

Fast Multi-agent Temporal-difference Learning: Homotopy Stochastic Primal-dual Method

Dongsheng Ding, Xiaohan Wei, Zhuoran Yang, Zhaoran Wang, Mihailo R. Jovanović
dongshed@usc.edu, xiaohanw@usc.edu, zy6@princeton.edu, zhaoran.wang@northwestern.edu, mihailo@usc.edu

BACKGROUND AND MOTIVATION

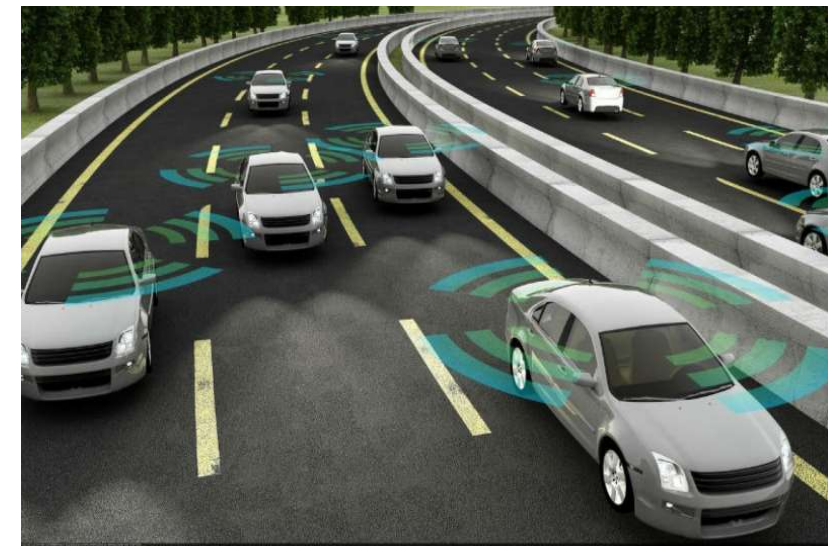
Multi-agent reinforcement learning



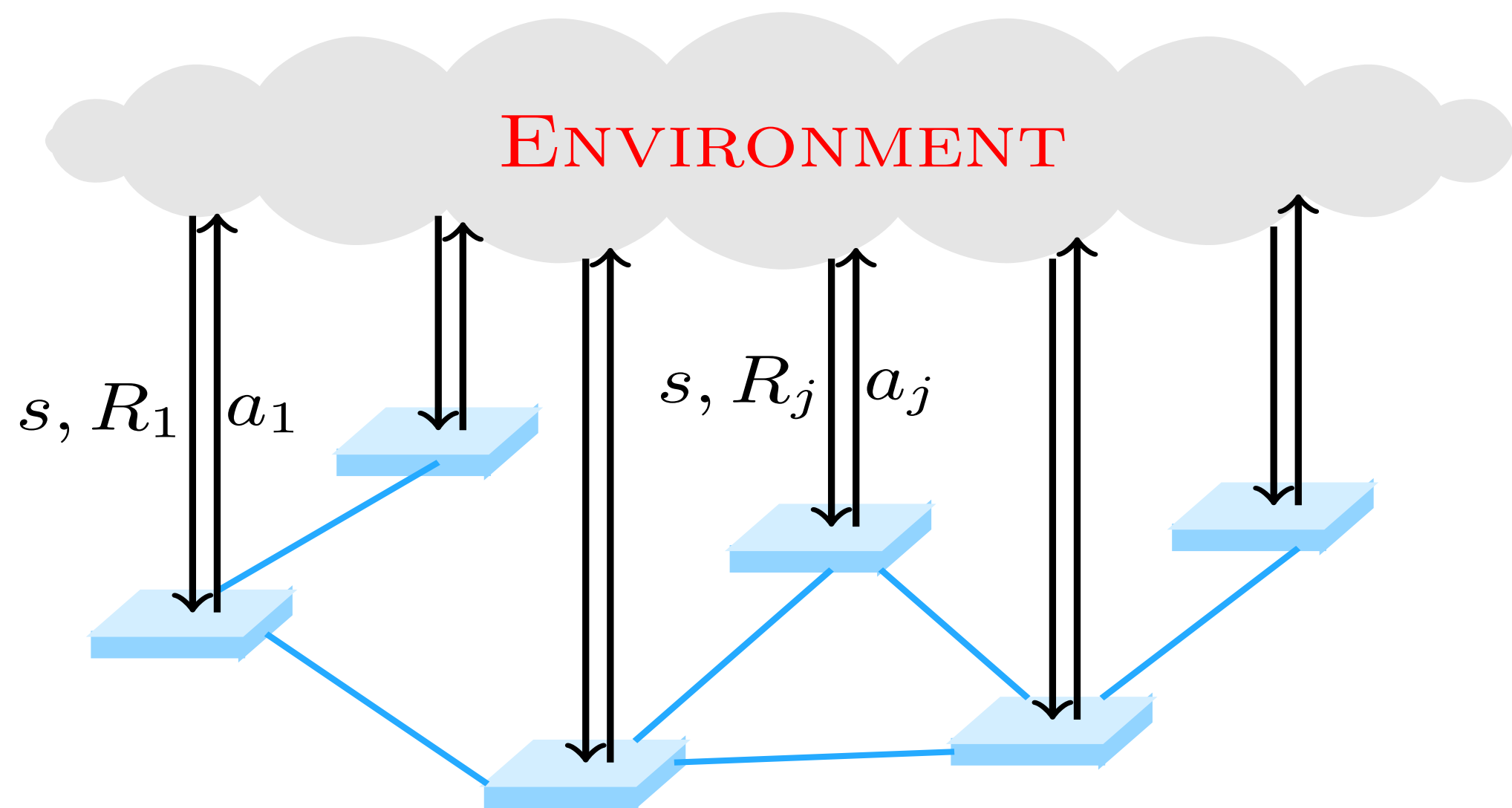
Power



Robotics



Transportation



$$R_c^\pi(s) = \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{a \sim \pi(\cdot|s)} [R_j(s, a)]$$

- Objective of policy evaluation

$$V^\pi(s) = \mathbb{E}[R_c^\pi(s_0) + \gamma R_c^\pi(s_1) + \gamma^2 R_c^\pi(s_2) + \dots | s_0 = s, \pi]$$

- Challenges

- Distributed samples over a network
- Markovian samples for a given policy

MULTI-AGENT TD LEARNING

Bellman error minimization

- Bellman equation with linear function approximation

$$\mathbf{V}_x = \gamma \mathbf{P}^\pi \mathbf{V}_x + \mathbf{R}_c^\pi$$

- $V_x(s) = \phi^T(s)x$ - linear approximation of $V^\pi(s)$
- $\mathbf{V}_x, \mathbf{R}_c^\pi$ - vectors of $V_x(s), R_c^\pi(s)$ for all states s
- \mathbf{P}^π - probability transition matrix

- Projected Bellman error minimization

Centralized problem	Decentralized problem
minimize $\frac{1}{2} \ Ax - b\ _{C^{-1}}^2$	minimize $\frac{1}{2N} \sum_{j=1}^N \ Ax - b_j\ _{C^{-1}}^2$
$b = \mathbb{E}_{s \sim \Pi} [\mathcal{R}_c^\pi(s) \phi(s)]$	$b_j = \mathbb{E}_{s \sim \Pi} [\mathcal{R}_j^\pi(s) \phi(s)]$

- $A = \mathbb{E}_{s \sim \Pi} [\phi(s)(\phi(s) - \gamma \phi(s'))^T]$ and $C = \mathbb{E}_{s \sim \Pi} [\phi(s)\phi(s)^T]$
- Π - unknown stationary distribution for a given policy

Decentralized stochastic saddle-point problem

- Dualization of the objective function

$$\|Ax - b_j\|_{C^{-1}}^2 = \max_{y_j \in \mathcal{Y}} \underbrace{y_j^T (Ax - b_j) - \frac{1}{2} y_j^T C y_j}_{\psi_j(x, y_j) = \mathbb{E}_{\xi \sim \Pi} [\Psi_j(x, y_j; \xi)]}$$

- Stochastic saddle-point problem

$$\min_{x \in \mathcal{X}} \max_{y_j \in \mathcal{Y}} \psi(x, y) := \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{\xi \sim \Pi} [\Psi_j(x, y_j; \xi)]$$

- Dependent samples, unknown distribution

- $\xi_t \sim P_t$ - samples from a Markov process P_t at time t
- P_t - Markov process converging to the unknown Π

HOMOTOPY PRIMAL-DUAL ALGORITHM

- Distributed dual averaging for aggregating local information
- Homotopy method for adaptive stepsize selection

Algorithm 1 Distributed Homotopy Primal-Dual (DHPD) Algorithm

Initialization: $x_{j,1}(1) = x'_{j,1}(1) = 0, y_{j,1}(1) = y'_{j,1}(1) = 0, \eta_1, T_1, K$

For $k = 1$ to K **do** ▷ for all agents $j \in \mathcal{V}$

(1) **For** $t = 1$ to $T_k - 1$ **do**

- Primal update

▷ Distributed dual averaging

$$x'_{j,k}(t+1) = \sum_{i=1}^N W_{ij} x'_{i,k}(t) - \eta_k \nabla_x \Psi_j(z_{j,k}(t); \xi_k(t))$$

$$x_{j,k}(t+1) = \mathcal{P}_X(x'_{j,k}(t+1))$$

- Dual update

▷ Local update

$$y'_{j,k}(t+1) = y'_{j,k}(t) + \eta_k \nabla_y \Psi_j(z_{j,k}(t); \xi_k(t))$$

$$y_{j,k}(t+1) = \mathcal{P}_Y(y'_{j,k}(t+1))$$

end for

(2) $(x_{j,k+1}(1), y_{j,k+1}(1)) = \left(\frac{1}{T_k} \sum_{t=1}^{T_k} x_{j,k}(t), \frac{1}{T_k} \sum_{t=1}^{T_k} y_{j,k}(t) \right)$

(3) $(x'_{j,k+1}(1), y'_{j,k+1}(1)) = (x_{j,k+1}(1), y_{j,k+1}(1))$

(4) $\eta_{k+1} = \eta_k/2, T_{k+1} = 2T_k$

▷ Adaptive stepsize

end for

Output: $(\hat{x}_{j,K}, \hat{y}_{j,K}) = \left(\frac{1}{T_K} \sum_{t=1}^{T_K} x_{j,K}(t), \frac{1}{T_K} \sum_{t=1}^{T_K} y_{j,K}(t) \right)$

- Optimality gap induced by $\hat{x}_{i,k}$

$$\epsilon(\hat{x}_{i,k}) = \frac{1}{2N} \sum_{j=1}^N \left(\|A\hat{x}_{i,k} - b_j\|_{C^{-1}}^2 - \|Ax^* - b_j\|_{C^{-1}}^2 \right)$$

FAST CONVERGENCE RATE

- Assumptions

- $\sup_{x \in \mathcal{X}, y \in \mathcal{Y}} \|(x, y)\|^2 \leq R^2$ - convex compact domain
- $\psi_j(x, y_j) - \rho_y$ -strongly concave; G -gradient bounded; L -gradient Lipschitz
- $\max_{y_j \in \mathcal{Y}} \psi_j(x, y_j) - \rho_x$ -strongly convex
- W - doubly stochastic communication matrix

- Claim:** For any $\eta_1 \geq 1/(4/\rho_y + 2/\rho_x)$, any T_1, K satisfying $T_1 \geq 1 + \lceil \log(\Gamma T) / |\log \rho| \rceil := \tau$ where $T = \sum_{k=1}^K T_k$, we have

$$\frac{1}{N} \sum_{j=1}^N \mathbb{E}[\epsilon(\hat{x}_{j,K})] = C_1 \frac{G(RL+G)\log^2(\sqrt{NT})}{T(1-\sigma_2(W))} + C_2 \frac{G(G+RL)(1+T_1)}{T}$$

- Fast convergence rate $O(\log(\sqrt{NT})/T)$
- Network dependence $\log^2(\sqrt{NT})/(1-\sigma_2(W))$
- Fast $1/T$ -mixing fast convergence

CASE STUDY OF MOUNTAIN CAR TASK

- SPD - stochastic primal-dual method

