Welcome to

DS595 Reinforcement Learning Prof. Yanhua Li

Time: 6:00pm –8:50pm W Zoom Lecture Spring 2022

Project 3 leader board

Тор	Date	Name	Score
1	4/6/2022	Hongchao Zhang	128
2	4/6/2022	Puru Upadhyay	82
3	4/6/2022	Khai Yi Chin	81
3	4/6/2022	Karter Krueger	81
5	4/6/2022	Sailesh Rajagopalan	78
5	4/6/2022	Steven Hyland	78
7	4/6/2022	Yiran Fang	74
8	4/6/2022	Zhentian Qian	67
9	4/6/2022	Anujay Sharma	66

Project 4 Progressive report due today

- https://github.com/yingxue-zhang/DS595-RL-Projects/tree/master/Project4
- Progressive report: April 13, 2022, Up to 5 pages, to be submitted to Canvas discussion board.
- Final Project:
 - Monday 4/25/2022 team project report is due
 - Wed 4/27/2022 Virtual Poster Session

- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update

	Reinforcement	Inverse		
	Learning	Reinforcement Learning		
Single Agent	Tabular representation of rewardModel-based controlModel-free control(MC, SARSA, Q-Learning)	Linear reward function learning Imitation learning Apprenticeship learning Inverse reinforcement learning MaxEnt IRL MaxCausalEnt IRL MaxRelEnt IRL		
	Function representation of reward 1. Linear value function approx (MC, SARSA, Q-Learning) 2. Value function approximation			
	<i>(Deep Q-Learning, Double DQN, prioritized DQN, Dueling DQN)</i> <i>3. Policy function approximation (Policy gradient,</i> PPO, TRPO) 4. Actor-Critic methods (A2C, A3C, Pathwise Derivative PG) Advanced topics in RL (Sparse Rewards)	Non-linear reward function learning Generative adversarial imitation learning (GAIL)		
		Adversarial inverse reinforcement learning (AIRL)		
	Review of Deep Learning <i>As bases for non-linear function</i> <i>approximation (used in 2-4).</i>	Review of Generative Adversarial nets As bases for non-linear IRL		
ltiple ents	Multi-Agent Reinforcement Learning Multi-agent Actor-Critic	Multi-Agent Inverse Reinforcement Learning MA-GAII		
Mu Ag	Applicatio	MA-AIRL AMA-GAIL		

- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update



• Actor-Critic (Learned both Value and Policy Functions)

Model-free RL Algorithms

- Value-based (Learned Value Function) Actor-Critic
 - Deep Q-Learning (DQN)
 - Double DQN
 - Dueling DQN
 - Prioritized DQN Multi-step DQN,
 - Noisy net DQN
 - Distributional DQN
 - DQN for continuous action space
- Policy-based (Learned Policy Function)
 - Basic Policy Gradient
 Algorithm
 - REINFORCE
 - Vanilla, PPO, TRPO, PPO2

A2C (Learned both Value and Policy Functions)
 A3C

Pathwise Derivative
 Policy Gradient



Sparse Reward

Reward Shaping

Reward Shaping



Reward Shaping

VizDoom

https://openreview.net/forum?id=Hk3 mPK5gg¬eId=Hk3mPK5gg

Parameters	Description	FlatMap	CIGTrack1
living	Penalize agent who just lives	-0.008 / action	
health_loss	Penalize health decrement	-0.05 / unit	
ammo_loss	Penalize ammunition decrement	-0.04 / unit	
health_pickup	Reward for medkit pickup	0.04 / unit	
ammo_pickup	Reward for ammunition pickup	0.1	5 / unit
dist_penalty	Penalize the agent when it stays	-0.03	/ action
dist_reward	Reward the agent when it moves	9e-5 / ui	nit distance



Reward Shaping





Get reward, when closer Need domain knowledge

https://openreview.net/pdf?id=Hk3mPK5gg

Curiosity



Intrinsic Curiosity Module







Intrinsic Curiosity Module



- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update

Sparse Reward

Curriculum Learning

Curriculum Learning

Starting from simple training examples, and then becoming harder and harder.

Facebook wins the VizDoom competition

	Class 0	Class 1	Class 2	Class 3	Class 4	Class 5	Class 6	Class 7
Speed	0.2	0.2	0.4	0.4	0.6	0.8	0.8	1.0
Health	40	40	40	60	60	60	80	100

Reverse Curriculum Generation



- \succ Given a goal state s_g .
- > Sample some states s_1 "close" to s_g
- > Start from states s_1 , each trajectory has reward $R(s_1)$

Reverse Curriculum Generation



- Delete s₁ whose reward is too large (already learned) or too small (too difficult at this moment)
- > Sample s_2 from s_1 , start from s_2

- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update

Sparse Reward

Hierarchical Reinforcement Learning

Hierarchical RL

This is a fake example. Don't take it seriously.



If lower agent cannot achieve the goal, the upper agent would get penalty.

https://arxiv.org/abs/1805.08180





https://arxiv.org/abs/1805.08180

- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update

Multi-Agent Systems are going to be everywhere





Autonomous Vehicles





New advances in Multi-Agent RL research - Hide-and-Seek



https://openai.com/blog/emergent-tool-use/ for more details. Just released on 9/17/2019 by OpenAI.

New advances in Multi-Agent RL research - DOTA



https://openai.com/blog/openai-five/ for more details. Released by OpenAI on 6/25/2018.

Multi-Agent Reinforcement Learning: Naïve Methods

Each agent learns independently, by viewing all other agents and the system as the environment.





Multi-Agent Reinforcement Learning: Naïve Methods











Each individual agent runs an RL algorithm, separately: Policy Gradient DQN Actor-Critic

Issues with the Naïve Methods

- Single agent
 - Model: Markov Decision Process (MDPs)
 - Goal: Maximize reward
 - Steady state: Maximized reward
 - Agents dependency: Independent

- Multiple agents
 - Model: Markov Games (MGs)
 - Goal: Maximize its own/or team reward
 - Steady state: Nash Equilibrium
 - Agents dependency: no agents can achieve a higher expected reward by unilaterally changing its own policy

Dependency across agents are missing!

- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update

MDPs vs MGs

- MDP: <S, A, P, γ, R>
 - State set: S
 - Action set: A
 - Transition Probability: P(s'|s, a)
 - Discount factor γ
 - Reward function: R R: S × A $\mapsto \mathbb{R}$

MG: <S, {A_i}, P, γ, {R_i}>

- State set: S (joint states)
- Action set: {A₁,...,A_N}
- Transition Probability:
- $P(s'|s, a_1,...,a_N)$
- Discount factor γ
- Reward function: $\{R_1, \dots, R_N\}$ $R_i : S \times A_i \longrightarrow \mathbb{R}$



 $\mathbf{Q(s,a)} \text{ net:} \nabla \bar{R}_{\theta} = \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{\tau_n} (r^n_t + \max_{a^n_{t+1}} Q(s^n_{t+1}, a^n_{t+1}) - \max_{a^n_t} Q(s^n_t, a^n_t)) \nabla_{\theta} \log \pi(a^n_t | s^n_t)$

Multi-Agent RL (Principle)



Figure Source: https://medium.com/brillio-data-science/improving-openai-multi-agent-actor-critic-rlalgorithm-27719f3cafd4



Updating each θ_i by policy gradient

$$\nabla \bar{R}_{\theta_{i}} = \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{\tau_{n}} \left(r^{n}_{t} + \max_{a^{n}_{i,t+1}} Q_{i}(s^{n}_{t+1}, a^{n}_{t+1}) - \max_{a^{n}_{i,t}} Q_{i}(s^{n}_{t}, a^{n}_{t}) \right) \nabla_{\theta} \log \pi_{i}(a^{n}_{i,t}|s^{n}_{t})$$

- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update



Updating θ by policy gradient

 $\nabla \bar{R}_{\theta} = \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{\tau_n} (r^n_t + \max_{a^n_{t+1}} Q(s^n_{t+1}, a^n_{t+1}) - \max_{a^n_t} Q(s^n_t, a^n_t)) \nabla_{\theta} \log \pi(a^n_t | s^n_t)$ with $\bar{R}_{\theta} = \mathbb{E}_{s \sim P, a \sim \pi} [Q(s, a) - V(s)] = \mathbb{E}_{s \sim P, a \sim \pi} [A(s, a)].$

 $\nabla \bar{R}_{\theta} = \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{\tau_n} \nabla_{a_t}^n Q(s_t^n, a_t^n) \nabla_{\theta} \mu_{\theta}(s_t^n), \text{ with } \bar{R}_{\theta} = \mathbb{E}_{s \sim P, a \sim \mu}[Q(s, a)]$ https://arxiv.org/pdf/1509.02971.pdf



Updating each θ_i by MA-DDPG

$$\nabla \bar{R}_{\boldsymbol{\theta}_{i}} = \frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{\tau_{n}} \nabla_{a^{n}_{i,t}} \boldsymbol{Q}_{i} (s^{n}_{t}, \boldsymbol{a}^{n}_{i,t}) \nabla_{\boldsymbol{\theta}_{i}} \mu_{i}(s^{n}_{t})$$

- Advanced RL Techniques
 - Sparse Reward
 - Reward shaping, Curiosity module
 - Curriculum learning,
 - Hierarchical RL
- Multi-Agent RL
 - Intro
 - Multi-Agent RL principle
 - Multi-Agent Deep Deterministic Policy Gradient (MA-DDPG)
- Project #4 progress update

	Reinforcement	Inverse			
	Learning	Reinforcement Learning			
Single Agent	Tabular representation of rewardModel-based controlModel-free control(MC, SARSA, Q-Learning)	Linear reward function learning Imitation learning Apprenticeship learning Inverse reinforcement learning			
	Function representation of reward1. Linear value function approx(MC, SARSA, Q-Learning)2. Value function approximation	MaxEnt IRL MaxCausalEnt IRL MaxRelEnt IRL			
	(Deep Q-Learning, Double DQN, prioritized DQN, Dueling DQN) 3. Policy function approximation (Policy gradient, PPO, TRPO) 4. Actor-Critic methods (A2C, A3C, Pathwise Derivative PG) Advanced topics in RL (Sparse Rewards)	Non-linear reward function learning Generative adversarial imitation learning (GAIL)			
		Adversarial inverse reinforcement learning (AIRL)			
	Review of Deep Learning <i>As bases for non-linear function</i> <i>approximation (used in 2-4).</i>	Review of Generative Adversarial nets As bases for non-linear IRL			
lltiple jents	Multi-Agent Reinforcement Learning Multi-agent Actor-Critic etc.	Multi-Agent Inverse Reinforcement Learning MA-GAIL			
Mu Ag	Applicatio	MA-AIRL AMA-GAIL			

Questions?