Successor Feature for Transfer in Games

3rd ACC Workshop on Recent Advancement of Human Autonomy Interaction and Integration

Sunny Amatya Robotics and Intelligent Systems Laboratory, The Polytechnic School,

> Arizona State University May 30th, 2023





Introduction and Motivation





For safe interactions AV needs to able to efficiently infer the intent of other vehicles while being able to adapt to new and unseen scenarios

Autonomous Driving as Incomplete Information Game

- ARIZONA STATE UNIVERSITY
- Interaction between Human and AV is general-sum dynamic game with incomplete information.
- Reward is a function of states, action and preference parameters
 - Reward $R_p(t) = c_{safety} + \theta_i c_{task}$
- In our previous research, we train at least 4 networks for all agent types (A, NA)
- Training the set up in new task is computationally expensive and time consuming
- Current in literature (Use the same value function for new game); this may not work when tasks are significantly different.
- Example: A strategy trained on aggressive agent may not work for non aggressive agent.

Joint Prob. (ΘxΘ)	Aggressive	Non-aggressive
Aggressive	0.1	0.5
Non-aggressive	0.3	0.1

₹

Human



Autonomous Driving as General Sum Game

- Step back and look into into complete information game
- Looking into Nash Q-learning (temporal difference)
- Reward is a function of states, action and preferences
 - Reward function

$$\boldsymbol{w}_i(s, a, s') = \boldsymbol{\phi}(s, a, s')^\top \mathbf{w}_i$$

- Key Question
 - Can you introduce successor feature in multi agent game?
 - Does **successor feature** aid in better initialization during training?
 - In current AV applications, does the proposed method fare better than the state of the art algorithms?



RIZONA

STATE UNIVERSITY

[1]Successor Feature for Transfer in RL

^{1.} Barreto, A., Dabney, W., Munos, R., Hunt, J.J., Schaul, T., van Hasselt, H.P. and Silver, D., 2017. Successor features for transfer in reinforcement learning. Advances in neural information processing systems, 30.

Transfer in RL? Problem Definition





- Environment is a set of MDPs
- Each MDP M_i is a task
- The only difference between the MDPs are reward function r_i:

$$r_i(s, a, s') = \boldsymbol{\phi}(s, a, s')^\top \mathbf{w}_i$$

- Path chosen differs on the preference eg (IRL)
- (features: coffee, food, distance, leisure)

What is successor feature?





- Exchange of information takes place whenever useful (Generalized Policy Improvement)
 - Already existing knowledge should be transferable
- Transfer is seamlessly integrated with the RL process. (Successor Feature)
 - This requires computing Q as function of weight (preference) and features

^{1.} Barreto, A., Dabney, W., Munos, R., Hunt, J.J., Schaul, T., van Hasselt, H.P. and Silver, D., 2017. Successor features for transfer in reinforcement learning. Advances in neural information processing systems, 30.





Mathematical Formulation

- Changes: two successor table for each agent
- Selection mechanism choice (f)
- Among all mechanism choose reactive policy that maximizes ego's reward
- Update SF table of each agent according to the actions generated from the new policy

```
ARIZONA STATE UNIVERSITY
```

_	earning	
in	put : discount factor γ , selection mechanism	
	f, task index k ,	
	learning rate α , action set b	
OL	itput : Successor features $\psi_i^{\pi_k}$	
1 fo	$\mathbf{r} \ t = 1 \ to \ T \ \mathbf{do}$	
2	for all $i \in N$ do	
3	if <i>Bernoulli</i> $(\epsilon) = 1$ then	
4	$a_i \leftarrow \text{Uniform}(A)$ //exploration	
5	else	
6	$a_i \leftarrow \arg\max_k \max_k \psi^{\pi_k}(s, a_1a_n) w_i$	
	//GPI	
7	end	
8	end	
9	simulate action a_1, a_n in state s	
0	observe rewards r_1, r_n and the next state s'	
1	$w_i = [r_i - \phi_i(s, a, s')^T \tilde{w_i}]\phi_i(s, a, s')$ //Learn	
	weight	
2	for all $i \in N$ do	
3	$Q_i(s', a_1,, a_n) = \max_k \psi_i^{\pi_k}(s, a_1a_n) w_i$	
	//extract $Q_i(s')$	
4	compute $V_i(s') \in$	
	NASH _i $(Q_1^*(s', a_1,, a_n), Q_n^*(s', a_1,, a_n)$	
5	$\pi(s',a_1,a_n)\in$	
	$f(Q_1(s', a_1,, a_n) Q_n(s', a_1,, a_n))$	
6	$V(s')=\pi(s',a_1,a_n)st Q(s,a_1,a_n)$	
7	Compute new action a'_1, a'_n	
8	$\psi_i^{\pi_k}(s,a_1a_n) = (1-lpha)\psi_i^{\pi_k}(s,a_1a_n) +$	
	$\alpha[\phi_i + \gamma \psi_i^{\pi_k}(s', a'_1a'_n)]$ // update successor	
	teature	
9	end	

Experiment Setup

Features

• Goal driven
$$R_p(t) = c_{safety} + \theta_i c_{task}$$

- Property induced by intent parameter
- Rewards (safety) : one step distance from other agent
- Rewards (task) : distance from goal

Simplification

- Single agent setup
- Discrete state [10 20]
- Discrete action [0, 1]
- Change in the intent of ego agent
- Other agent is moving forward

$$R_p(t) = \phi w$$
 $R_p(t) = [c_{safety} \ c_{task}] \bullet \begin{bmatrix} 1 \\ \theta_i \end{bmatrix}$



ARIZONA STATE UNIVERSITY

Preliminary Result: Transfer in Aggressiveness

ARIZONA STATE UNIVERSITY

10

- Single ego agent test
- Testing change of aggressiveness from the generated tabular form
- 50000 iteration with change of aggressiveness every 5000 iterations
- Reduced loss in the system specially for aggressive agents.





Baseline Result

ARIZONA STATE UNIVERSITY

[[16, 10], [10, 16]]	[[16, 10], [10, 16]]
[[16, 10], [10, 15]] Non Aggressive AV	[[16, 10], [10, 15]]
[[16, 10], [10, 14]]	[[16, 10], [10, 14]]
[[16_10] [10_13]] Aggressive H	[[16, 10], [10, 13]]
[[16, 10], [10, 10]]	[[16, 10], [10, 12]]
	[[16, 10], [10, 11]]
	[[16, 10], [10, 10]]
[[16, 10], [10, 10]]	[[16, 10], [10, 9]]
[[16, 10], [10, 9]]	nashSFQl 0
NashQLearning 0	Player 0 follows the policy : stay-stay-stay-stay-stay-stay-stay of length 7
Player 0 follows the policy : stay-stay-stay-stay-stay-stay of length 7	Player 1 follows the policy : up-up-up-up-up-up of length 7
Player 1 follows the policy : up-up-up-up-up-up-up of length 7	game over
	100%
	[[15, 10], [10, 16]]
100%]	[[14, 10], [10, 16]]
100%[]]] 15000/15000 [00:27<00:00, 545.17it/s]	[[13, 10], [10, 16]]
[[15, 10], [10, 16]]	[[12, 10], [10, 15]]
[[15, 10], [10, 15]] NON Aggressive AV	[[11, 10], [10, 15]]
[[15, 10], [10, 15]] Non Aggressive H	[[10, 10], [10, 14]]
NashOLearning 1	[[9, 10], [10, 13]]
Player A follows the policy m-o-d-e-lf-a-i-l-e-dt-oc-o-p-v-e-p-d-e-	nashSFQl 1
Player 4 follows the policy and do 2 follows the policy of	Player 0 follows the policy : up-up-up-up-up of length 6
Player 1 follows the policy : m-o-d-e-lf-a-1-l-e-dt-oc-o-n-v-e-r-g-e-	Player 1 follows the policy : stay-stay-up-stay-up-up of length 6

• Non aggressive ego agent is able to identify other non aggressive agent and provide required policy with proposed algorithm where baseline algorithm fails



Thank you



Sunny Amatya (<u>Samatya@asu.edu</u>) Wenlong Zhang (<u>Wenlong.Zhang@asu.edu</u>) <u>https://home.riselab.info/</u> <u>https://home.riselab.info/nri.html</u>



This material is based on the work supported in part by the National Science Foundation (NSF) under NSF Award Number **925403**. Any opinions, findings and conclusions, or recommendations expressed in this material are those of the author(s) and, do not necessarily reflect those of the NSF.