

Value Function Transfer for Deep Multi-Agent Reinforcement Learning Based on N-Step Returns

Liu et al.

International Joint Conference on Artificial Intelligence(IJCAI) 2019

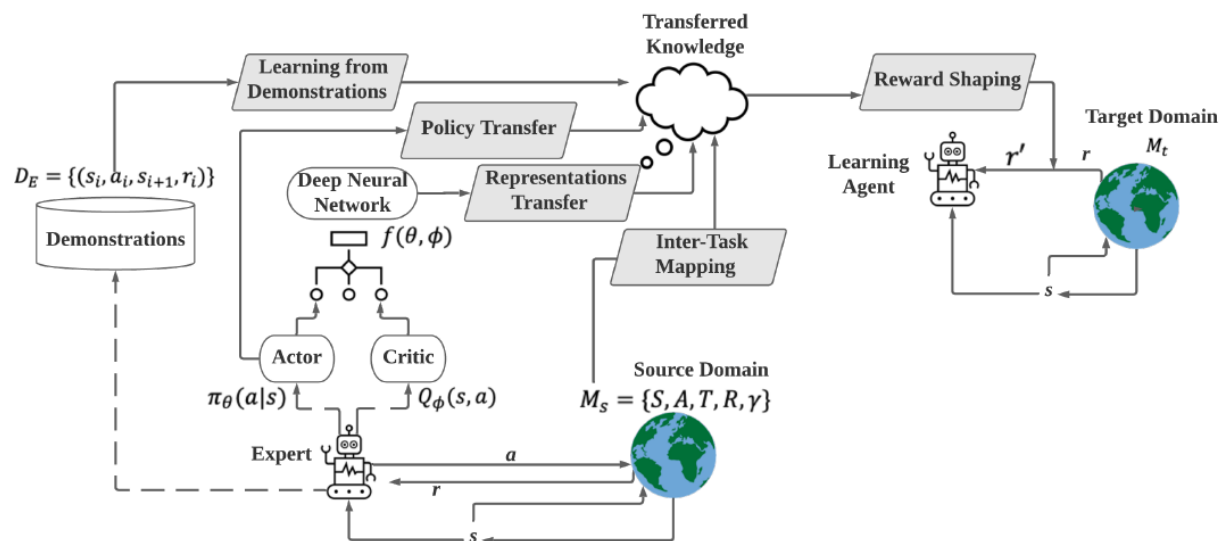
Speaker : Min Joon Kim

Feb 28th, 2024

Introduction

Transfer learning(TL)

- **(Definition)** Transfer learning is a technique in [machine learning](#) in which knowledge learned from a task is re-used in order to boost performance on a related task.
- **TL in RL**
 - (1) Learning from demonstrations (2) Policy Transfer (3) Representations transfer
 - (4) Inter-task Mapping (5) Reward shaping



"Transfer learning in deep reinforcement learning: A survey" 2023

Problem definition

- Value function transfer for deep multi-agent RL
 - Want to accelerate the multi-agent learning process with single agent knowledge.

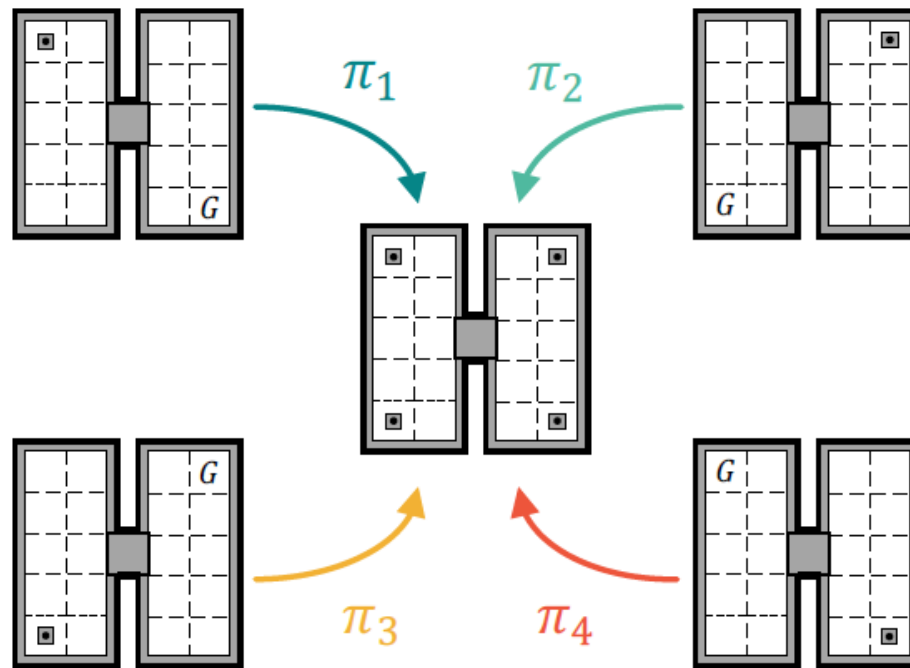


Figure 1: Transfer single-agent knowledge to multi-agent system

Key ideas

- **Propose the two knowledge transfer methods**
 - Direct value function transfer
 - N-step Return(NSR) based value function transfer
 - Using NSR value representing the MDP similarity between single-agent and multi-agent environment
 - Experiments are conducted in grid world, multi-agent particle environment and Ms. Pac-Man game.
 - Sparse interaction environment (independent agents)

Related works

Year	Paper	Note
2010	Learning multi-agent state space representations	Knowledge transfer in multi-agent / <u>Tabular domain</u>
2011	Transfer learning for multi-agent coordination	
2015	Learning in multi-agent systems with sparse interactions by knowledge transfer and game abstraction.	
		MDP similarity
2018	Value-decomposition networks for cooperative multi-agent learning based on team reward.	VDN
2018	QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning.	QMIX

- **Previous works**

- Can be adopted only in tabular domain
- Not much studies in value function transfer

- **Contribution**

- Proposed value function transfer methods
- Suggest a new MDP similarity metric using N-step return value

Direct Value Function Transfer

- Direct Value Function Transfer Network Architecture
 - Model A : Single-agent expert policy
 - Model B : Multi-agent network(QMIX)

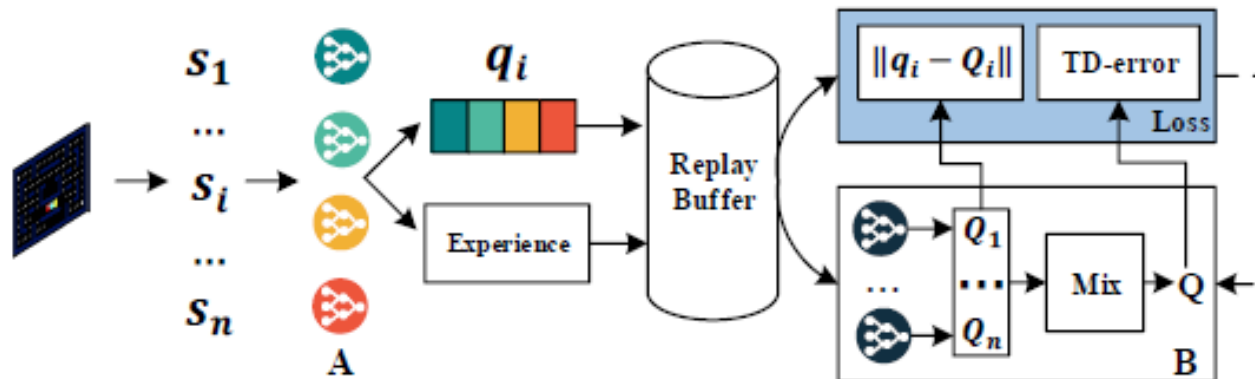


Figure 2: Direct Value Function Transfer Network Architecture. Model A represents single-agent expert policy network and model B represents multi-agent network.

Direct Value Function Transfer

- Pseudo code

Algorithm 1: Direct Value Function Transfer

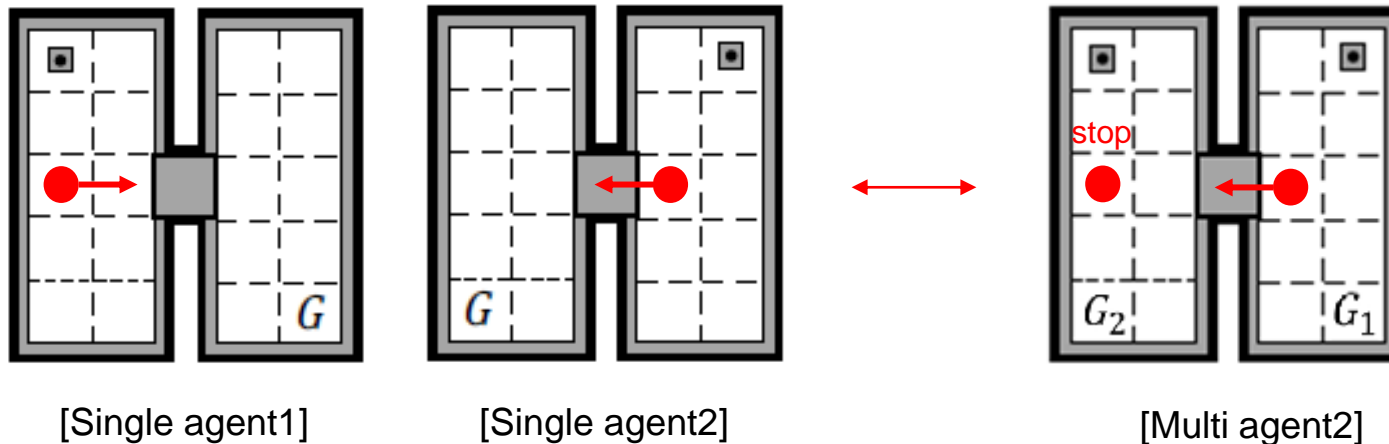
Input: local value function $q_i(s_i, a)$ for each agent i , discount factor γ , exploration factor ϵ

```
1 Initialization.  $Q(s, \vec{a}) \leftarrow \theta, \hat{Q}(s, \vec{a}) \leftarrow \theta', Q_i(s_i, \vec{a}) \leftarrow \theta_i;$ 
2 foreach episode do
3   Initialize state  $s;$ 
4   repeat
5     foreach agent  $i$  do
6        $(s_i, a_i) \leftarrow$  the component of agent  $i$  in  $(s, \vec{a});$ 
7        $a_i \leftarrow$  max  $q_i$  index with  $\epsilon$ -greedy policy;
8      $\vec{a} \leftarrow [a_1, \dots, a_n];$ 
9     store experience  $(s, \vec{a}, r, s', done, [q_1, \dots, q_n]);$ 
10     $s \leftarrow s';$ 
11    Sample training experience from buffer;
12    if not done then
13       $y = r + \gamma \hat{Q}_{max}(s, \vec{a}; \theta');$ 
14    else
15       $y = r;$    Q-value from single agent(expert) – Q-value for agent  $i$ 
16       $L(\theta) = \alpha \sum_1^N (q_i(s_i, a_i) - Q_i(s_i, a_i; \theta_i))^2 + (1 -$ 
17         $\alpha)(y - Q(s, \vec{a}; \theta))^2;$  (TD loss for the global Q value)
18      Update  $\theta$  by a gradient method w.r.t.  $L(\theta);$ 
19      Every  $C$  steps reset  $\hat{Q} = Q;$ 
20  until  $s = terminal;$ 
```

Store experience with Single-agents policy

Limitation of Direct Value Function Transfer

- Limitation and MDP similarity idea
 - What if the state is different between single-agent and multi-agent environment?
 - It causes negative transfer.



- So, they want to transfer only when the MDP is similar.
 - MDP similarity calculation : single agent Q value - N-step reward(multi-agent)

N-step Return based Value Function Transfer

▪ Pseudo code

Algorithm 2: NSR-based Value Function Transfer

Input: local value function $q_i(s_i, a)$ and single-agent policy π_i for each agent i , N ($N \geq 1$), discount factor γ , exploration factor ϵ, δ

```
1 Initialization. NSR value function  $\hat{\mathcal{R}}_{\pi_i}^N \leftarrow \psi_i$ , Replay Buffer  $\mathcal{B}_i$  for each agent  $i$  ;
2 foreach episode do
3   Initialize state  $s_t, t = 0$ ;
4   repeat
5     foreach agent  $i$  do
6        $(s_{i,t}, a_i) \leftarrow$  the state of agent  $i$  in  $(s_t, \vec{a})$ ;
7        $a_i \leftarrow$  max  $q_i$  index with probability  $\delta$ ;
8        $y_i = r_{t-N+1} + \gamma r_{t-N} + \dots + \gamma^{N-1} r_t$ ;
9       Store  $(s_{i,t-N+1}, y_i)$  in  $\mathcal{B}_i$ ;
10      Update  $\psi_i$  by a gradient method w.r.t.
11         $(y_i - \hat{\mathcal{R}}_{\pi_i}(s_{i,t-N+1}; \psi))^2$ ;
12       $t = t + 1$ 
13    until  $s_t = \text{terminal}$ ;
```

- Update the N-step return prediction network (Multi-agent environment)

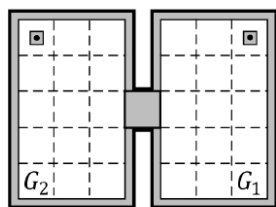
```
13 foreach episode do
14   Initialize state  $s$ ;
15   repeat
16      $(s_i, a_i) \leftarrow$  the state of agent  $i$  in  $(s, \vec{a})$ ;
17     if  $\max_i |\hat{\mathcal{R}}_{\pi_i}(s_i) - q_i(s_i, \pi_i(s_i))| \leq \tau$  then
18        $a_i \leftarrow$  argmax  $q_i(s_i, a_i)$  for each agent  $i$ ;
19     else
20        $a_i \leftarrow$  argmax  $Q_i(s, a_i)$  with  $\epsilon$  greedy for each agent  $i$ ;
21     Learning by multi-agent method;
22   until  $s = \text{terminal}$ ;
```

- Similarity threshold τ
- Transfer (models are similar)
- No transfer (models are different)

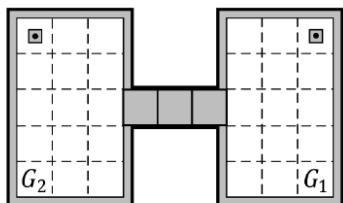
Experiments

Experiment settings

- Each single agent expert policy is trained by DQN.
- 3 environments (Grid world, Multi-agent Particle, Ms. Pacman)
- Similarity threshold $\tau = 2, 3, 5, 7$

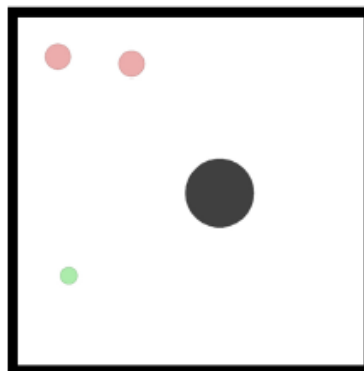


(a) map 1



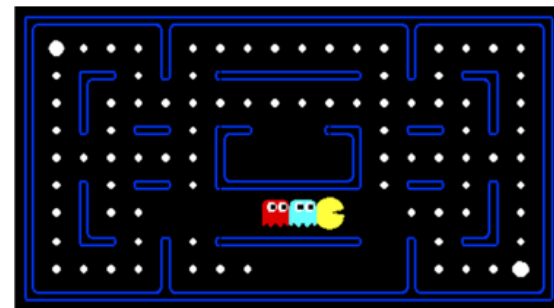
(b) map 2

Grid world



(a) predator-prey

Multi-agent Particle Environment (MPE)

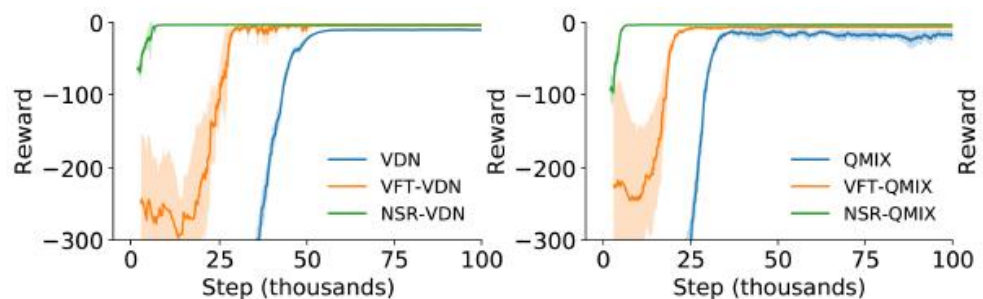


Ms. Pac-Man

Results

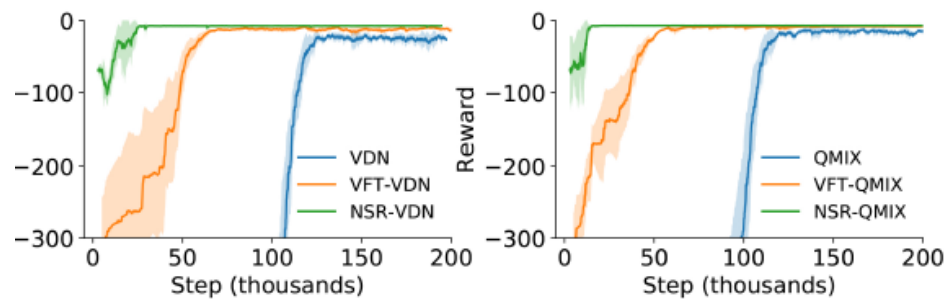
▪ Grid world

- NSR-based transfer method showed more stable and faster convergence than others.
- QMIX was slightly better than VDN algorithm.
- About double steps are needed for the second environment(map2) learning.



(a) map1-VDN

(b) map1-QMIX

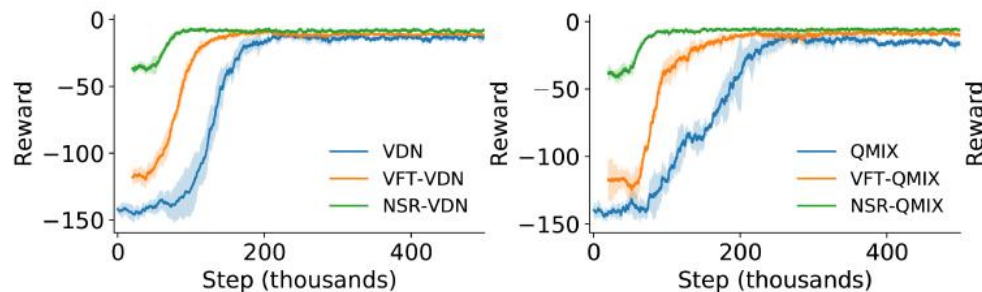


(c) map2-VDN

(d) map2-QMIX

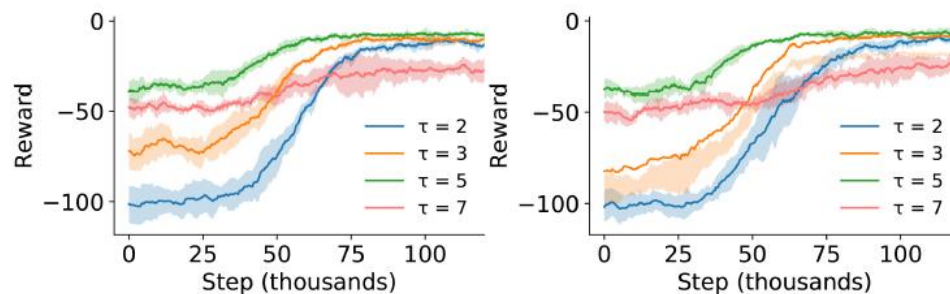
Results

- Multi-agent Particle Environment(MPE)
 - NSR-based transfer method showed much faster convergence than others.
 - Proposed method converges in 90,000 steps, while the VDN needs 200,000 steps.
 - They found the appropriate threshold($\tau = 5$) experimentally in figure (g), (h)



(e) MPE-VDN

(f) MPE-QMIX



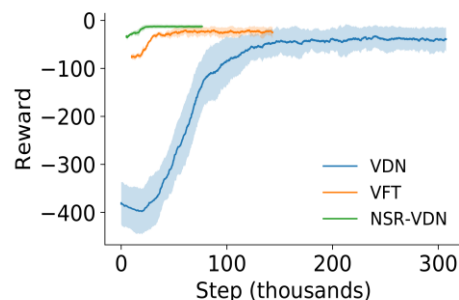
(g)

(h)

Results

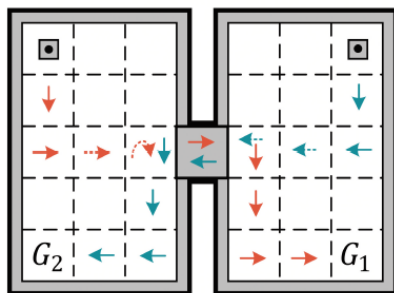
- Ms. Pac-Man environment

- Proposed method showed best performance, but there is not much improvement compared to the direct value function transfer method.

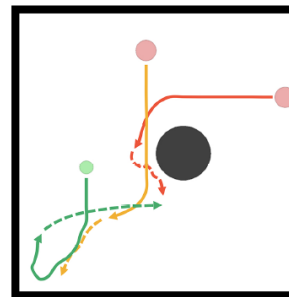


- Learned policy (Grid world, MPE environment)

- Learned by single-agent(solid arrow), Learned by multi-agent(dotted arrow)



(b) Policy Display



(b) Policy Display

Conclusion

▪ **Conclusion**

- They suggested value function transfer method for the Multi-agent RL.
- Proposed method showed faster and stable convergence than MARL algorithms.
- Achieved about 2 to 5 times faster convergence than existing models.

▪ **Limitations**

- Only for sparse interactions (independent agents)
- Similarity calculation makes the additional burden in the training steps.