CVaR, i.e., $\max_b (b - \frac{1}{\tau} J^*(b))$

Another loop for threshold b ?

Algorithm 1 UCB-driven CVaR learning

Input: Regularizer λ_n , parameter α_n , models $\mathcal{M} = \{(\mu, \phi) : \mu \in \Psi, \phi \in \Phi\}$, number of iterations N .

- 1: **for** $b \in \mathcal{B}$ **do**
- 2: **for** episode $n = 1, \dots, N$ **do**
- 3: Collect a tuple $(s, a, s', a', \tilde{s})$
- 4: Update datasets: $\mathcal{D}_n = \mathcal{D}_{n-1} + \{(s, a, s')\}$, $\mathcal{D}'_n = \mathcal{D}'_{n-1} + \{(s', a', \tilde{s})\}$
- 5: Learn representations via MLE

$$\hat{P}_n = \arg \max_{\mu, \phi \in \mathcal{M}} \mathbb{E} \ln \mu^\top(s') \phi(s, a)$$

- 6: Update empirical covariance $\hat{\Sigma}_n = \sum_{s, a \in \mathcal{D}_n} \hat{\phi}_n(s, a) \hat{\phi}_n(s, a)^\top + \lambda_n I_n$
- 7: Set the exploration bonus:

$$\hat{b}_n = ???$$

- 8: Update policy $\pi_n = \arg \max_{\pi} V_{\hat{P}_n, r + \hat{b}_n}$
 - 9: **end for**
 - 10: Store $\pi_1(|b), \dots, \pi_N(|b)$
 - 11: **end for**
- Output:** $b = \arg \max_b (b - \frac{1}{\tau} \mathbb{E}_{s \sim d_0} V^\pi(s, b))$ and $\pi_1(|b), \dots, \pi_N(|b)$
-

Final remarks

- Concentrations should still hold
- Performance difference lemma changes because of the threshold parameter, recall that $V^\pi(s, b) = E [\max(0, b - \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t)) | s_0 = s, \pi, b]]$
- Optimism results does not apply, need new techniques
- The relation between model error and potential function is built upon the MDP transition which still holds, might not change much but need double check.