

SMIX(λ): Enhancing Centralized Value Functions for Cooperative Multi-Agent Reinforcement Learning (Supplementary)

Chao Wen, Xinghu Yao, Yuhui Wang, Xiaoyang Tan

A. Proof of Theorem 1

Theorem 1. For QMIX, if $\frac{\partial Q_{tot}}{\partial Q^i} \geq 0$ for $i \in \{1, 2, \dots, N\}$, then we have

$$\max_{\mathbf{a}} Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) = Q_{tot}(\boldsymbol{\tau}, \operatorname{argmax}_{a^1} Q^1(\tau^1, a^1), \dots, \operatorname{argmax}_{a^N} Q^N(\tau^N, a^N)) \quad (1)$$

Proof. For QMIX, we have

$$Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) := \hat{Q}(Q_1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N)),$$

where \hat{Q} is the *mixing network* and $\mathbf{a} = (a_1, \dots, a_N)$. Similarly, we have $\frac{\partial \hat{Q}}{\partial Q^i} \geq 0$ and

$$Q_{tot}(\boldsymbol{\tau}, \operatorname{argmax}_{a^1} Q^1(\tau^1, a^1), \dots, \operatorname{argmax}_{a^N} Q^N(\tau^N, a^N)) := \hat{Q}\left(\max_{a^1} Q^1(\tau^1, a^1), \dots, \max_{a^N} Q^N(\tau^N, a^N)\right).$$

Since $\frac{\partial \hat{Q}}{\partial Q^i} \geq 0$, given $(\bar{a}^2, \dots, \bar{a}^N)$, we have

$$\hat{Q}(Q^1(\tau^1, a^1), Q^2(\tau^2, \bar{a}^2), \dots, Q^N(\tau^N, \bar{a}^N)) \leq \hat{Q}\left(\max_{a^1} Q^1(\tau^1, a^1), Q^2(\tau^2, \bar{a}^2), \dots, Q^N(\tau^N, \bar{a}^N)\right) \text{ for any } a^1.$$

Similarly, given $(\bar{a}^1, \dots, \bar{a}^{k-1}, \bar{a}^{k+1}, \dots, \bar{a}^N)$, we have

$$\hat{Q}(Q^1(\tau^1, \bar{a}^1), \dots, Q^k(\tau^k, a^k), \dots, Q^N(\tau^N, \bar{a}^N)) \leq \hat{Q}\left(Q^1(\tau^1, \bar{a}^1), \dots, \max_{a^k} Q^k(\tau^k, a^k), \dots, Q^N(\tau^N, \bar{a}^N)\right) \text{ for any } a^k.$$

Finally, for any (a^1, \dots, a^N) , we have

$$\begin{aligned} & \hat{Q}(Q^1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N)) \\ & \leq \hat{Q}\left(\max_{a^1} Q^1(\tau^1, a^1), \max_{a^2} Q^2(\tau^2, a^2), Q^3(\tau^3, a^3), \dots, Q^N(\tau^N, a^N)\right) \\ & \leq \hat{Q}\left(\max_{a^1} Q^1(\tau^1, a^1), \max_{a^2} Q^2(\tau^2, a^2), \dots, \max_{a^N} Q^N(\tau^N, a^N)\right). \end{aligned}$$

Therefore, we obtain

$$\max_{a^1, \dots, a^N} \hat{Q}(Q^1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N)) = \hat{Q}\left(\max_{a^1} Q^1(\tau^1, a^1), \max_{a^2} Q^2(\tau^2, a^2), \dots, \max_{a^N} Q^N(\tau^N, a^N)\right),$$

which is the specific form of (1) for QMIX. □

B. Proof of Theorem 2

Theorem 2. Suppose we update the value function from $Q_n^{SMIX}(\boldsymbol{\tau}_t, \mathbf{a}_t) = Q_n^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t)$, where n represents the n -th update. Let $\epsilon = \max_{\boldsymbol{\tau}} \|\boldsymbol{\pi}(\cdot|\boldsymbol{\tau}) - \boldsymbol{\mu}(\cdot|\boldsymbol{\tau})\|_1$, $M = \max_{\boldsymbol{\tau}, \mathbf{a}} |Q_n^{Q(\lambda)}(\boldsymbol{\tau}, \mathbf{a})|$. Then, the error between $Q_{n+1}^{SMIX}(\boldsymbol{\tau}_t, \mathbf{a}_t)$ and $Q_{n+1}^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t)$ can be bounded by the expression:

$$|Q_{n+1}^{SMIX}(\boldsymbol{\tau}_t, \mathbf{a}_t) - Q_{n+1}^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t)| \leq \frac{\epsilon\gamma}{1 - \lambda\gamma} M. \quad (2)$$

Proof. First, we have,

$$\begin{aligned} |\delta_t^\pi - \hat{\delta}_t^\pi| &= |\gamma \mathbb{E}_{\boldsymbol{\mu}} Q_n^{SMIX}(\boldsymbol{\tau}_{t+1}, \cdot) - \gamma \mathbb{E}_{\boldsymbol{\pi}} Q_n^{Q(\lambda)}(\boldsymbol{\tau}_{t+1}, \cdot)| \\ &= \gamma \left| \sum_{\mathbf{a}} \boldsymbol{\mu}(\mathbf{a}|\boldsymbol{\tau}_{t+1}) Q_n^{SMIX}(\boldsymbol{\tau}_{t+1}, \cdot) - \sum_{\mathbf{a}} \boldsymbol{\pi}(\mathbf{a}|\boldsymbol{\tau}_{t+1}) Q_n^{Q(\lambda)}(\boldsymbol{\tau}_{t+1}, \cdot) \right| \leq \gamma \epsilon M. \end{aligned}$$

Thus,

$$\begin{aligned}
|Q_{n+1}^{\text{SMIX}}(\tau_t, \mathbf{a}_t) - Q_{n+1}^{Q(\lambda)}(\tau_t, \mathbf{a}_t)| &= \left| \mathbb{E}_{\mu} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) \delta_k^{\pi} \right] - \mathbb{E}_{\mu} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) \hat{\delta}_k^{\pi} \right] \right| \\
&= \left| \mathbb{E}_{\mu} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) (\delta_k^{\pi} - \hat{\delta}_k^{\pi}) \right] \right| \leq \left| \mathbb{E}_{\mu} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) (\gamma \epsilon M) \right] \right| \\
&\leq \mathbb{E}_{\mu} \left[\frac{1}{1 - \lambda \gamma} (\gamma \epsilon M) \right] = \frac{\epsilon \gamma}{1 - \lambda \gamma} M.
\end{aligned}$$

Therefore, the expression (2) holds. \square

C. Algorithm

Algorithm 1 Training Procedure for SMIX(λ)

- 1: Initialize the behavior network with parameters θ , the target network with parameters θ^- , empty replay buffer \mathcal{D} to capacity $N_{\mathcal{D}}$, training batch size b
 - 2: **for** each training episode **do**
 - 3: **for** each episode **do**
 - 4: **for** $t = 1$ to $T - 1$ **do**
 - 5: Obtain the partial observation $\mathbf{o}_t = \{o_t^1, \dots, o_t^N\}$ for all agents and global state s_t
 - 6: Select action a_t^i according to ϵ -greedy policy w.r.t agent i 's decentralized value function Q^i for $i = 1, \dots, N$
 - 7: Execute joint action $\mathbf{a}_t = \{a^1, a^2, \dots, a^N\}$ in the environment
 - 8: Obtain the global reward r_{t+1} , the next partial observation \mathbf{o}_{t+1}^i for each agent i and next global state s_{t+1}
 - 9: **end for**
 - 10: Store the episode in \mathcal{D} , replacing the oldest episode if $|\mathcal{D}| \geq N_{\mathcal{D}}$
 - 11: **end for**
 - 12: Sample a batch of b episodes $\sim \text{Uniform}(\mathcal{D})$
 - 13: Calculate λ -return targets $y_i^{\text{tot}} = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$, where $G_t^{(n)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n \mathbb{E}_{\pi} Q(\tau_{t+n}, \mathbf{a}_{t+n}; \theta^-)$
 - 14: Update θ by minimizing $\sum_{t=1}^{T-1} \sum_{i=1}^b [(y_i^{\text{tot}} - Q_{\text{tot}}^{\pi}(\tau, \mathbf{a}; \theta))^2]$
 - 15: Replace target parameters $\theta^- \leftarrow \theta$ every C episodes
 - 16: **end for**
-

D. Implementation Details

The agent network consists of a 64-dimensional GRU (Chung et al. 2014). One 64-dimensional fully connected layer with ReLU activation function before GRU is applied for processing the input. The layer after GRU is a fully connected layer of 64 units, which outputs the decentralized state-action values $Q^i(\tau^i, \cdot)$ of agent i . All agent networks share parameters for reducing the number of parameters to be learned. Thus the agent's one-hot index i is concatenated onto each agent's observations. Agent's previous action is also concatenated to the input. The target network is updated after every 200 training episodes. We use RMSprop optimizer with learning rate 0.0005 and $\alpha = 0.99$ without weight decay or momentum during training. We perform independent ϵ -greedy for exploration. ϵ is annealed linearly from 1.0 to 0.05 across the first 50k timesteps for all experiments. The discount factor is set to $\gamma = 0.99$. The architecture of the mixing network is the same as in Rashid et al. (2018). Actually, all the settings of SMIX(λ)'s parameters are kept the same as those of QMIX for fair comparison, except that SMIX(λ) has an extra parameter λ and uses smaller buffer size.

E. Additional Results

We provide quantitative comparisons of our methods and their counterparts in Table 1.

Algorithms	3m		8m		2s3z		3s5z		2s_vs_1sc		3s_vs_3z	
	mean \pm std	median	mean \pm std	median	mean \pm std	median	mean \pm std	median	mean \pm std	median	mean \pm std	median
SMIX(λ)	99 (± 0)	99	91 (± 3)	90	90 (± 4)	91	61 (± 11)	62	94 (± 5)	96	84 (± 14)	88
QMIX	95 (± 3)	95	90 (± 3)	89	81 (± 7)	81	16 (± 12)	11	39 (± 19)	45	15 (± 20)	9
SMIX(λ)-COMA	93 (± 8)	97	92 (± 2)	93	44 (± 18)	47	0 (± 0)	0	97 (± 4)	100	0 (± 0)	0
COMA	92 (± 2)	93	90 (± 2)	91	24 (± 6)	24	0 (± 0)	0	77 (± 11)	78	0 (± 0)	0
SMIX(λ)-VDN	98 (± 0)	98	94 (± 3)	93	78 (± 14)	79	29 (± 12)	26	96 (± 2)	97	67 (± 25)	83
VDN	95 (± 2)	95	86 (± 5)	87	64 (± 16)	71	1 (± 2)	0	86 (± 8)	88	27 (± 9)	27
SMIX(λ)-IQL	91 (± 4)	94	80 (± 5)	79	32 (± 8)	31	0 (± 0)	0	92 (± 6)	94	35 (± 21)	31
IQL	83 (± 9)	86	59 (± 15)	58	14 (± 10)	13	0 (± 0)	0	51 (± 22)	54	5 (± 4)	6

Table 1: Mean, standard deviation, and median of test win rate percentages after training for 1 million timesteps in six different scenarios. The highest results are in bold. Our methods (SMIX(λ), SMIX(λ)-COMA, SMIX(λ)-VDN, SMIX(λ)-IQL) perform on par or significantly better than their counterparts in terms of the learning speed across all scenarios.

References

- Chung, J.; Gulcehre, C.; Cho, K.; and Bengio, Y. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. In *Advances in Neural Information Processing Systems*.
- Rashid, T.; Samvelyan, M.; Witt, C. S.; Farquhar, G.; Foerster, J.; and Whiteson, S. 2018. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, 4292–4301.