

POLICY EVALUATION BACKGROUND

Markov Reward Process

- S : State Space.
- A : Action Space.
- π : Policy $S \rightarrow A$.
- $P(s'|s, a)$: Transition Probability.
- $R(s)$: Reward function.

Value Function

- Long term payoff of a given policy.
- Def: $V^\pi(s) = E[r_1 + \gamma r_2 + \dots | s_1 = s]$.
- Linear Policy Evaluation:

$$\widehat{V}(s) = \phi(s)^\top \theta$$

- $\phi(s)$: d -dimensional feature vector.
- $\theta \in \mathbb{R}^d$: coefficients to be learned.

OBJECTIVE FUNCTION AND REFORMULATION

Data: trajectory of n steps: $(s_1, r_1, s_2, r_2, \dots, s_n, r_n)$ with $\phi_t = \phi(s_t)$.

Goal: Minimize Empirical MSPBE (Mean Square Projected Bellman Error):

$$\min_{\theta} \frac{1}{2} \|A\theta - b\|_{C^{-1}}^2 + \frac{\rho}{2} \|\theta\|_2^2$$

$$A = \frac{1}{n} \sum_{t=1}^n A_t = \phi_t (\phi_t - \phi_{t+1})^\top, b = \frac{1}{n} \sum_{t=1}^n b_t = r_t \phi_t, C = \frac{1}{n} \sum_{t=1}^n \phi_t \phi_t^\top$$

Challenges:

- The objective does NOT have the finite sum structure (preclude modern methods).
- Even computing the gradient requires $O(nd^2)$ complexity.

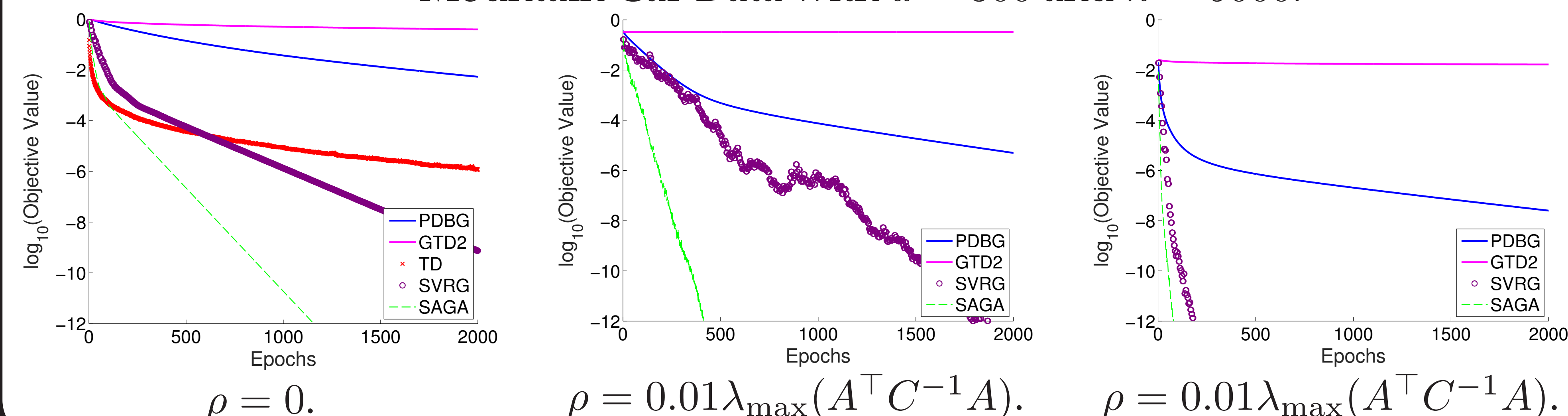
Key Observation: minimizing Regularized Empirical MSPBE is equivalent to

$$\begin{aligned} & \min_{\theta} \max_w \left\{ w^\top (b - A\theta) - \frac{1}{2} \|w\|_C^2 + \frac{\rho}{2} \|\theta\|_2^2 \right\} \\ & = \min_{\theta} \max_w \frac{1}{n} \sum_{t=1}^n \left\{ w^\top (b_t - A_t \theta) - \frac{1}{2} (w^\top C_t)^2 + \frac{\rho}{2} \|\theta\|_2^2 \right\} \end{aligned}$$

- Has the FINITE SUM structure.
- Stochastic gradient $B_t(\theta, w) = \begin{bmatrix} \rho\theta - A_t w \\ A_t \theta - b_t + C_t w \end{bmatrix}$ can be computed in $O(d)$ time and $B(\theta, w) = \sum_{t=1}^n B_t(\theta, w)$ can be computed in $O(nd)$ time.

EXPERIMENT ON MOUNTAIN CAR DATA

Mountain Car Data with $d = 300$ and $n = 5000$.



ALGORITHMS, COMPLEXITIES, EXTENSIONS

Algorithm: Primal-Dual-Batch-Gradient (PDBG)

Input: initial point (θ, w) , step sizes σ_θ and σ_w , and number of epochs M .

$$\text{For } i = 1 \text{ to } M: \begin{bmatrix} \theta \\ w \end{bmatrix} \leftarrow \begin{bmatrix} \theta \\ w \end{bmatrix} - \begin{bmatrix} \sigma_\theta & 0 \\ 0 & \sigma_w \end{bmatrix} B(\theta, w)$$

Algorithm: SVRG for Policy Evaluation

Input: initial point (θ, w) , step sizes $\{\sigma_\theta, \sigma_w\}$, outer iterations M , and inner iterations N .

For $m = 1$ **to** M :

Initialize $(\tilde{\theta}, \tilde{w}) = (\theta, w)$ and compute $B(\tilde{\theta}, \tilde{w})$.

For $j = 1$ **to** N :

Sample an index t_j from $\{1, \dots, n\}$ and update using variance-reduced gradient:

$$\begin{bmatrix} \theta \\ w \end{bmatrix} \leftarrow \begin{bmatrix} \theta \\ w \end{bmatrix} - \begin{bmatrix} \sigma_\theta & 0 \\ 0 & \sigma_w \end{bmatrix} \left(B_{t_j}(\theta, w) + B(\tilde{\theta}, \tilde{w}) - B_{t_j}(\tilde{\theta}, \tilde{w}) \right).$$

Algorithm: SAGA for Policy Evaluation

Input: initial point (θ, w) , step sizes σ_θ and σ_w , and number of iterations M .

Compute each $g_t = B_t(\theta, w)$ for $t = 1, \dots, n$ and $B = B(\theta, w) = \frac{1}{n} \sum_{t=1}^n g_t$.

For $m = 1$ **to** M :

Sample an index t_m from $\{1, \dots, n\}$, compute $h_{t_m} = B_{t_m}(\theta, w)$ and update:

$$\begin{bmatrix} \theta \\ w \end{bmatrix} \leftarrow \begin{bmatrix} \theta \\ w \end{bmatrix} - \begin{bmatrix} \sigma_\theta & 0 \\ 0 & \sigma_w \end{bmatrix} (B + h_{t_m} - g_{t_m}), B \leftarrow B + \frac{1}{n} (h_{t_m} - g_{t_m}), g_{t_m} \leftarrow h_{t_m}.$$

Algorithm	Total Complexity $O(\cdot)$	
SVRG / SAGA (this paper)	$nd(1 + \frac{\kappa_1}{n}) \log(1/\epsilon)$	$\kappa_1, \kappa_2, \kappa_3$ are condition numbers.
GTD2 (Sutton et al, 2009)	$d\kappa_2/\epsilon$	
PDBG (this paper)	$nd\kappa_3 \log(1/\epsilon)$	
LSTD (Boyan, 2002)	nd^2	

Our algorithms extend directly to:

- off-policy PE: data sampled from a different distribution, i.e., "covariate shift".
- PE with eligibility traces, which smooth the gap between MSPBE and MSE.

The only difference is in how A, B, C are formed.

CONCLUSION AND FUTURE DIRECTIONS

Conclusions:

1. A saddle-point formulation of (batch) policy evaluation.
2. First-order algorithms with linear convergence rate

Future Directions:

1. Extension to nonlinear value-function approximation.
2. Extension to the controlled case (i.e., policy optimization).