

VIPeR: Provably Efficient Algorithm for Offline RL with Neural Function Approximation

Thanh Nguyen-Tang & Raman Arora

Department of Computer Science, Johns Hopkins University

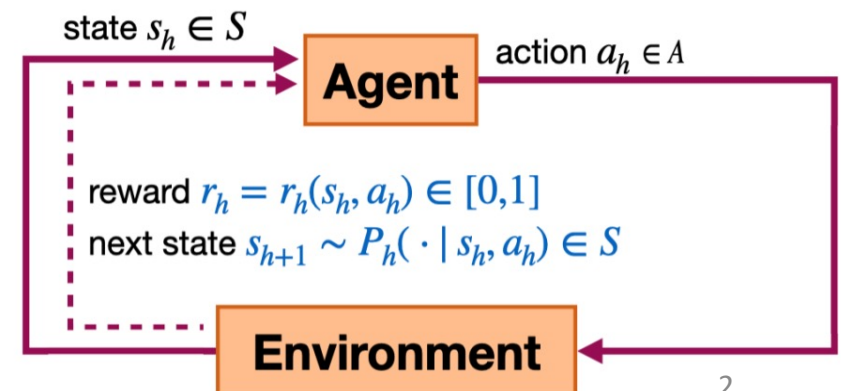


Episodic MDP

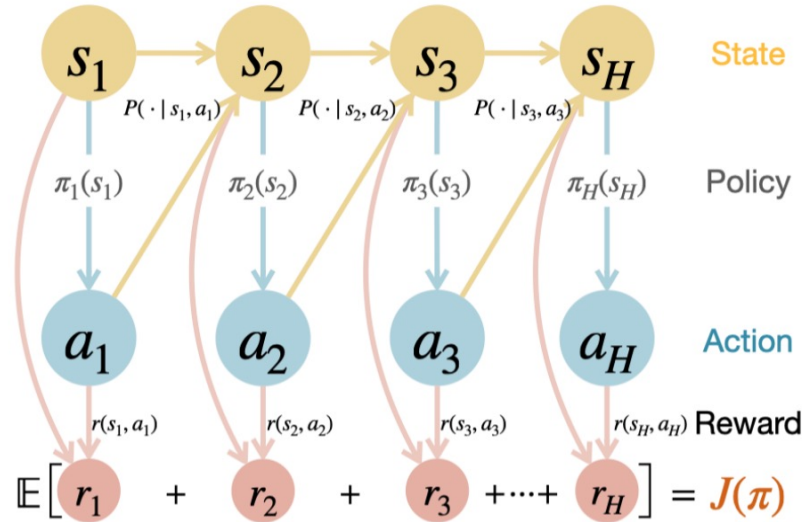
- Episodic time-inhomogeneous Markov decision process

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, P, r, d_1)$$

- State space \mathcal{S} : finite but exponentially large
- Action space \mathcal{A} : finite but exponentially large
- Episode length H : finite
 - Agent interacts with MDP for H steps and then restart the episode
- Unknown Transition kernels $P = (P_1, \dots, P_H)$, where $P_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$
- Unknown Reward functions $r = (r_1, \dots, r_H)$, where $r_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$
- Initial state distribution $d_1 \in \Delta(\mathcal{S})$



Episodic MDP



- A policy $\pi = \{\pi_h\}_{h \in [H]}$ where $\pi_h: \mathcal{S} \rightarrow \Delta(\mathcal{A})$
- Action-value functions:

$$Q_h^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{i=1}^H r_i \mid (s_h, a_h) = (s, a) \right]$$

- Value functions

$$V_h^\pi(s) = \mathbb{E}_\pi \left[\sum_{i=1}^H r_i \mid s_h = s \right]$$

- The optimal policy π^* maximizes V_1^π

Offline RL

- Offline dataset: collected a priori, $\mathcal{D} = \{(s_h^t, a_h^t, r_h^t)\}_{h \in [H]}^{t \in [K]}$
 - $a_h^t \sim \mu_h(\cdot | s_h^t)$, $s_{h+1}^t \sim P_h(\cdot | s_h^t, a_h^t)$
 - μ is the behavior policy
 - K : # number of episodes
- No further interactions with MDP
- Learning objective:

$$\text{SubOpt}(\hat{\pi}; s_1) = V_1^*(s_1) - V_1^{\hat{\pi}}(s_1)$$

where $\hat{\pi} = \text{OfflineRLAlgo}(\mathcal{D}, \mathcal{F})$, \mathcal{F} is some function class (e.g., neural networks)

PEVI algorithm (Jin et al., ICML'21)

- **Low-rank MDPs:**

$$r_{h(s,a)} = \phi_h(s,a)^T w_h, P_h(s'|s,a) = \phi_h(s,a)^T v_h(s')$$

- **Pessimism design: Computing LCB (lower confidence bound)**

$$\hat{Q}_h(s,a) = \phi_h(s,a)^T \hat{w}_h - \beta \|\phi_h(s,a)\|_{\Lambda_h^{-1}}$$

- $\Lambda_h = \lambda I + \sum_{k=1}^K \phi_h(s_h^k, a_h^k) \phi_h(s_h^k, a_h^k)^T$
- \hat{w}_h : ERM

- Then, extract $\hat{\pi}_h$ greedy w.r.t. \hat{Q}_h
- PEVI is both **statistically efficient** and **computationally efficient** in low-rank MDPs
- PEVI require data coverage **only over an optimal policy**

[Jin et al., ICML'21] “Is pessimism provably efficient for offline rl?”

Offline RL with Neural Function Approximation

However, in many **practical** settings, MDPs do **not** admit any **linear** structures

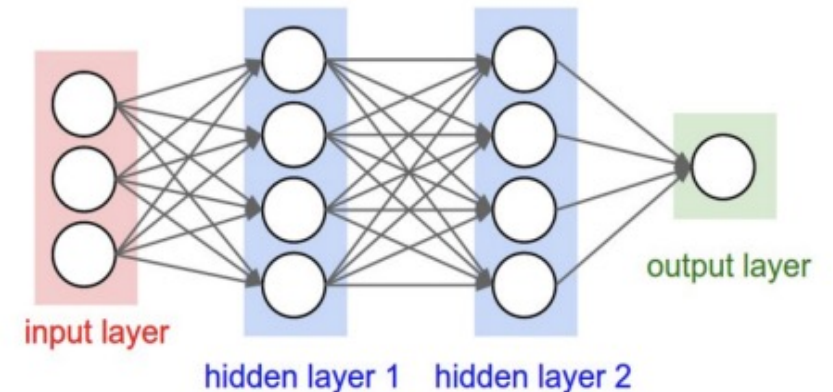
→ Desirable to use differentiable function approximator such as **neural networks**

Background: Neural networks

- Action-value functions are approximated by a two-layer neural net

$$f(\mathbf{x}; W) = \frac{1}{\sqrt{m}} \sum_{i=1}^m b_i \cdot \text{ReLU}(\mathbf{w}_i^T \mathbf{x})$$

- $\mathbf{x} = (s, a) \in \mathcal{X} = \mathcal{S} \times \mathcal{A} \in \mathbb{R}^d$
- m : network width
- $b_i \sim \text{Unif}(\{-1, 1\})$: fixed during training
- $W = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m) \in \mathbb{R}^{md}$: trainable
- W_0 : initialization of W
 - $\mathbf{w}_i \sim \mathcal{N}(0, \frac{I_d}{d})$



NeuraLCB (Nguyen-Tang et al., ICLR)

LCB + neural function approximation

- Use neural network $f(s, a; W)$ to approximate the action-value function

- Computing **LCB**:

$$\hat{Q}_h(s, a) = f(s, a; \hat{W}) - \beta \|\nabla_W f(s, a; \hat{W})\|_{\Lambda_h^{-1}}$$

- \hat{W} : found by gradient descent
- $\nabla_W f(s, a; \hat{W})$: the gradient of the neural network w.r.t. its parameters
- $\Lambda_h = \lambda I + \sum_{k=1}^K \nabla_W f(s_h^k, a_h^k; \hat{W}) \left(\nabla_W f(s_h^k, a_h^k; \hat{W}) \right)^T$ is the empirical covariance matrix
- Complexity of computing LCB: $(\# \text{ of network parameters})^2$

VIPeR: Implicit pessimism via reward perturbing

- End of Episode: $\tilde{Q}_{H+1} \leftarrow 0$
- Bootstrapping (backward induction): for $h = H, H - 1, \dots, 1$
 - Perturb the rewards with independent Gaussian noises $\xi_h^{k,i} \sim \mathcal{N}(0, \sigma^2)$ to create **M perturbed datasets** $\tilde{\mathcal{D}}_h^1, \tilde{\mathcal{D}}_h^2, \dots, \tilde{\mathcal{D}}_h^M$

$$\tilde{\mathcal{D}}_h^i = \{ \underbrace{(s_h^k, a_h^k)}_{\text{input}}, \underbrace{r_h^k + \tilde{V}_{h+1}(s_{h+1}^k) + \xi_h^{k,i}}_{\text{perturbed output}} \}_{k \in [K]}$$

- Train a different neural network in each $\tilde{\mathcal{D}}_h^i$ to get $f(s, a; W^i)$
- Pessimistically estimate
$$\tilde{Q}_h(s, a) = \min_{i \in [M]} f(s, a; W^i)$$
- Optimize $\tilde{\pi}_h(s) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \tilde{Q}_h(s, a)$ and $\tilde{V}_h(s) = \max_a \tilde{Q}_h(s, a)$
- Move to step **h-1** and repeat

Value suboptimality of VIPeR

- Conditions:
 - Adaptively collected data with coverage of only an optimal policy
 - Completeness assumption
 - The ensemble size M is **polylogarithmically** large
 - The noise level $\sigma = \tilde{O}(H d_{\text{eff}})$
 - Small learning rate and sufficiently large training iterations
 - The network width m is polynomially large
- Then,

$$\text{SubOpt}(\tilde{\pi}; s_1) \lesssim \sigma \mathbb{E}_{\pi^*} \left[\sum_{h=1}^H \|\nabla_w f(s_h, a_h; \hat{W})\|_{\Lambda_h^{-1}} \right]$$

- Same guarantee as **NeuraLCB**, but **no** need to compute LCB

Comparison

work	bound	i.i.d?	explorative data?	finite spectrum?	matrix inverse?	opt
Jin et al. (2021)	$\tilde{O}\left(\frac{d_{lin}^{3/2} H^2}{\sqrt{K}}\right)$	no	yes	yes	yes	analytical
Yang et al. (2020)	$\tilde{O}\left(\frac{H^2 \sqrt{\tilde{d}^2 + \tilde{d}\tilde{n}}}{\sqrt{K}}\right)$	no	–	no	yes	oracle
Xu & Liang (2022)	$\tilde{O}\left(\frac{\tilde{d} H^2}{\sqrt{K}}\right)$	yes	yes	yes	yes	oracle
This work	$\tilde{O}\left(\frac{\kappa H^{5/2} \tilde{d}}{\sqrt{K}}\right)$	no	no	no	no	GD

Table 1: SOTA results for offline RL with function approximation. The third and fourth columns ask if the corresponding result needs the data to be i.i.d, and well-explored, resp.; the fifth column asks if the induced RKHS needs to have a finite spectrum, the sixth column asks if the algorithm needs to inverse a covariance matrix and the last column presents the optimizer being used. Here \tilde{n} is the log covering number.

Experiment: Neural Contextual Bandits

- **LinLCB**: LCB + linear
- **Lin-VIPeR**: Reward perturbing + linear
- **NeuralGreedy**: FQI + neural
- **NeuraLCB**: LCB + neural
- **Neural-VIPeR**: reward perturbing + neural

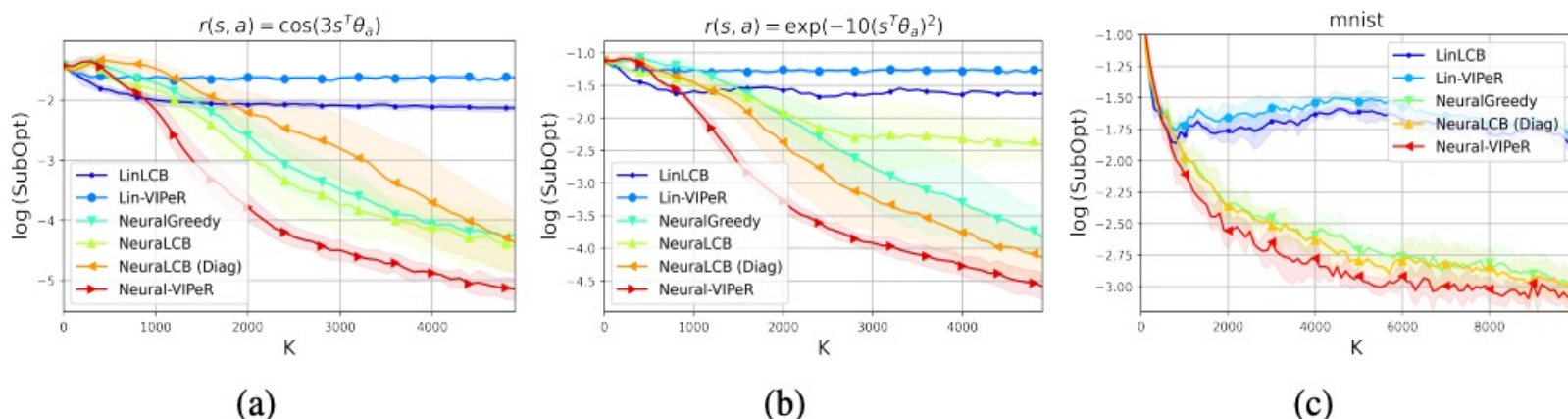


Figure 3: Sub-optimality (on log-scale) vs. sample size (K) for neural contextual bandits with following reward functions: (a) $r(s, a) = \cos(3s^T \theta_a)$, (b) $r(s, a) = \exp(-10(s^T \theta_a)^2)$, and (c) MNIST.

Runtime efficiency

- **NeuraLCB** spends $O(K^2)$ time in action selection
- **Neural-VIPeR** spends $O(1)$ in action selection

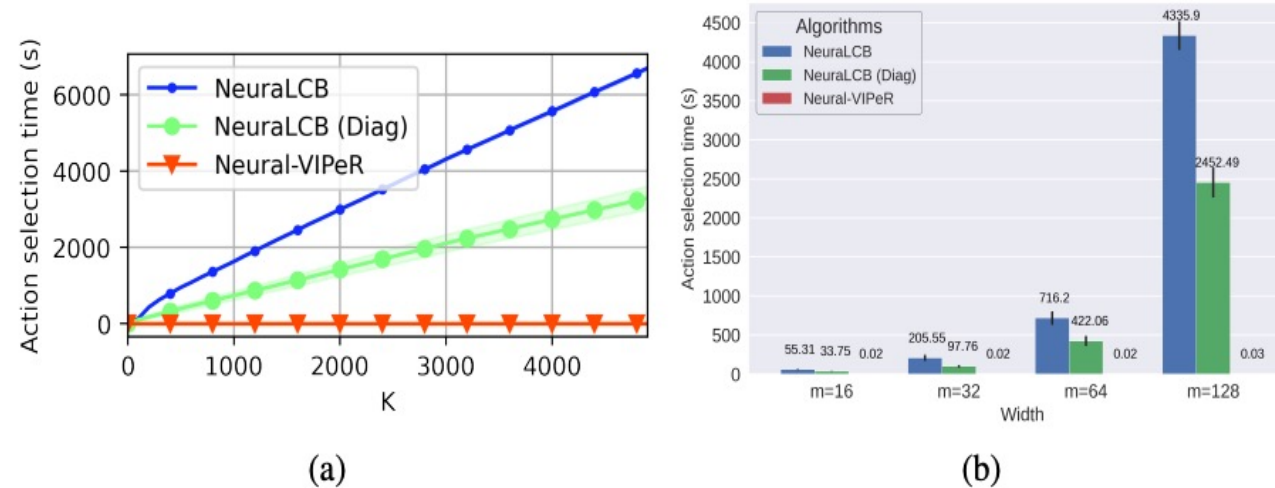


Figure 4: Elapsed time (in seconds) for action selection in the contextual bandits problem with $r(s, a) = 10(s^T \theta_a)^2$: (a) Runtime of action selection versus the number of (offline) data points K , and (b) runtime of action selection versus the network width m (for $K = 500$).

Controlling pessimism

M controls the level of pessimism and nicely correlates with the decrease of subopt in practice

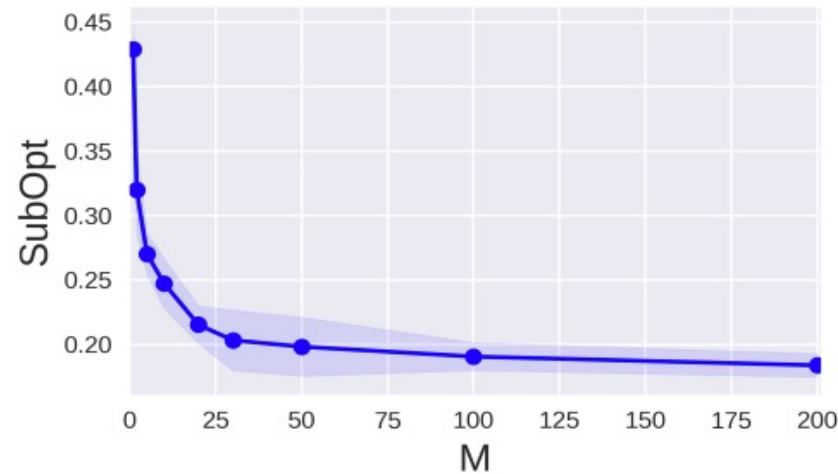


Figure 5: Sub-optimality of NeuralPER versus different values of M .

Result in D4RL

- **BEAR** (Kumar et al., 2019): use MMD distance to constraint policy to the offline data
- **UWAC** (Wu et al., 2021): improves BEAR using dropout uncertainty
- **CQL** (Kumar et al., 2020) that minimizes Q-values of OOD actions
- **MOPO** (Yu et al., 2020): uses model-based uncertainty via ensemble dynamics
- **TD3-BC** (Fujimoto & Gu, 2021): uses adaptive behavior cloning
- **PBRL** (Bai et al., 2022): use uncertainty quantification via disagreement of bootstrapped Q-functions

		BEAR	UWAC	CQL	MOPO	TD3-BC	PBRL	VIPeR
Random	HalfCheetah	2.3 \pm 0.0	2.3 \pm 0.0	17.5 \pm 1.5	35.9 \pm 2.9	11.0 \pm 1.1	11.0 \pm 5.8	14.5 \pm 2.1
	Hopper	3.9 \pm 2.3	2.7 \pm 0.3	7.9 \pm 0.4	16.7 \pm 12.2	8.5 \pm 0.6	26.8 \pm 9.3	31.4 \pm 0.0
	Walker2d	12.8 \pm 10.2	2.0 \pm 0.4	5.1 \pm 1.3	4.2 \pm 5.7	1.6 \pm 1.7	8.1 \pm 4.4	20.5 \pm 0.5
Medium	HalfCheetah	43.0 \pm 0.2	42.2 \pm 0.4	47.0 \pm 0.5	73.1 \pm 2.4	48.3 \pm 0.3	57.9 \pm 1.5	58.5 \pm 1.1
	Hopper	51.8 \pm 4.0	50.9 \pm 4.4	53.0 \pm 28.5	38.3 \pm 34.9	59.3 \pm 4.2	75.3 \pm 31.2	99.4 \pm 6.2
	Walker2d	-0.2 \pm 0.1	75.4 \pm 3.0	73.3 \pm 17.7	41.2 \pm 30.8	83.7 \pm 2.1	89.6 \pm 0.7	89.6 \pm 1.2
Medium Replay	HalfCheetah	36.3 \pm 3.1	35.9 \pm 3.7	45.5 \pm 0.7	69.2 \pm 1.1	44.6 \pm 0.5	45.1 \pm 8.0	45.0 \pm 8.6
	Hopper	52.2 \pm 19.3	25.3 \pm 1.7	88.7 \pm 12.9	32.7 \pm 9.4	60.9 \pm 18.8	100.6 \pm 1.0	100.2 \pm 1.0
	Walker2d	7.0 \pm 7.8	23.6 \pm 6.9	81.8 \pm 2.7	73.7 \pm 9.4	81.8 \pm 5.5	77.7 \pm 14.5	83.1 \pm 4.2
Medium Expert	HalfCheetah	46.0 \pm 4.7	42.7 \pm 0.3	75.6 \pm 25.7	70.3 \pm 21.9	90.7 \pm 4.3	92.3 \pm 1.1	94.2 \pm 1.2
	Hopper	50.6 \pm 25.3	44.9 \pm 8.1	105.6 \pm 12.9	60.6 \pm 32.5	98.0 \pm 9.4	110.8 \pm 0.8	110.6 \pm 1.0
	Walker2d	22.1 \pm 44.9	96.5 \pm 9.1	107.9 \pm 1.6	77.4 \pm 27.9	110.1 \pm 0.5	110.1 \pm 0.3	109.8 \pm 0.5
Expert	HalfCheetah	92.7 \pm 0.6	92.9 \pm 0.6	96.3 \pm 1.3	81.3 \pm 21.8	96.7 \pm 1.1	92.4 \pm 1.7	97.4 \pm 0.9
	Hopper	54.6 \pm 21.0	110.5 \pm 0.5	96.5 \pm 28.0	62.5 \pm 29.0	107.8 \pm 7	110.5 \pm 0.4	110.8 \pm 0.4
	Walker2d	106.6 \pm 6.8	108.4 \pm 0.4	108.5 \pm 0.5	62.4 \pm 3.2	110.2 \pm 0.3	108.3 \pm 0.3	108.3 \pm 0.2
Average		38.78 \pm 10.0	50.41 \pm 2.7	67.35 \pm 9.1	53.3 \pm 16.3	67.55 \pm 3.8	74.37 \pm 5.3	78.2 \pm 1.9

Table 2: Average normalized score and standard deviation of all algorithms over five seeds in the Gym domain in the “v2” dataset of D4RL (Fu et al., 2020). The scores for all the baselines are from Table 1 of Bai et al. (2022). The highest scores are highlighted.

Summary

We now have a **provably efficient & computationally efficient algorithm** for **neural** function approximation with polynomial sample and runtime under *mild data coverage*

Algorithm: perturbed rewards + (Stochastic) gradient descent

Sample complexity: $\text{SubOpt}(\tilde{\pi}) = \tilde{O}\left(\frac{H^{2.5} \cdot \kappa \cdot d_{\text{eff}}}{\sqrt{K}}\right)$

- K : # samples
- H : horizon
- κ : single-policy concentration coefficient
- d_{eff} : effective dimension

Future Work

- Match the lower bound in terms of κ ?
- Avoid the need of training an ensemble of M models?
- General-purpose offline RL with any regression oracle?
- Can randomized value iteration obtain a horizon-free rate?
- Extension to the POMDP?
- Model-free posterior sampling for offline RL with optimal frequentist rates?

Thank you

More details at our paper: <https://openreview.net/forum?id=WOquZTLCBO1>