

PROJECT REPORT

MULTI-AGENT REINFORCEMENT LEARNING VIA ADAPTIVE KALMAN TEMPORAL DIFFERENCE AND SUCCESSOR REPRESENTATION

Mohammad Salimibeni | Arash Mohammadi | Parvin Malekzadeh and Konstantinos N. Plataniotis

RLinControl | Fall 2024

Professor: Dr. Saeed Shamaghdari
Seyede Setare Khosravi | Mobina Lashgari

TABLE OF CONTENTS

Abstract

Introduction

Problem Formulation

MAK-TD Framework

MAK-SR Framework

Experimental Results

Simulation Results

Conclusions

ABSTRACT

Reinforcement Learning Challenges

- Fixed Reward Function
- Overfitting, High sensitivity and ...

Kalman an Idea to Solve The Challenges

- Online Learning and ...

INTRODUCTION_{MARI}

Questions!

INTRODUCTION_{MARI}

Challenges!

Problem Formulation

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

Problem Formulation

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

SINGLE AGENT REINFORCEMENT LEARNING

$$Q_{\pi}(\mathbf{s}, a) = \mathbb{E} \left\{ \sum_{k=0}^T \gamma^k r_k \mid \mathbf{s}_0 = \mathbf{s}, a_0 = a, a_k = \pi(\mathbf{s}_k) \right\}$$

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

SINGLE AGENT REINFORCEMENT LEARNING

$$Q_{\pi}(\mathbf{s}, a) = \mathbb{E} \left\{ \sum_{k=0}^T \gamma^k r_k \mid \mathbf{s}_0 = \mathbf{s}, a_0 = a, a_k = \pi(\mathbf{s}_k) \right\}$$

$$a_k = \arg \max_{a \in \mathcal{A}} Q_{\pi^*}(\mathbf{s}_k, a).$$

Problem Formulation

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

TD LEARNING

$$Q_{\pi^*}(\mathbf{s}_k, a_k) = Q_{\pi^*}(\mathbf{s}_k, a_k) + \alpha \left(r_k + \gamma \max_{a \in \mathcal{A}} Q_{\pi^*}(\mathbf{s}_{k+1}, a) - Q_{\pi^*}(\mathbf{s}_k, a_k) \right)$$

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

TD LEARNING

$$Q_{\pi^*}(\mathbf{s}_k, a_k) = Q_{\pi^*}(\mathbf{s}_k, a_k) + \alpha \left(r_k + \gamma \max_{a \in \mathcal{A}} Q_{\pi^*}(\mathbf{s}_{k+1}, a) - Q_{\pi^*}(\mathbf{s}_k, a_k) \right)$$

$$a_k = \arg \max_{a \in \mathcal{A}} Q_{\pi^*}(\mathbf{s}_k, a)$$

MULTI AGENT SETTING

Agent i , for $(1 \leq i \leq N)$

$$\mathbb{S} = \{\mathcal{S}^{(1)}, \dots, \mathcal{S}^{(N)}\}$$

$$\mathbb{A} = \{\mathcal{A}^{(1)}, \dots, \mathcal{A}^{(N)}\}$$

$$\mathbb{Z} = \{\mathcal{Z}^{(1)}, \dots, \mathcal{Z}^{(N)}\}$$

$$r^{(i)} : \mathbb{S} \times \mathcal{A}^{(i)} \rightarrow \mathbb{R}$$

$$R^{(i)} = \sum_{t=0}^T \gamma^t (r^{(i)})^t$$

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

MULTI AGENT SR

$$M_{\pi^{(i)}}(s^{(i)}, s'^{(i)}, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \mathbb{1}[s_k^{(i)} = s'^{(i)}] \mid s_0^{(i)} = s^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$M_{\pi^{(i)}}^{\text{new}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) = M_{\pi^{(i)}}^{\text{old}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) + \\ \alpha \left(\mathbb{1}[s_k^{(i)} = s'^{(i)}] + \gamma M_{\pi^{(i)}}(s_{k+1}^{(i)}, s'^{(i)}, a_{k+1}^{(i)}) - M_{\pi^{(i)}}^{\text{old}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) \right)$$

$$Q_{\pi^{(i)}}(s_k^{(i)}, a_k^{(i)}) = \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} M(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) R^{(i)}(s'^{(i)}, a_k^{(i)})$$

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

MULTI AGENT SR

$$M_{\pi^{(i)}}(s^{(i)}, s'^{(i)}, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \mathbb{1}[s_k^{(i)} = s'^{(i)}] | s_0^{(i)} = s^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$M_{\pi^{(i)}}^{\text{new}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) = M_{\pi^{(i)}}^{\text{old}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) + \\ \alpha \left(\mathbb{1}[s_k^{(i)} = s'^{(i)}] + \gamma M_{\pi^{(i)}}(s_{k+1}^{(i)}, s'^{(i)}, a_{k+1}^{(i)}) - M_{\pi^{(i)}}^{\text{old}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) \right)$$

$$Q_{\pi^{(i)}}(s_k^{(i)}, a_k^{(i)}) = \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} M(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) R^{(i)}(s'^{(i)}, a_k^{(i)})$$

- Single Agent Reinforcement Learning
- Off-Policy Temporal Difference (TD) Learning
- Multi-Agent Setting
- Multi-Agent SR

MULTI AGENT SR

$$M_{\pi^{(i)}}(s^{(i)}, s'^{(i)}, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \mathbb{1}[s_k^{(i)} = s'^{(i)}] | s_0^{(i)} = s^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$M_{\pi^{(i)}}^{\text{new}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) = M_{\pi^{(i)}}^{\text{old}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) + \\ \alpha \left(\mathbb{1}[s_k^{(i)} = s'^{(i)}] + \gamma M_{\pi^{(i)}}(s_{k+1}^{(i)}, s'^{(i)}, a_{k+1}^{(i)}) - M_{\pi^{(i)}}^{\text{old}}(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) \right)$$

$$Q_{\pi^{(i)}}(s_k^{(i)}, a_k^{(i)}) = \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} M(s_k^{(i)}, s'^{(i)}, a_k^{(i)}) R^{(i)}(s'^{(i)}, a_k^{(i)})$$

MAK-TD FRAMEWORK

$$Q_{\pi^{(i)*}}(\mathbf{s}_k^{(i)}, a_k^{(i)}) \approx r_k^{(i)} + \gamma \max_{a^{(i)} \in \mathcal{A}} Q_{\pi^{(i)*}}(\mathbf{s}_{k+1}^{(i)}, a^{(i)})$$

$$r_k^{(i)} = Q_{\pi^{(i)*}}(\mathbf{s}_k^{(i)}, a_k^{(i)}) - \gamma \max_{a^{(i)} \in \mathcal{A}} Q_{\pi^{(i)*}}(\mathbf{s}_{k+1}^{(i)}, a^{(i)}) + v_k^{(i)}$$

v_k is modeled as a zero-mean normal distribution with variance of $R^{(i)}$

MAK-TD FRAMEWORK

$$Q_{\pi^{(i)*}}(\mathbf{s}_k^{(i)}, a_k^{(i)}) \approx r_k^{(i)} + \gamma \max_{a^{(i)} \in \mathcal{A}} Q_{\pi^{(i)*}}(\mathbf{s}_{k+1}^{(i)}, a^{(i)})$$

$$r_k^{(i)} = Q_{\pi^{(i)*}}(\mathbf{s}_k^{(i)}, a_k^{(i)}) - \gamma \max_{a^{(i)} \in \mathcal{A}} Q_{\pi^{(i)*}}(\mathbf{s}_{k+1}^{(i)}, a^{(i)}) + v_k^{(i)}$$

v_k is modeled as a zero-mean normal distribution with variance of $R^{(i)}$

MAK-TD FRAMEWORK

$$Q_{\pi^{(i)}}(\mathbf{s}_k^{(i)}, a_k^{(i)}) = \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$r_k^{(i)} = \left[\boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(\mathbf{s}_{k+1}^{(i)}, a^{(i)})^T \right] \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

$$\mathbf{h}_k^{(i)} = \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(\mathbf{s}_{k+1}^{(i)}, a^{(i)})$$

$$r_k^{(i)} = [\mathbf{h}_k^{(i)}]^T \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

MAK-TD FRAMEWORK

$$Q_{\pi^{(i)}}(\mathbf{s}_k^{(i)}, a_k^{(i)}) = \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$r_k^{(i)} = \left[\boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(\mathbf{s}_{k+1}^{(i)}, a^{(i)})^T \right] \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

$$\mathbf{h}_k^{(i)} = \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(\mathbf{s}_{k+1}^{(i)}, a^{(i)})$$

$$r_k^{(i)} = [\mathbf{h}_k^{(i)}]^T \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

MAK-TD FRAMEWORK

$$Q_{\pi^{(i)}}(s_k^{(i)}, a_k^{(i)}) = \boldsymbol{\phi}(s_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$r_k^{(i)} = \left[\boldsymbol{\phi}(s_k^{(i)}, a_k^{(i)})^T - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(s_{k+1}^{(i)}, a^{(i)})^T \right] \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

$$\mathbf{h}_k^{(i)} = \boldsymbol{\phi}(s_k^{(i)}, a_k^{(i)}) - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(s_{k+1}^{(i)}, a^{(i)})$$

$$r_k^{(i)} = [\mathbf{h}_k^{(i)}]^T \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

MAK-TD FRAMEWORK

$$Q_{\pi^{(i)}}(s_k^{(i)}, a_k^{(i)}) = \boldsymbol{\phi}(s_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$r_k^{(i)} = \left[\boldsymbol{\phi}(s_k^{(i)}, a_k^{(i)})^T - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(s_{k+1}^{(i)}, a^{(i)})^T \right] \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

$$\mathbf{h}_k^{(i)} = \boldsymbol{\phi}(s_k^{(i)}, a_k^{(i)}) - \gamma \max_{a^{(i)} \in \mathcal{A}} \boldsymbol{\phi}(s_{k+1}^{(i)}, a^{(i)})$$

$$r_k^{(i)} = [\mathbf{h}_k^{(i)}]^T \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

MAK-TD FRAMEWORK

$$r_k^{(i)} = [\mathbf{h}_k^{(i)}]^T \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

$$\mathbf{x}_k = \mathbf{F}_{k-1} \mathbf{x}_{k-1} + \mathbf{G}_{k-1} \mathbf{u}_{k-1} + \mathbf{w}_{k-1}$$

$$y_k = \mathbf{H}_k \mathbf{x}_k + v_k$$

$$E(\mathbf{w}_k) = E(v_k) = 0 \quad E(\mathbf{w}_k \mathbf{w}_j^T) = \mathbf{Q}_k \delta_{k-j} \quad E(v_k v_j^T) = R_k \delta_{k-j} \quad E(v_k \mathbf{w}_j^T) = 0$$

MAK-TD FRAMEWORK

$$x_k = F_{k-1}x_{k-1} + G_{k-1}u_{k-1} + w_{k-1}$$

$$y_k = H_k x_k + v_k$$

$$r_k^{(i)} = [h_k^{(i)}]^T \theta_k^{(i)} + v_k^{(i)}$$

$$E(w_k) = E(v_k) = 0 \quad E(w_k w_j^T) = Q_k \delta_{k-j} \quad E(v_k v_j^T) = R_k \delta_{k-j} \quad E(v_k w_j^T) = 0$$

MAK-TD FRAMEWORK

$$x_k = F_{k-1}x_{k-1} + G_{k-1}u_{k-1} + w_{k-1}$$

$$y_k = H_k x_k + v_k$$

$$r_k^{(i)} = [h_k^{(i)}]^T \theta_k^{(i)} + v_k^{(i)}$$

$$E(w_k) = E(v_k) = 0 \quad E(w_k w_j^T) = Q_k \delta_{k-j} \quad E(v_k v_j^T) = R_k \delta_{k-j} \quad E(v_k w_j^T) = 0$$

$$K_k = P_k^- H_k^T (H_k P_k^- H_k^T + R_k)^{-1} = P_k^+ H_k^T R_k^{-1}$$

$$\hat{x}_k^+ = \hat{x}_k^- + K_k (y_k - H_k \hat{x}_k^-)$$

$$P_k^+ = (I - K_k H_k) P_k^- = (I - K_k H_k) P_k^- (I - K_k H_k)^T + K_k R_k K_k^T$$

MAK-TD FRAMEWORK

$$r_k^{(i)} = [\mathbf{h}_k^{(i)}]^T \boldsymbol{\theta}_k^{(i)} + v_k^{(i)}$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k)^{-1} = \mathbf{P}_k^+ \mathbf{H}_k^T \mathbf{R}_k^{-1}$$

$$\hat{\mathbf{x}}_k^+ = \hat{\mathbf{x}}_k^- + \mathbf{K}_k (y_k - \mathbf{H}_k \hat{\mathbf{x}}_k^-)$$

$$\mathbf{P}_k^+ = (\mathbf{I} - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_k^- = (\mathbf{I} - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_k^- (\mathbf{I} - \mathbf{K}_k \mathbf{H}_k)^T + \mathbf{K}_k \mathbf{R}_k \mathbf{K}_k^T$$

$$\mathbf{K}_k^{j(i)} = \mathbf{P}_{(\boldsymbol{\theta}, k|k-1)}^{(i)} \mathbf{h}_k^{(i)} (\mathbf{h}_k^{T(i)} \mathbf{P}_{(\boldsymbol{\theta}, k|k-1)}^{(i)} \mathbf{h}_k^{(i)} + \mathbf{R}^{j(i)})^{-1}$$

$$\hat{\boldsymbol{\theta}}_k^{j(i)} = \hat{\boldsymbol{\theta}}_{(k|k-1)}^{(i)} + \mathbf{K}_k^{j(i)} (r_k^{(i)} - \mathbf{h}_k^{T(i)} \hat{\boldsymbol{\theta}}_{(k|k-1)}^{(i)})$$

$$\mathbf{P}_{\boldsymbol{\theta}, k}^{j(i)} = (\mathbf{I} - \mathbf{K}_k^{j(i)} \mathbf{h}_k^{T(i)}) \mathbf{P}_{(\boldsymbol{\theta}, k|k-1)}^{T(i)} (\mathbf{I} - \mathbf{K}_k^{j(i)} \mathbf{h}_k^{T(i)}) + \mathbf{K}_k^{j(i)} \mathbf{R}^{j(i)} \mathbf{K}_k^{j(i)T}$$

MAK-TD FRAMEWORK

$$x_k = F_{k-1}x_{k-1} + G_{k-1}u_{k-1} + w_{k-1}$$

$$\bar{x}_k = F_{k-1}\bar{x}_{k-1} + G_{k-1}u_{k-1}$$

$$P_k = F_{k-1}P_{k-1}F_{k-1}^T + Q_{k-1}$$

MAK-TD FRAMEWORK

$$\hat{\boldsymbol{\theta}}_k^{(i)} = \sum_{j=1}^M \omega^{j(i)} \hat{\boldsymbol{\theta}}_k^{j(i)}$$

$$\mathbf{P}_{\boldsymbol{\theta},k}^{(i)} = \sum_{j=1}^M \omega^{j(i)} \left(\mathbf{P}_{\boldsymbol{\theta},k}^{j(i)} + (\hat{\boldsymbol{\theta}}^{j(i)} - \hat{\boldsymbol{\theta}}^{(i)})(\hat{\boldsymbol{\theta}}^{j(i)} - \hat{\boldsymbol{\theta}}^{(i)})^T \right)$$

MAK-TD FRAMEWORK

$$\boldsymbol{\phi}(\mathbf{s}_k^{(i)}) = [\phi_1(\mathbf{s}_k^{(i)}), \phi_2(\mathbf{s}_k^{(i)}), \dots, \phi_{N_b-1}(\mathbf{s}_k^{(i)}), \phi_{N_b}(\mathbf{s}_k^{(i)})]^T$$

$$\phi_n(\mathbf{s}_k^{(i)}) = \exp\left\{\frac{-1}{2}(\mathbf{s}_k^{(i)} - \boldsymbol{\mu}_n^{(i)})^T \boldsymbol{\Sigma}_n^{(i)-1} (\mathbf{s}_k^{(i)} - \boldsymbol{\mu}_n^{(i)})\right\}$$

where $\boldsymbol{\mu}_n^{(i)}$ and $\boldsymbol{\Sigma}_n^{(i)}$ are the mean and covariance of $\phi_n(\mathbf{s}_k^{(i)})$, for $(1 \leq n \leq N_b)$

MAK-TD FRAMEWORK

$$\boldsymbol{\phi}(\mathbf{s}_k^{(i)}) = [\phi_1(\mathbf{s}_k^{(i)}), \phi_2(\mathbf{s}_k^{(i)}), \dots, \phi_{N_b-1}(\mathbf{s}_k^{(i)}), \phi_{N_b}(\mathbf{s}_k^{(i)})]^T$$

$$\phi_n(\mathbf{s}_k^{(i)}) = \exp\left\{-\frac{1}{2}(\mathbf{s}_k^{(i)} - \boldsymbol{\mu}_n^{(i)})^T \boldsymbol{\Sigma}_n^{(i)-1} (\mathbf{s}_k^{(i)} - \boldsymbol{\mu}_n^{(i)})\right\}$$

where $\boldsymbol{\mu}_n^{(i)}$ and $\boldsymbol{\Sigma}_n^{(i)}$ are the mean and covariance of $\phi_n(\mathbf{s}_k^{(i)})$, for $(1 \leq n \leq N_b)$

MAK-TD FRAMEWORK

$$\boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) = [\phi_{1,a_1}(\mathbf{s}_k^{(i)}), \dots, \phi_{N_b,a_1}(\mathbf{s}_k^{(i)}), \phi_{1,a_2}(\mathbf{s}_k^{(i)}), \dots, \phi_{N_b,a_{D(i)}}(\mathbf{s}_k^{(i)})]^T$$

$$\boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) = [0, \dots, 0, \phi_1(\mathbf{s}_k^{(i)}), \dots, \phi_N(\mathbf{s}_k^{(i)}), 0, \dots, 0]^T$$

MAK-TD FRAMEWORK

$$L_k^{(i)} = (\boldsymbol{\phi}^T(\mathbf{s}_k^{(i)}, a_k) \boldsymbol{\theta}_k^{(i)} - r_k^{(i)})^2$$

$$\Delta \boldsymbol{\mu}^{(i)} = -\frac{\partial L_k^{(i)}}{\partial \boldsymbol{\mu}^{(i)}} = -\frac{\partial L_k^{(i)}}{\partial Q_{\pi^*(i)}} \frac{\partial Q_{\pi^*(i)}}{\partial \boldsymbol{\phi}^{(i)}} \frac{\partial \boldsymbol{\phi}^{(i)}}{\partial \boldsymbol{\mu}^{(i)}}$$

and $\Delta \boldsymbol{\Sigma}^{(i)} = -\frac{\partial \boldsymbol{\Sigma}_k^{(i)}}{\partial \boldsymbol{\mu}^{(i)}} = -\frac{\partial L_k^{(i)}}{\partial Q_{\pi^*(i)}} \frac{\partial Q_{\pi^*(i)}}{\partial \boldsymbol{\phi}^{(i)}} \frac{\partial \boldsymbol{\phi}^{(i)}}{\partial \boldsymbol{\Sigma}^{(i)}}$



MAK-TD FRAMEWORK

$$L_k^{(i)} = (\boldsymbol{\phi}^T(\mathbf{s}_k^{(i)}, a_k) \boldsymbol{\theta}_k^{(i)} - r_k^{(i)})^2$$

$$\Delta \boldsymbol{\mu}^{(i)} = -\frac{\partial L_k^{(i)}}{\partial \boldsymbol{\mu}^{(i)}} = -\frac{\partial L_k^{(i)}}{\partial Q_{\pi^*(i)}} \frac{\partial Q_{\pi^*(i)}}{\partial \boldsymbol{\phi}^{(i)}} \frac{\partial \boldsymbol{\phi}^{(i)}}{\partial \boldsymbol{\mu}^{(i)}}$$

and $\Delta \boldsymbol{\Sigma}^{(i)} = -\frac{\partial \boldsymbol{\Sigma}_k^{(i)}}{\partial \boldsymbol{\mu}^{(i)}} = -\frac{\partial L_k^{(i)}}{\partial Q_{\pi^*(i)}} \frac{\partial Q_{\pi^*(i)}}{\partial \boldsymbol{\phi}^{(i)}} \frac{\partial \boldsymbol{\phi}^{(i)}}{\partial \boldsymbol{\Sigma}^{(i)}}$



MAK-TD FRAMEWORK

$$\begin{aligned} a_k^{(i)} &= \arg \max_a \left(\mathbf{h}_k^{(i)}(\mathbf{s}_k^{(i)}, a^{(i)}) R^{-1(i)} \mathbf{h}_k^{T(i)}(\mathbf{s}_k^{(i)}, a^{(i)}) \right) \\ &= \arg \max_a \left(\mathbf{h}_k^{(i)}(\mathbf{s}_k^{(i)}, a^{(i)}) \mathbf{h}_k^{T(i)}(\mathbf{s}_k^{(i)}, a^{(i)}) \right). \end{aligned}$$

Algorithm 1 THE PROPOSED MAK-TD FRAMEWORK

```

1: Learning Phase:
2: Set  $\theta_0, P_{\theta,0}, F, \mu_{n,i_d}, \Sigma_{n,i_d}$  for  $n = 1, 2, \dots, N$  and  $i_d = 1, 2, \dots, D$ 
3: Repeat (for each episode):
4:   Initialize  $s_k$ 
5:   Repeat (for each agent  $i$ ):
6:     While  $s_k^{(i)} \neq s_T$  do:
7:        $a_k^{(i)} = \arg \max_a \left( h_k^{(i)}(s_k^{(i)}, a^{(i)}) h_k^{T(i)}(s_k^{(i)}, a^{(i)}) \right)$ 
8:       Take action  $a_k^{(i)}$ , observe  $s_{k+1}^{(i)}, r_k^{(i)}$ 
9:       Calculate  $\phi^{(i)}(s^{(i)}, a^{(i)})$  via Equations (22) and (23)
10:       $h_k^{(i)}(s_k^{(i)}, a_k^{(i)}) = \phi^{(i)}(s_k^{(i)}, a_k^{(i)}) - \gamma \arg \max_a \phi^{(i)}(s_{k+1}^{(i)}, a^{(i)})$ 
11:       $\hat{\theta}_{(k|k-1)}^{(i)} = F^{(i)} \hat{\theta}_k^{(i)}$ 
12:       $P_{(\theta,k|k-1)}^{(i)} = F^{(i)} P_{\theta,k-1}^{(i)} F^{T(i)} + Q^{(i)}$ 
13:      for  $j = 1 : M$  do:
14:         $k_k^{j(i)} = P_{(\theta,k|k-1)}^{(i)} h_k^{(i)} (h_k^{T(i)} P_{(\theta,k|k-1)}^{(i)} h_k^{(i)} + R^{j(i)})^{-1}$ 
15:         $\hat{\theta}_k^{j(i)} = \hat{\theta}_{(\theta,k|k-1)}^{(i)} + k_k^{j(i)} (r_k^j - h_k^{T(i)} \hat{\theta}_{(k|k-1)}^{(i)})$ 
16:         $P_{\theta,k}^{(i)} = (I - K_k^{j(i)} h_k^{T(i)}) P_{(\theta,k|k-1)}^{(i)} (I - K_k^{j(i)} h_k^{T(i)})^T + K_k^{j(i)} R^j K_k^{jT(i)}$ 
17:      end for

```

```

18:      Compute the value of  $c$  and  $w^{j^{(l)}}$  by using  $\sum_{j=1}^M w^{j^{(l)}} = 1$  and Equation (19)
19:       $\hat{\theta}_k^{(i)} = \sum_{j=1}^M w^{j^{(i)}} \hat{\theta}_k^{j^{(i)}}$ 
20:       $P_{\theta_k}^{(i)} = \sum_{j=1}^M \omega^{j^{(i)}} \left( P_{\theta_k}^{j^{(i)}} + (\hat{\theta}^{j^{(i)}} - \hat{\theta}^{(i)})(\hat{\theta}^{j^{(i)}} - \hat{\theta}^{(i)})^T \right)$ 
21:      RBFs Parameters Update:
22:       $L_k^{(i)} = (\phi^T(s_k^{(i)}, a_k) \theta_k^{(i)} - r_k^{(i)})^2$ 
23:      if  $L_k^{(i)\frac{1}{2}} (\theta_k^{(i)T} \phi(\cdot)) > 0$  then:
24:          Update  $\Sigma_{n,a_d}$  via Equation (29)
25:      else:
26:          Update  $\mu_{n,a_d}$  via Equation (30)
27:      end if
28:  end while
29: Testing Phase:
30: Repeat (for each trial episode):
31:   While  $s_k \neq s_T$  do:
32:    Repeat (for each agent):
33:      $a_k = \arg \max_a \phi(s_k, a)^T \theta_k$ 
34:     Take action  $a_k$ , and observe  $s_{k+1}, r_k$ 
35:     Calculate Loss  $S_k$  for all agents
36:   End While

```

MAK-SR FRAMEWORK

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}^{(i)}, :, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) | \mathbf{s}_0^{(i)} = \mathbf{s}^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$r^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) \approx \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$Q(\mathbf{s}_k^{(i)}, a_k^{(i)}) = \boldsymbol{\theta}_k^{(i)T} \mathbf{M}(\mathbf{s}_k^{(i)}, :, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) \approx \mathbf{M}_k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}^{\text{new}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) = \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) + \alpha(\boldsymbol{\phi}^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) + \gamma \mathbf{M}_{\pi^{(i)}}(\mathbf{s}_{k+1}^{(i)}, :, a_{k+1}^{(i)}) - \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}))$$

MAK-SR FRAMEWORK

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}^{(i)}, :, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) | \mathbf{s}_0^{(i)} = \mathbf{s}^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$r^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) \approx \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$Q(\mathbf{s}_k^{(i)}, a_k^{(i)}) = \boldsymbol{\theta}_k^{(i)T} \mathbf{M}(\mathbf{s}_k^{(i)}, :, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) \approx \mathbf{M}_k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}^{\text{new}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) = \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) + \alpha(\boldsymbol{\phi}^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) + \gamma \mathbf{M}_{\pi^{(i)}}(\mathbf{s}_{k+1}^{(i)}, :, a_{k+1}^{(i)}) - \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}))$$

MAK-SR FRAMEWORK

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}^{(i)}, :, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) | \mathbf{s}_0^{(i)} = \mathbf{s}^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$r^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) \approx \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$Q(\mathbf{s}_k^{(i)}, a_k^{(i)}) = \boldsymbol{\theta}_k^{(i)T} \mathbf{M}(\mathbf{s}_k^{(i)}, :, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) \approx \mathbf{M}_k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}^{\text{new}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) = \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) + \alpha(\boldsymbol{\phi}^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) + \gamma \mathbf{M}_{\pi^{(i)}}(\mathbf{s}_{k+1}^{(i)}, :, a_{k+1}^{(i)}) - \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}))$$

MAK-SR FRAMEWORK

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}^{(i)}, :, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) | \mathbf{s}_0^{(i)} = \mathbf{s}^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$r^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) \approx \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$Q(\mathbf{s}_k^{(i)}, a_k^{(i)}) = \boldsymbol{\theta}_k^{(i)T} \mathbf{M}(\mathbf{s}_k^{(i)}, :, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) \approx \mathbf{M}_k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}^{\text{new}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) = \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) + \alpha(\boldsymbol{\phi}^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) + \gamma \mathbf{M}_{\pi^{(i)}}(\mathbf{s}_{k+1}^{(i)}, :, a_{k+1}^{(i)}) - \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}))$$

MAK-SR FRAMEWORK

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}^{(i)}, :, a^{(i)}) = \mathbb{E} \left[\sum_{k=0}^T \gamma^k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)}) | \mathbf{s}_0^{(i)} = \mathbf{s}^{(i)}, a_0^{(i)} = a^{(i)} \right]$$

$$r^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) \approx \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})^T \boldsymbol{\theta}_k^{(i)}$$

$$Q(\mathbf{s}_k^{(i)}, a_k^{(i)}) = \boldsymbol{\theta}_k^{(i)T} \mathbf{M}(\mathbf{s}_k^{(i)}, :, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) \approx \mathbf{M}_k \boldsymbol{\phi}(\mathbf{s}_k^{(i)}, a_k^{(i)})$$

$$\mathbf{M}_{\pi^{(i)}}^{\text{new}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) = \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}) + \alpha(\boldsymbol{\phi}^{(i)}(\mathbf{s}_k^{(i)}, a_k^{(i)}) + \gamma \mathbf{M}_{\pi^{(i)}}(\mathbf{s}_{k+1}^{(i)}, :, a_{k+1}^{(i)}) - \mathbf{M}_{\pi^{(i)}}^{\text{old}}(\mathbf{s}_k^{(i)}, :, a_k^{(i)}))$$

MAK-SR FRAMEWORK

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M^{\text{new}}(s_k^{(i)}, :, a_k^{(i)}) - \gamma M(s_{k+1}^{(i)}, :, a_{k+1}^{(i)}) + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M_k \left[\underbrace{\phi(s_k^{(i)}, a_k^{(i)}) - \gamma \phi(s_{k+1}^{(i)}, a_{k+1}^{(i)})}_{g_k^{(i)}} \right] + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = (g_k^{(i)T} \otimes I) m_k^{(i)} + n_k^{(i)}$$

$$m_{k+1}^{(i)} = m_k^{(i)} + \mu_k^{(i)}$$

MAK-SR FRAMEWORK

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M^{\text{new}}(s_k^{(i)}, :, a_k^{(i)}) - \gamma M(s_{k+1}^{(i)}, :, a_{k+1}^{(i)}) + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M_k \left[\underbrace{\phi(s_k^{(i)}, a_k^{(i)}) - \gamma \phi(s_{k+1}^{(i)}, a_{k+1}^{(i)})}_{g_k^{(i)}} \right] + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = (g_k^{(i)T} \otimes I) m_k^{(i)} + n_k^{(i)}$$

$$m_{k+1}^{(i)} = m_k^{(i)} + \mu_k^{(i)}$$

MAK-SR FRAMEWORK

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M^{\text{new}}(s_k^{(i)}, :, a_k^{(i)}) - \gamma M(s_{k+1}^{(i)}, :, a_{k+1}^{(i)}) + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M_k \underbrace{\left[\phi(s_k^{(i)}, a_k^{(i)}) - \gamma \phi(s_{k+1}^{(i)}, a_{k+1}^{(i)}) \right]}_{g_k^{(i)}} + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = (g_k^{(i)T} \otimes I) m_k^{(i)} + n_k^{(i)}$$

$$m_{k+1}^{(i)} = m_k^{(i)} + \mu_k^{(i)}$$

MAK-SR FRAMEWORK

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M^{\text{new}}(s_k^{(i)}, :, a_k^{(i)}) - \gamma M(s_{k+1}^{(i)}, :, a_{k+1}^{(i)}) + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = M_k \left[\underbrace{\phi(s_k^{(i)}, a_k^{(i)}) - \gamma \phi(s_{k+1}^{(i)}, a_{k+1}^{(i)})}_{g_k^{(i)}} \right] + n_k^{(i)}$$

$$\hat{\phi}(s_k^{(i)}, a_k^{(i)}) = (g_k^{(i)T} \otimes I) m_k^{(i)} + n_k^{(i)}$$

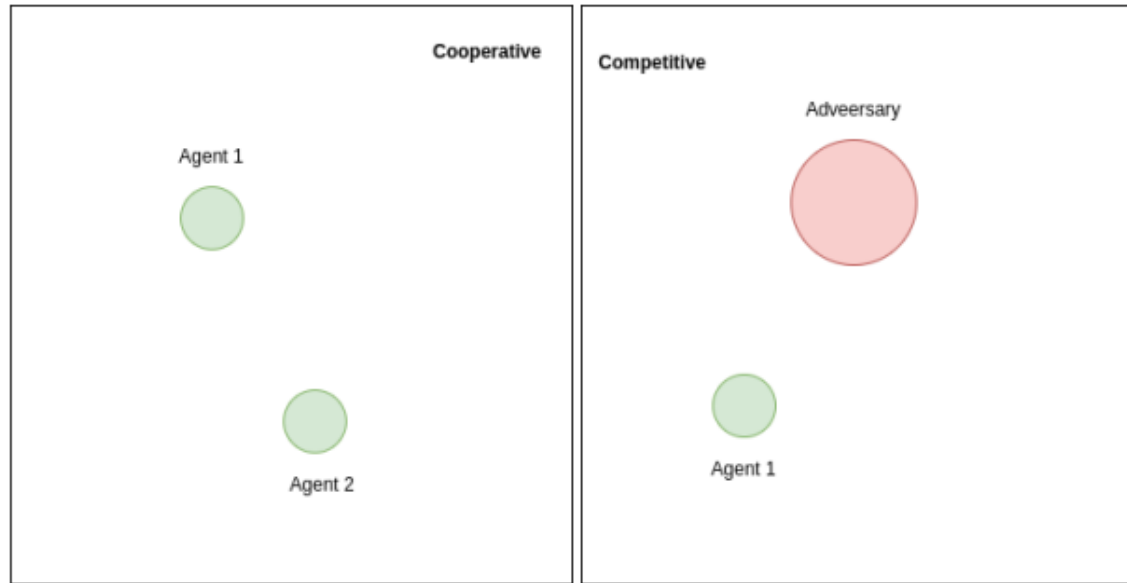
$$m_{k+1}^{(i)} = m_k^{(i)} + \mu_k^{(i)}$$

Algorithm 2 THE PROPOSED MAK-SR FRAMEWORK

- 1: **Learning Phase:**
 - 2: **Initialize:** $\theta_0, P_{\theta,0}, m_0, P_{M,0}, \mu_n$, and Σ_n for $n = 1, 2, \dots, N$
 - 3: **Parameters:** $Q_\theta, Q_M, \lambda_\mu, \lambda_\Sigma$, and $\{R_\theta^j, R_M^j\}$ for $j = 1, 2, \dots, M$
 - 4: **Repeat** (for each episode):
 - 5: Initialize s_k
 - 6: **Repeat** (for each agent i):
 - 7: **While** $s_k^{(i)} \neq s_T$ **do:**
 - 8: Reshape m_k into $L \times L$ to construct 2-D matrix M_k .
 - 9: $a_k^{(i)} = \arg \max_a \left(g_k^{(i)}(s_k^{(i)}, a) g_k^{(i)T}(s_k^{(i)}, a^{(i)}) \right)$
 - 10: Take action $a_k^{(i)}$, observe $s_{k+1}^{(i)}$ and $r_k^{(i)}$.
 - 11: Calculate $\phi(s_k^{(i)}, a_k^{(i)})$ via Equations (23) and (25).
 - 12: **Update reward weights vector:** Perform MMAE to update $\theta_k^{(i)}$.
 - 13: **Update SR weights vector:** Perform KF on Equations (40) and (41) to update $m_k^{(i)}$.
 - 14: **Update RBFs parameters:** Perform RGD on the loss function L_k to update Σ_n and μ_n .
 - 15: **end while**
-

Experimental Results

Experimental Results



(a)

(b)



(c)

(d)

Experimental Results

Table 1. Total loss averaged across all the episodes and for all the four implemented scenarios.

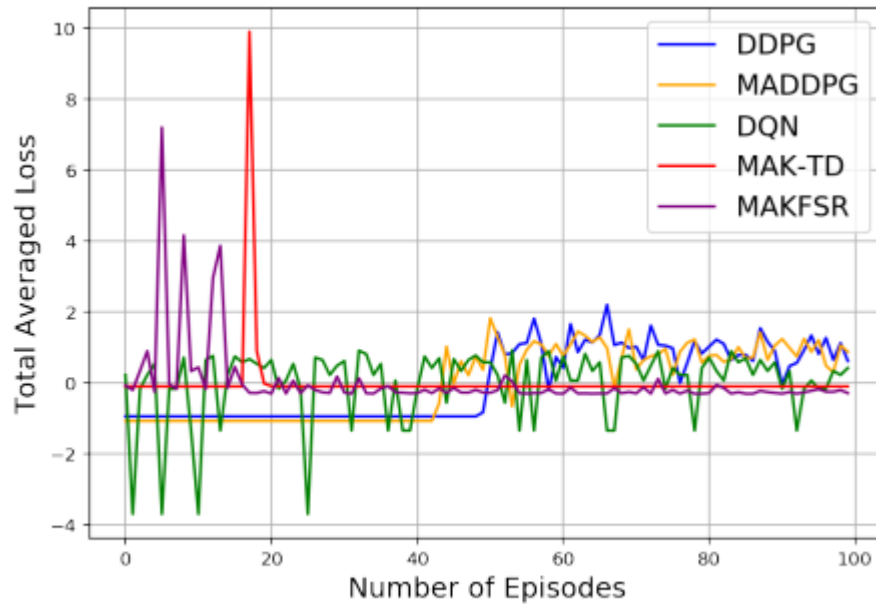
Environment	MAK-SR	MAK-TD	MADDPG	DDPG	DQN
Cooperation	8.93	2.4088	9649.84	10,561.16	10.93
Competition	0.43	4.9301	10,158.18	10,710.37	107.39
Predator–Prey 1v2	0.005	1.9374	6816.34	6884.33	8.21
Predator–Prey 2v1	8.87	1.2421	7390.18	6882.2	10.24

Experimental Results

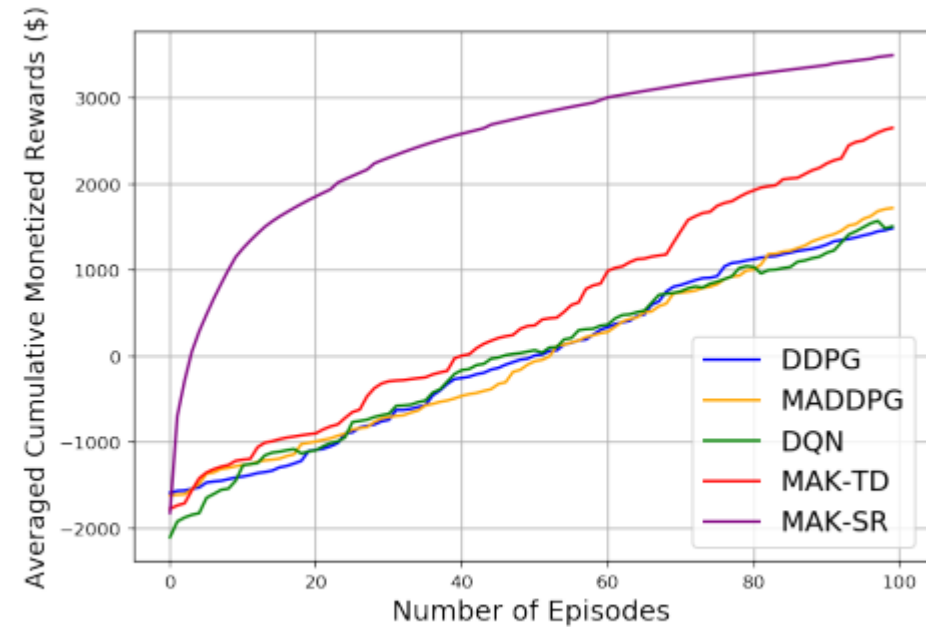
Table 2. Total received reward by the agents averaged for all the four implemented scenarios.

Environment	MAK-SR	MAK-TD	MADDPG	DDPG	DQN
Cooperation	−16.0113	−23.0113	−69.28	−66.29	−39.96
Competition	−0.778	−13.358	−63.30	−61.34	−14.49
Predator–Prey 1v2	−0.0916	−13.432	−46.17	−20.53	−23.451
Predator–Prey 2v1	−0.081	−17.0058	−55.69	−49.41	−44.32

Experimental Results

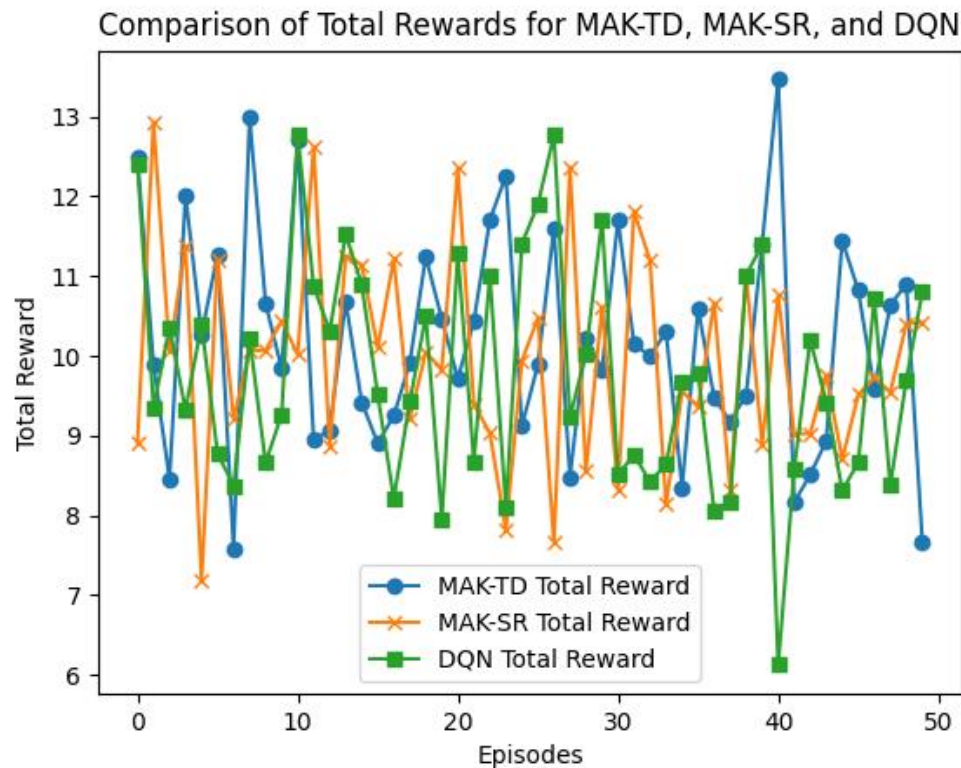


(a)



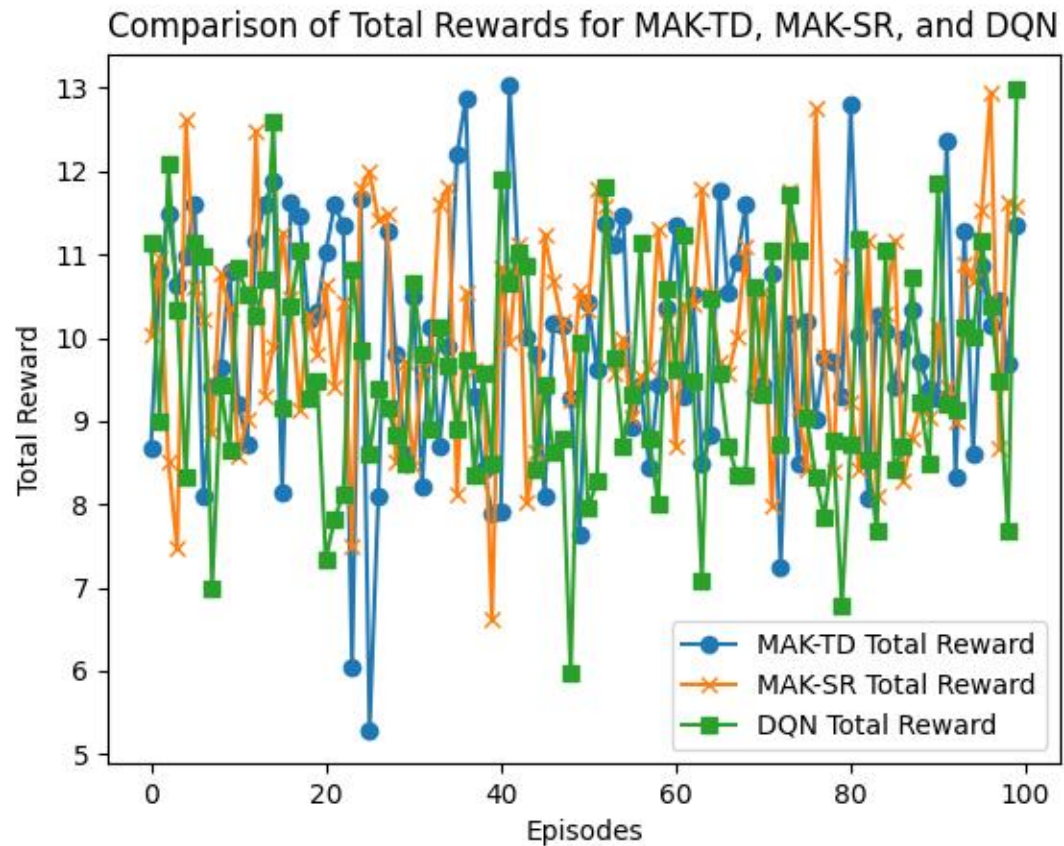
(b)

Simulation Results



Final Total Rewards after 50 Episodes:
MAK-TD Total Reward: 510.10918848452957
MAK-SR Total Reward: 498.1273978229948
DQN Total Reward: 488.53212117234517

Simulation Results



Final Total Rewards after 100 Episodes:
MAK-TD Total Reward: 996.1441199350996
MAK-SR Total Reward: 1003.7453232643612
DQN Total Reward: 955.8240584992651

Key Achievements:

- MAK-TD
- MAK-SR

CONCLUSIONS

THANK YOU!

Winter 2025
