

A new Potential-Based Reward Shaping for Reinforcement Learning Agent

Babak Badnava¹, Mona Esmaeili², Nasser Mozayani³, and Payman Zarkesh-Ha²

¹University of Kansas, Lawrence, KS 66045, USA.

²University of New Mexico, Albuquerque, NM 87131, USA.

³Iran University of Science and Technology, Tehran, 16846, Iran.

Abstract—Potential-based reward shaping (PBRS) is a particular category of machine learning methods that aims to improve the learning speed of a reinforcement learning agent by extracting and utilizing extra knowledge while performing a task. There are two steps in the transfer learning process: extracting knowledge from previously learned tasks and transferring that knowledge to use it in a target task. The latter step is well discussed in the literature, with various methods being proposed for it, while the former has been explored less. With this in mind, the type of knowledge that is transmitted is very important and can lead to considerable improvement. Among the literature of both transfer learning and potential-based reward shaping, a subject that has never been addressed is the knowledge gathered during the learning process itself. In this paper, we presented a novel potential-based reward shaping method that attempted to extract knowledge from the learning process. The proposed method extracts knowledge from episodes' cumulative rewards. The proposed method has been evaluated in the Arcade learning environment, and the results indicate an improvement in the learning process in both the single-task and the multi-task reinforcement learner agents.

Index Terms—Potential-based Reward Shaping, Reinforcement Learning, Reward Shaping, Knowledge Extraction

I. INTRODUCTION

In reinforcement learning (RL) problems, an agent learns to maximize a reward signal, which may be time-delayed [1]. The RL framework has gained much success in handling complex problems in last years. However, mastering difficult tasks is often slow. Thus most of the RL researcher focuses on improving the speed of the learning process by using the knowledge that provided by an expert or a heuristic function.

Transfer learning (TL) is one of the approaches that try to improve the speed of learning by using knowledge extracted from previously learned tasks [2]. Reward shaping (RS) is also one of the methods that have been used for transferring such knowledge or any other type of knowledge.

In problems that could be modeled using RL, many tasks have kind of sparse reward signal. For example, there are some tasks that an agent will receive nothing until it gets to a goal or it gets nothing until some events occur in the environment. The reward function of a task in an environment is independent of time, and the agent tries to maximize the sum of its discounted reward in an episode, while from an episode to another episode of learning, the agent could use a knowledge extracted from past episodes to reinforce the reward signal. There is much information that can be used to reinforce the reward signal. For example, after one episode of learning has done, the agent could compare its performance with the best episode of the learning or the worst one, and by this comparison, it could improve the learning process of itself.

II. BACKGROUND

A. Reinforcement learning

RL is a group of machine learning methods, which can be used for training an agent using environmental feedback. An RL agent in an environment can be modeled as a Markov decision process(MDP) [1]. An MDP is a tuple like (S, A, T, γ, R) , and in this tuple: S is a set of states that appear in the environment; A is a set of agent's possible actions; $T : S \times A \times S \rightarrow [0, 1]$ is a function which gives probability of reaching state s' while agent currently is in state s and chose action a to perform; $R : S \times A \times S \rightarrow \mathbb{R}$ is a function which gives amount of reward the agent will receives when performed the action a in state s and transferred to state s' ; γ called discount factor which indicates how future rewards are important for the agent.

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

The agent wants to maximize equation 1, which is called the expected discounted return, during its lifetime [1]. There are many methods for training RL agents.

Methods such as Q-learning[3] and SARSA which estimate a value for each (*state, action*) pair and by using this estimated value compute a policy that can maximize the agent’s discounted return. There are also some methods that use policy gradient [4–6] to approximate a policy and use this policy for the agent’s decision making.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (2)$$

B. Reward shaping

Reward shaping (RS) refers to allowing an agent to train on an artificial reward signal rather than environmental feedback [7]. However, this artificial signal must be a potential function. Otherwise, the optimal policy will be different for the new MDP [8]. Authors of [8] prove that F is a potential function if there exists a real-valued function $\phi : S \rightarrow \mathbb{R}$ such that for all $s \in S - \{s_0\}, a \in A, s' \in S$:

$$F(s, s') = \gamma \phi(s') - \phi(s) \quad (3)$$

Consequently by using this potential function, the optimal policy in the new MDP ($\mathcal{M}' = (S, A, T, \gamma, R + F)$) remains optimal for the old MDP ($\mathcal{M} = (S, A, T, \gamma, R)$). In the new MDP, equation 3, and all of the extended equations from equation 3, γ is unchanged and must have the same value as in the old MDP. There are also other works in the literature that extend the potential function. Authors of [9] have shown that the potential function can be a function of the joint state-action space. Next, Authors of [10] have shown that the potential function could be a dynamic function and that might change over the time while all the properties of potential based reward shaping (PBRS) remain constant. In addition to [9] and [10], [11] combines these two extensions and show that any function in the form of equation 6 could be a potential function.

$$F(s, a, s', a') = \gamma \phi(s', a') - \phi(s, a) \quad (4)$$

$$F(s, t, s', t') = \gamma \phi(s', t') - \phi(s, t) \quad (5)$$

$$F(s, a, t, s', a', t') = \gamma \phi(s', a', t') - \phi(s, a, t) \quad (6)$$

By using reward shaping Q-learning update rule will turn into equation 7. In equation 7, F could be any of equations represented in 3, 4, 5, or 6. Equation 7 will give us this opportunity to enhance the learning by providing some extra information about the problem for the agent. This additional knowledge can come from any source like human knowledge of the problem, some heuristic function, or the knowledge that extracted by the agent.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + F + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (7)$$

III. RELATED WORKS

The first method, which should be considered, is the Multi-grid RL with RS, [12] which propose a method for learning a potential function in an online manner. Authors of [12] estimate a value function during the learning process and by using this value function shape a potential function. The estimated value function is the value of an abstract state. State abstraction could be applied by any method.

$$V(z) = (1 - \alpha_v) V(z) + \alpha_v (r_v + \gamma_v^t(z')) \quad (8)$$

$$F(s, s') = \gamma_v V(z') - V(z) \quad (9)$$

Another work that has been done in the literature is [13], which proposed a potential function based on a plan. In [13], the agent gets an extra reward based on progress in the plan. In equation 10, z is an abstract state that can be any state of the plan, and the function $step(z)$ returns the time step at which given the abstract state z appears during the process of executing the plan.

$$\phi(s) = step(z) \quad (10)$$

Another work that is an extended version of [13], for multi-agent reinforcement learning, is [14]. In [14] two methods have been proposed, which can be used to shape the reward signal of agents and improve their learning process.

Another work that has been presented in the literature of transfer learning is a method that transfers the agent’s policy using reward shaping. Authors of [15] assume that there is a mapping function, which maps target task to source task. By using this mapping function, they define a potential function for shaping the reward signal. The proposed potential function as shown in equation 11 is defined based on the policy of source task. In equation 11, \mathcal{X}_S is a mapping function from target task state space to source task state space, and \mathcal{X}_A is a mapping function from target task action space to source task action space, and π is the policy of agent in the source task.

$$\phi(s, a) = \pi(\mathcal{X}_S(s), \mathcal{X}_A(a)) \quad (11)$$

In [15], RS has been used as a knowledge-transfer procedure, and a learned policy from another task has been used as knowledge. There are many other works in the literature that proposed a potential based reward shaping method for multi-agent reinforcement learning (MARL) and single-agent reinforcement learning (SARL). [16] proposed a method for RS, while there is a demonstration of a task. [16] proposed a method that uses a similarity measure for calculating the similarity between state-action pairs. Another work, which assumes a situation like that assumed in [16], is [17]. Authors of [17] instead

of using a similarity measure used an Inverse reinforcement learning approach to approximate a reward function for shaping the potential function. Another technique that is presented for MARL is [18]. In [18] difference reward used to shape a potential function. The difference reward helps the agent to learn the impact of its action on the environment by not considering the impact of other agents' action.

IV. PROPOSED POTENTIAL BASED REWARD SHAPING METHOD

Given these points and motivations, we go on to describe our approach to reinforce reward signal by using a knowledge gained from the learning process. So, we want a reward function that will change every time the agent make progress in the task. This reward function must encourage or punish the agent according to the best and worst episode until now. The reward function must also handle sparse reward signals. The function that represented in equation 12 has properties that we want. As we know from [8], changing reward function might change the optimal policy for the current task. Hence, we use reward shaping to manipulate the reward function in order to consider the agent's improvement during the learning process.

$$\phi(s, a, t) = \begin{cases} 0 & R(s, a) = 0 \\ 1 + \frac{R^{ep} - R_u^{ep}(t)}{R_u^{ep}(t) - R_l^{ep}(t)} & O.W \end{cases} \quad (12)$$

In equation 12: $R(s, a)$ is the immediate reward, which is a sparse reward signal; R^{ep} is the sum of rewards in the current episode, which we call it *episode reward*; $R_u^{ep}(t)$ is the maximum value of *episode reward* until now; and $R_l^{ep}(t)$ is the minimum value of *episode reward* until now. This function returns zero, while the reward signal of the environment has no information in it; the case to control the effect of sparse-reward. This function reinforces the reward signal by measuring a distance from a fixed point. We could use this approach with any kind of learning algorithm to boost the learning process. As presented in algorithm 1, after each episode of learning the agent will change parameters of potential function if needed.

How our proposed method works? By decomposition of equation 13 and considering equation 12, we can analyze how the proposed method works from a state to another state; the equation 17 is the result of this decomposition. During the learning process of an agent in an environment with sparse reward, three situations will appear by transmission from a state to another state: most of the time, the agent receives no reward and the episode reward will not change; sometimes the agent will receive a positive reward signal and it increases the

Algorithm 1: Q-learning + proposed PBRS method

```

 $\forall s, a \ Q(s, a) = 0$ 
 $s \leftarrow$  start state
 $max \ episode \ reward \leftarrow -\infty$ 
 $min \ episode \ reward \leftarrow \infty$ 
 $replay \ buffer \leftarrow null$ 
while True do
   $episode \ reward = 0$ 
  while  $s$  is not terminal do
     $a \leftarrow$ 
    Choose an action based on a policy and given state
     $r, s' \leftarrow$ 
    do action, get reward and observe next state
     $episode \ reward \leftarrow episode \ reward + r$ 
    if  $r = 0$  then
       $\phi_s \leftarrow 0$ 
    else
       $\phi_s \leftarrow 1 + \frac{episode \ reward - max \ episode \ reward}{max \ episode \ reward - min \ episode \ reward}$ 
    end
    append  $replay \ buffer$  ( $s, a, r, \phi_s, s'$ )
     $s \leftarrow s'$ 
    if update priode or  $s$  is terminal then
      for  $i \leftarrow$  length of  $replay \ buffer - 1$  to
        0 by -1 do
           $s, a, r, \phi_s, s' \leftarrow replay \ buffer[i]$ 
           $\rightarrow, \rightarrow, \phi_{s'}, \leftarrow \leftarrow$ 
           $replay \ buffer[i+1]$ 
           $F \leftarrow \gamma \phi_{s'} - \phi_s$ 
           $Q(s, a) \leftarrow Q(s, a) +$ 
           $\alpha [r + F + \gamma \max_a Q(s', a) - Q(s, a)]$ 
        end
       $replay \ buffer \leftarrow null$ 
    end
  end
  if  $episode \ reward > max \ episode \ reward$ 
  then
     $max \ episode \ reward \leftarrow episode \ reward$ 
  else if
     $episode \ reward < min \ episode \ reward$ 
  then
     $min \ episode \ reward \leftarrow episode \ reward$ 
  end
end

```

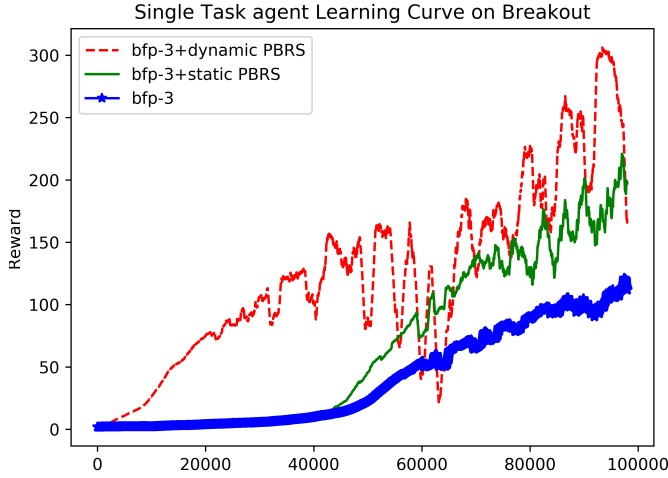


Fig. 1: single task agent learning curve on Breakout

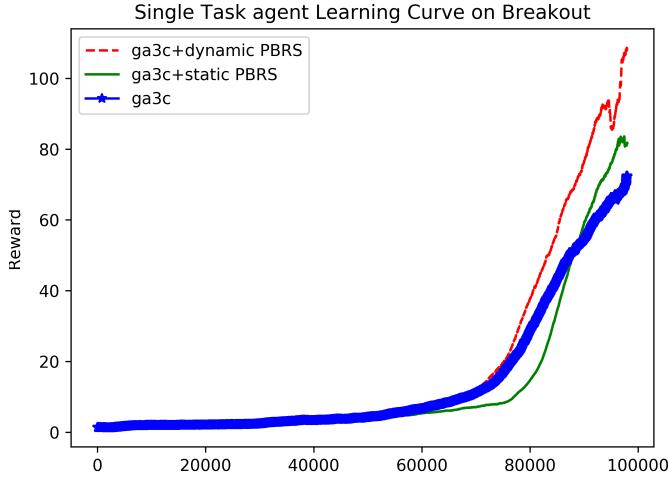


Fig. 2: single task agent learning curve on Breakout

episode reward; sometimes it receives a negative reward signal which causes a decrease in episode reward.

First of all, let's assume that by transmission from a state to another state no change in episode reward has been made. So, as we can see in equation 18, the potential function will be zero or negative. This will give the agent an information about the problem that means by remaining in the current state, it will not receive any reward and it also may lose something. Hence, this negative feedback will force the agent to avoid remaining in a state and to make an improvement from a state to another state. The more the agent stays in a state, it will receive more punishment. The more the agent in the vicinity of the goal and hesitate in a state, it will receive more punishment. It will lead to more exploration near the target and avoid sticking to sub-optimal goals.

How about the case when there is an increase in the *episode reward*? If all of the future rewards are equally valuable, the potential function will return zero every

time this situation appears to the agent. Whenever the episode reward has increased, as we can see in equation 19, the value of the potential function will be positive, and this positive value reinforces the environment reward.

Another situation that the agent might face is a decrease in the episode reward. Decreasing in the episode reward might be greater than a threshold or less than a threshold. While this decline is greater than a threshold (equation 21), the agent will receive a negative potential signal, which means the agent made a mistake or it spends its energy in a wrong way. Whenever the decline is less than the threshold (equation 22), it will still receive a positive potential signal, which provides the flexibility of spending more energy on improving the agent performance in the way of achieving the goal.

$$F(s, s') = \gamma\phi(s') - \phi(s) \quad (13)$$

$$= \gamma - 1 + \gamma \frac{R_u^{ep}(s') - R_u^{ep}(t)}{R_u^{ep}(t) - R_l^{ep}(t)} - \frac{R^{ep}(s) - R_u^{ep}(t)}{R_u^{ep}(t) - R_l^{ep}(t)} \quad (14)$$

$$= \gamma - 1 + \frac{\gamma R^{ep}(s') - R^{ep}(s) + R_u^{ep}(t) - \gamma R_u^{ep}(t)}{R_u^{ep}(t) - R_l^{ep}(t)} \quad (15)$$

$$= \gamma - 1 + \frac{\gamma R^{ep}(s') - R^{ep}(s) + R_u^{ep}(t)(1 - \gamma)}{R_u^{ep}(t) - R_l^{ep}(t)} \quad (16)$$

$$= \frac{\gamma R^{ep}(s') - R^{ep}(s) + R_l^{ep}(t)(1 - \gamma)}{R_u^{ep}(t) - R_l^{ep}(t)} \quad (17)$$

$$R^{ep}(s') = R^{ep}(s) \Rightarrow F(s, s') = \frac{(R_l^{ep}(t) - R^{ep}(t))(1 - \gamma)}{R_u^{ep}(t) - R_l^{ep}(t)} \leq 0 \quad (18)$$

$$\gamma R^{ep}(s') \geq R^{ep}(s) \Rightarrow F(s, s') > 0 \quad (19)$$

$$\gamma R^{ep}(s') < R^{ep}(s) \wedge \quad (20)$$

$$\gamma R^{ep}(s') - R^{ep}(s) > -R_l^{ep}(t)(1 - \gamma) \Rightarrow F(s, s') > 0 \quad (21)$$

$$\gamma R^{ep}(s') < R^{ep}(s) \wedge$$

$$\gamma R^{ep}(s') - R^{ep}(s) \leq -R_l^{ep}(t)(1 - \gamma) \Rightarrow F(s, s') \leq 0 \quad (22)$$

Extending to multi-task agents The presented potential function helps the learning methods to improve their performance by using knowledge extracted from previous episodes. Sometimes the agent needs to learn multiple tasks simultaneously. With this in mind, we extended the proposed method to be used in multi-task

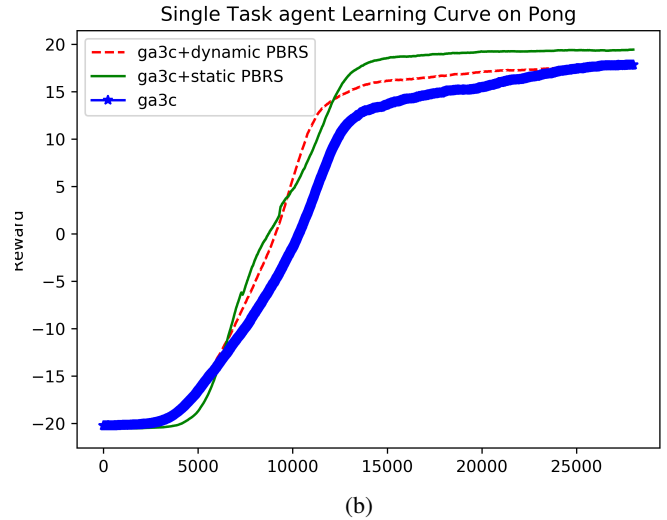
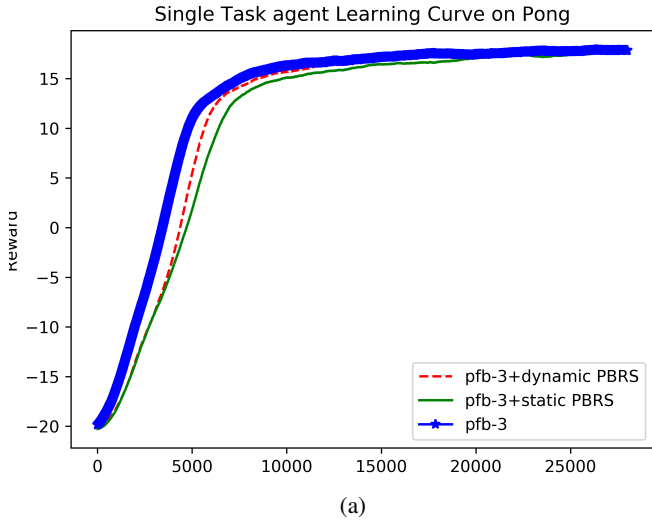


Fig. 3: single task agent learning curve on Pong

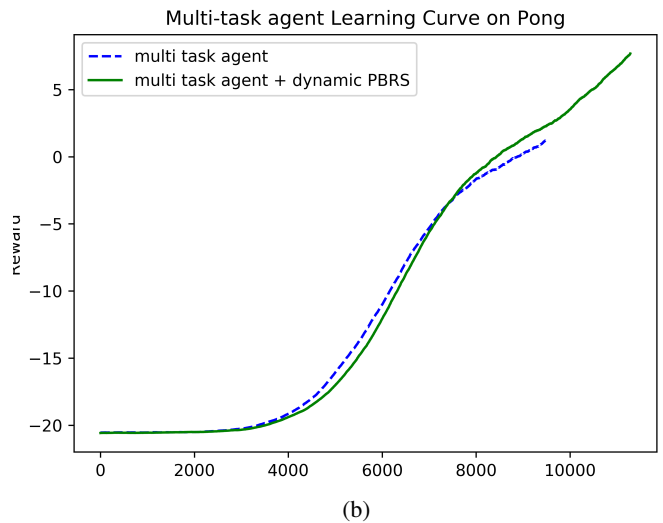
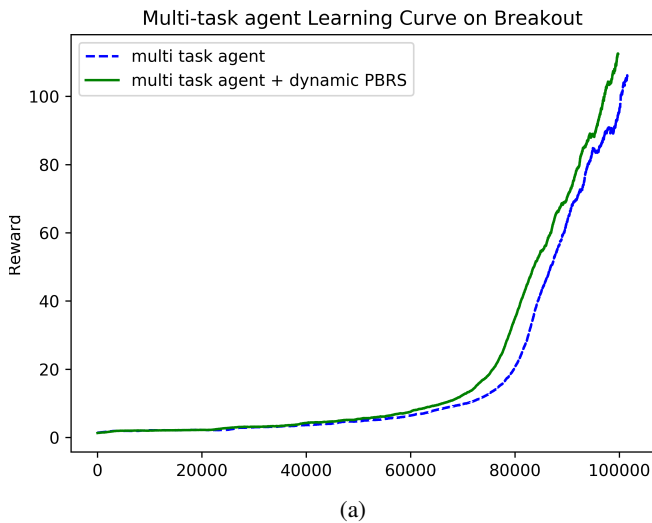


Fig. 4: multi-task agent learning curve on Pong and Breakout

RL. According to the equation 23, which is an extension of 12, all parameters of the potential function need to be calculated separately for each of the tasks. By using equation 23, each of the tasks has a different potential function. So, the policy of each task remains optimal. In equation 23, the i is an identifier for each of the tasks.

$$\phi(s, a, t, i) = \begin{cases} 0 & R(s, a) = 0 \\ 1 + \frac{R_u^{ep}(i) - R_{it}^{ep}(t, i)}{R_u^{ep}(t, i) - R_t^{ep}(t, i)} & O.W \end{cases} \quad (23)$$

V. RESULTS

In order to demonstrate the usefulness of the proposed approach, we perform experiments in two domains: Breakout and Pong games from arcade learning environment (ALE) [19]. According to the [19]:

“ALE provides an interface to hundreds of Atari 2600 game environments, each one different, interesting, and designed to be a challenge for human players. ALE presents significant research challenges for reinforcement learning, model learning, model-based planning, imitation learning, transfer learning, and intrinsic motivation.”

We used two different baselines to compare our proposed method with them. [20] used as a first baseline and [2] as a second baseline. [20] is an implementation of [21], which uses a policy gradient method to approximate the optimal policy of a specific task. As a second baseline, we implement [2] that presented in the literature of transfer learning. We evaluate our work in two different phases: first of all, we evaluate our method during the learning process, and then we evaluate the performance of the final policy. We tested our method



Fig. 5: multi-task agent policy evaluation on Pong and Breakout

with two different assumptions: the first assumption is that we know the maximum and minimum values of the episode reward in each task, and the second assumption is that we have no information about the episode reward of the task and this values will be obtained during the learning process.

Breakout: Breakout is an arcade game, which an agent handles a paddle to hit a ball trying to destroy more bricks while preventing the ball to cross the paddle. We trained an agent in this environment for 1000000 episodes. Figure 1 shows the learning curve of agents that trained on the Breakout game. As we see in figure 1 area under the curve of our method is greater than the baseline method, and the final performance is also better than the baseline.

Pong: Pong is also an arcade game, which an agent has the responsibility of moving a paddle to hit a ball. The agent will receive -1 reward if it loses the ball and will receive +1 if the opponent loses the ball. We trained an agent in this environment for 30000 episodes. Figure 3 shows the learning curve of agents that trained in the Pong game. As you see in figure 3 area under the curve of our method is greater than one of the baseline methods. However, in general, the proposed method in this environment has not been able to work very well and has a little improvement.

Multi-task Agent: We also implemented a multi-task agent, which learns how to play games of Pong and Breakout and then trains the agent for 115k episodes. We reinforced the reward signal using equation 23. Figure 4 and figure 5 demonstrate results of our experiments for a multi-task agent. Figure 4 demonstrates the learning curve of the multi-task agent for each of the games, while figure 5 demonstrates the performance of the learned-

policy. Figure 5 is the result of running the agent with the learned-policy for 100 episode and then take the average rewards of episodes. As we can see, the multi-task agent, which using reinforced reward signal, does perform better on average than a multi-task agent, which does not use reinforced reward signal.

VI. CONCLUSION

First of all, in this paper, we studied the literature on transfer learning and potential-based reward shaping. According to this study, there was no method that tries to extract knowledge from the learning process. Hence, we introduced a novel way of extracting knowledge from the learning process, and we used reward shaping as a knowledge transferring method.

The proposed method guides the agent toward a goal by looking at the value of episode reward. If the agent goes toward the goal, the reward signal will be reinforced. This method can be used with any learning algorithm, and it supposed to be efficient wherever applied. We implemented our method using [20] and compare the results with two different baselines in two different environments, and then we extend our method to be used in multi-task agents. In most of the experiments, the results were promising. We also tested our method in a multi-task agent, which tries to learn games of Pong and Breakout simultaneously, and we saw that our method leads to improvement.

REFERENCES

- [1] Sutton, Richard S. and Barto, Andrew G., "Introduction to Reinforcement Learning," *MIT Press*, 1998.
- [2] De la Cruz, Gabriel and Du, Yunshu and Irwin, James and Taylor, Matthew, "Initial Progress in Transfer for

- Deep Reinforcement Learning Algorithms,” in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2016.
- [3] Watkins, Christopher J.C.H. and Dayan, Peter, “Technical Note: Q-Learning,” *Machine Learning*, 1992.
- [4] Williams, Ronald J., “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine Learning*, 1992.
- [5] Sutton, Richard S. and McAllester, David and Singh, Satinder and Mansour, Yishay, “Policy Gradient Methods for Reinforcement Learning with Function Approximation,” in *Proceedings of the 12th International Conference on Neural Information Processing Systems*, 1999.
- [6] Peter L. Bartlett and Jonathan Baxter, “Infinite-Horizon Policy-Gradient Estimation,” *Journal of Artificial Intelligence Research*, 2001.
- [7] Taylor, Matthew E. and Stone, Peter, “Transfer Learning for Reinforcement Learning Domains: A Survey,” *JOURNAL OF MACHINE LEARNING RESEARCH*, 2009.
- [8] Ng, Andrew Y. and Harada, Daishi and Russell, Stuart J., “Policy Invariance Under Reward Transformations: Theory and Application to Reward Shaping,” in *Proceedings of the Sixteenth International Conference on Machine Learning*, 1999.
- [9] Wiewiora, Eric and Cottrell, Garrison and Elkan, Charles, “Principled Methods for Advising Reinforcement Learning Agents,” in *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, 2003.
- [10] Devlin, Sam and Kudenko, Daniel, “Dynamic Potential-based Reward Shaping,” in *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems - Volume 1*, 2012.
- [11] Harutyunyan, Anna and Devlin, Sam and Vrancx, Peter and Nowe, Ann, “Expressing Arbitrary Reward Functions As Potential-based Advice,” in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [12] Grześ, Marek and Kudenko, Daniel, “Multigrid Reinforcement Learning with Reward Shaping,” in *Proceedings of the 18th International Conference on Artificial Neural Networks, Part I*, 2008.
- [13] M. Grzes and D. Kudenko, “Plan-based reward shaping for reinforcement learning,” in *2008 4th International IEEE Conference Intelligent Systems*, 2008.
- [14] Sam Devlin and Daniel Kudenko, “Plan-based reward shaping for multi-agent reinforcement learning,” *Knowledge Eng. Review*, 2016.
- [15] Brys, Tim and Harutyunyan, Anna and Taylor, Matthew E. and Nowé, Ann, “Policy Transfer Using Reward Shaping,” in *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, 2015.
- [16] Brys, Tim and Harutyunyan, Anna and Suay, Halit Bener and Chernova, Sonia and Taylor, Matthew E. and Nowé, Ann, “Reinforcement Learning from Demonstration Through Shaping,” in *Proceedings of the 24th International Conference on Artificial Intelligence*, 2015.
- [17] Suay, Halit Bener and Brys, Tim and Taylor, Matthew E. and Chernova, Sonia, “Learning from Demonstration for Shaping Through Inverse Reinforcement Learning,” in *Proceedings of the 2016 International Conference on Autonomous Agents; Multiagent Systems*, 2016.
- [18] Sam Devlin and Logan Michael Yliniemi and Daniel Kudenko and Kagan Tumer, “Potential-based difference rewards for multiagent reinforcement learning,” in *AA-MAS*, 2014.
- [19] Bellemare, Marc G. and Naddaf, Yavar and Veness, Joel and Bowling, Michael, “The Arcade Learning Environment: An Evaluation Platform for General Agents,” *J. Artif. Int. Res.*, 2013.
- [20] Babaeizadeh, Mohammad and Frosio, Iuri and Tyree, Stephen and Clemons, Jason and Kautz, Jan, “Reinforcement Learning through Asynchronous Advantage Actor-Critic on a GPU,” in *5th International Conference on Learning Representations*, 2017.
- [21] Volodymyr Mnih and Adria Puigdomenech Badia and Mehdi Mirza and Alex Graves and Timothy Lillicrap and Tim Harley and David Silver and Koray Kavukcuoglu, “Asynchronous Methods for Deep Reinforcement Learning,” in *Proceedings of The 33rd International Conference on Machine Learning*, 2016.