TDDE13 LE3 HT2023 Multi-agent Learning

Fredrik Heintz

Dept. of Computer Science Linköping University

fredrik.heintz@liu.se @FredrikHeintz **Outline:**

- Reinforcement learning
- Deep reinforcement learning
- Multi-agent reinforcement learning



Classes of Learning Problems Supervised Learning Unsupervised Learning

Data: (x, y) x is data, y is label

Goal: Learn function to map $x \rightarrow y$

Data: *x x* is data, no labels!

Goal: Learn underlying structure

Goal: Maximize future rewards over many time steps

Reinforcement Learning

Data: state-action pairs

Apple example:

This thing is an apple.

Apple example:



Apple example:

Eat this thing because it will keep you alive.



From Supervised to Reinforcement Learning -Learning How to Act



Humorous reminder from IEEE Spectrum: The DARPA 2015 Humanoid Challenge "Fail Compilation" To be fair, this is the state of the art: <u>https://youtu.be/NR32ULxbjYc</u>

- Can we use supervised learning to **learn** how to act?
- E.g. engineering robot behavior can be fragile and time consuming
 - Things humans do without thinking require **extremely detailed instructions** for a robot. Even robust locomotion is hard.



Learning How to Act

- Yes, one can learn a mapping from problem state (e.g. position) to action
 - As in all supervised learning, this requires a teacher
 - Sometimes called "imitation learning"
- However, **supervised learning** with robots can get tedious as providing examples of correct behaviour is difficult to automate
- Can we remove the human from the loop?
 - 1. An **automated teacher** like a **planning or optimal control** algorithm can generate supervised examples **if it has a model of the environment**
 - Mordatch et al, <u>https://www.youtube.com/watch?v=IxrnToJOs40</u>
 - LiU's research with real nano-quadcopters (deep ANN on-board the microcontroller)
 - 2. Reinforcement learning attempts to generalize this to learning from scratch in completely **unknown environments**



Reinforcement Learning Basic Concept

• Reinforcement Learning is learning what to do – how to map situations to actions – so as to maximum a numerical reward.

Reinforcement Learning: An introduction Sutton & Barto

- Rather than learning from explicit training data, or discovering patterns in static data, reinforcement learning discovers the best option (highest reward) from trial and error.
- Inverse Reinforcement Learning
 - Learn reward function by observing an expert
 - "Apprenticeship learning"
 - E.g. Abbeel et al. *Autonomous Helicopter Aerobatics through Apprenticeship Learning*







2023-11-24



Agent: takes actions.





2023-11-24

Environment: the world in which the agent exists and operates.





Action: a move the agent can make in the environment.

Action space A: the set of possible actions an agent can make in the environment





Observations: of the environment after taking actions.





State: a situation which the agent perceives.





Reward: feedback that measures the success or failure of the agent's action.



















A Reinforcement Learning Problem

- The environment
- The reinforcement function *r*(*s*,*a*)
 - Pure delay reward and avoidance problems
 - Minimum time to goal
 - Games
- The value function *V*(*s*)
 - Policy $\pi: S \to A$
 - Value $V^{\pi}(s) := \Sigma_i \gamma^i r_{t+i}$
- Find the optimal policy π* that maximizes V^{π*}(s) for all states s.





Goal: Learn to choose actions that maximize $r_0 + \gamma r_1 + \gamma^2 r_2 + ...$, where $0 < \gamma < 1$



RL Value Function - Example

A minimum time to goal world





Markov Decision Processes

Assume:

- finite set of states *S*, finite set of actions *A*
- at each discrete time agent observes state $s_t \in S$ and chooses action $a_t \in A$
- then receives immediate reward r_t
- and state changes to s_{t+1}
- Markov assumption: $s_{t+1} = \delta(s_t, a_t)$ and $r_t = r(s_t, a_t)$
 - i.e. r_t and s_{t+1} depend only on current state and action
 - functions δ and r may be non-deterministic
 - functions δ and r not necessarily known to the agent



MDP Example





Defining the Q-Function

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots$$

Total reward, R_t , is the discounted sum of all rewards obtained from time t

$$Q(\mathbf{s}_t, \mathbf{a}_t) = \mathbb{E}[R_t | s_t, \mathbf{a}_t]$$

The Q-function captures the **expected total future reward** an agent in state, *s*, can receive by executing a certain action, *a*



How to Take Actions Given a Q-Function $Q(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t]$ (state, action)

Ultimately, the agent needs a **policy** $\pi(s)$, to infer the **best action to take** at its state, s

Strategy: the policy should choose an action that maximizes future reward

$$\pi^*(s) = \operatorname{argmax}_{a} Q(s, a)$$



The Q-Function

Optimal policy:

- $\pi^*(s) = \operatorname{argmax}_a[r(s,a) + \gamma V^*(\delta(s,a))]$
- Doesn't work if we don't know r and δ .

The Q-function:

- $Q(s,a) := r(s,a) + \gamma V^*(\delta(s,a))$
- $\pi^*(s) = \operatorname{argmax}_a Q(s,a)$





Q(s,a)



The Q-Function

- Note Q and V* closely related: $V^*(s) = \max_{a'}Q(s,a')$
- Therefore Q can be written as: $Q(s_t, a_t) := r(s_t, a_t) + \gamma V^*(\delta(s_t, a_t)) = r(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, a')$
- If Q^{\wedge} denote the current approximation of Q then it can be updated by: $Q^{\wedge}(s,a) := r + \gamma \max_{a'} Q^{\wedge}(s',a')$



- Value-Based:
 - Learn value function
 - Implicit policy (e.g. greedy selection)
 - Example: Deep Q Networks (DQN)
- Policy-Based:
 - No value function
 - Learn explicit (stochastic) policy
 - Example: Stochastic Policy Gradients
- Actor-Critic:
 - Learn value function
 - Learn policy using value function
 - Example: Asynchronous Advantage Actor Critic (A3C)





Reinforcement Learning Algorithms

Value Learning

Find Q(s, a) $a = \underset{a}{\operatorname{argmax}} Q(s, a)$

Policy Learning

Find $\pi(s)$ Sample $a \sim \pi(s)$



Reinforcement Learning Algorithms

Value Learning

Find Q(s, a) $a = \underset{a}{\operatorname{argmax}} Q(s, a)$

Policy Learning

Find $\pi(s)$ Sample $a \sim \pi(s)$



2023-11-24 27

Q-Learning for Deterministic Worlds

For each *s*, *a* initialize table entry $Q^{(s,a)} := 0$.

Observe current state *s*.

Do forever:

- 1. Select an action *a* and execute it
- 2. Receive immediate reward r
- 3. Observe the new state *s*'
- 4. Update the table entry for $Q^{(s,a)}$: $Q^{(s,a)} := r + \gamma \max_{a'} Q^{(s',a')}$

5. s := s'



Q-Learning Example





Q-Learning Continued

- Exploration
 - Selecting the best action
 - Probabilistic choice
- Improving convergence
 - Update sequences
 - Remember old state-action transitions and their immediate reward
- Non-deterministic MDPs
- Temporal Difference Learning



2023-11-24 30

Deep Q-Learning (DQN)





How can we use deep neural networks to model Q-functions?



What happens if we take all the best actions? Maximize target return \rightarrow train the agent















Deep Q Network Summary

Use NN to learn Q-function and then use to infer the optimal policy, $\pi(s)$



Send action back to environment and receive next state



DQN Atari Results




DQN Atari Results





Downsides of Q-Learning

Complexity:

- · Can model scenarios where the action space is discrete and small
- · Cannot handle continuous action spaces

Flexibility:

 Policy is deterministically computed from the Q function by maximizing the reward → cannot learn stochastic policies

To address these, consider a new class of RL training algorithms: Policy gradient methods



Reinforcement Learning Algorithms





2023-11-24 40

Deep Q Networks

DQN: Approximate Q-function and use to infer the optimal policy, $\pi(s)$





Policy Gradient (PG): Key Idea

DQN: Approximate Q-function and use to infer the optimal policy, $\pi(s)$

Policy Gradient: Directly optimize the policy $\pi(s)$





Discrete vs Continuous Action Spaces

Discrete action space: which direction should I move?





Discrete vs Continuous Action Spaces

Discrete action space: which direction should I move?





Policy Gradient (PG): Key Idea

Policy Gradient: Enables modeling of continuous action space





Training Policy Gradients: Case Study

Reinforcement Learning Loop:

OBSERVATIONSState changes: S_{t+1}
Reward: r_t Reward: r_t Action: a_t Action: a_t ACTIONS

Agent:vehicleState:camera, lidar, etcAction:steering wheel angleReward:distance traveled



Case Study – Self-Driving Cars

- I. Initialize the agent
- 2. Run a policy until termination
- 3. Record all states, actions, rewards
- 4. Decrease probability of actions that resulted in low reward
- 5. Increase probability of actions that resulted in high reward





- I. Initialize the agent
- 2. Run a policy until termination
- 3. Record all states, actions, rewards
- 4. Decrease probability of actions that resulted in low reward
- 5. Increase probability of actions that resulted in high reward





- I. Initialize the agent
- 2. Run a policy until termination
- 3. Record all states, actions, rewards
- 4. Decrease probability of actions that resulted in low reward
- 5. Increase probability of actions that resulted in high reward





- I. Initialize the agent
- 2. Run a policy until termination
- 3. Record all states, actions, rewards
- 4. Decrease probability of actions that resulted in low reward
- 5. Increase probability of actions that resulted in high reward





Training Algorithm

- I. Initialize the agent
- 2. Run a policy until termination
- 3. Record all states, actions, rewards
- 4. Decrease probability of actions that resulted in low reward
- 5. Increase probability of actions that resulted in high reward

log-likelihood of action

$$\mathbf{loss} = -\log \mathbf{P}(a_t | s_t) \mathbf{R}_t$$

reward

Gradient descent update:

$$w' = w - \nabla \mathbf{loss}$$

$$w' = w + \nabla \log P(a_t | s_t) R_t$$

Policy gradient!



Actor-Critic Methods

- The main components in policy gradient are the policy model and the value function.
- Actor-critic methods learn the value function in addition to the policy.
 - Knowing the value function assists the policy update, such as by reducing gradient variance in vanilla policy gradients.

Actor-critic methods consist of two models, which may optionally share parameters:

- Critic updates the value function parameters w and depending on the algorithm it could be actionvalue $Q_w(a|s)$ or state-value $V_w(s)$.
- Actor updates the policy parameters heta for $\pi_{ heta}(a|s)$, in the direction suggested by the critic.



Actor-Critic Methods

- 1. Initialize s, heta, w at random; sample $a \sim \pi_{ heta}(a|s)$.
- 2. For t = 1 ... T:
 - 1. Sample reward $r_t \sim R(s,a)$ and next state $s' \sim P(s'|s,a)$;
 - 2. Then sample the next action $a' \sim \pi_{ heta}(a'|s')$;
 - 3. Update the policy parameters: $heta \leftarrow heta + lpha_ heta Q_w(s,a)
 abla_ heta \ln \pi_ heta(a|s)$;
 - 4. Compute the correction (TD error) for action-value at time t:

$$\delta_t = r_t + \gamma Q_w(s',a') - Q_w(s,a)$$

and use it to update the parameters of action-value function:

- $w \leftarrow w + lpha_w \delta_t
 abla_w Q_w(s,a)$
- 5. Update $a \leftarrow a'$ and $s \leftarrow s'$.

Two learning rates, α_{θ} and α_{w} , are predefined for policy and value function parameter updates respectively.



Deep Deterministic Policy Gradient (DDPG)

- A model-free off-policy actorcritic algorithm, combining DPG and DQN.
- DQN stabilizes the learning of Q-function by experience replay and the frozen target network.
- DQN works in discrete space, and DDPG extends it to continuous space with the actor-critic framework while learning a deterministic policy.

Algorithm 1 DDPG algorithm

Randomly initialize critic network $Q(s, a | \theta^Q)$ and actor $\mu(s | \theta^\mu)$ with weights θ^Q and θ^μ . Initialize target network Q' and μ' with weights $\theta^{Q'} \leftarrow \theta^Q$, $\theta^{\mu'} \leftarrow \theta^\mu$ Initialize replay buffer Rfor episode = 1, M do Initialize a random process \mathcal{N} for action exploration Receive initial observation state s_1 for t = 1, T do Select action $a_t = \mu(s_t | \theta^\mu) + \mathcal{N}_t$ according to the current policy and exploration noise Execute action a_t and observe reward r_t and observe new state s_{t+1} Store transition (s_t, a_t, r_t, s_{t+1}) in RSample a random minibatch of N transitions (s_i, a_i, r_i, s_{i+1}) from RSet $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta^{\mu'}) | \theta^{Q'})$ Update critic by minimizing the loss: $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2$ Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^{\mu}} J \approx \frac{1}{N} \sum_{i} \nabla_a Q(s, a | \theta^Q) |_{s=s_i, a=\mu(s_i)} \nabla_{\theta^{\mu}} \mu(s | \theta^{\mu}) |_{s_i}$$

Update the target networks:

 $\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$ $\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$

end for end for



Reinforcement Learning - Neural Networks as Function Approximators

- To tackle a high-dimensional state space or continous states we can use a neural network as function approximator
- Lunar Lander experiment
 - 8 continous/discrete states
 - XY-Pos, XY-Vel, Rot, Rot-rate, Leg1/Leg2 ground contact
 - 4 discrete actions
 - Left thrust
 - Right thrust
 - Main engine thrust
 - NOP
 - Rewards
 - Move from top to bottom of the screen (+ ~100-140)
 - Land between the posts (+100)
 - Put legs on ground (+10 per leg)
 - Penalties
 - Using main engine thrust (-0.3 per frame)
 - Crashing (-100)
- Solved using Stochastic Policy Gradients







Reinforcement Learning Neural Networks as Function Approximators





AlphaGo Beats Top Human Player (2016)





MuZero: Learning Dynamics for Planning (2020)



AlphaGo becomes the first program to master Go using neural networks and tree search (Jan 2016, Nature)



Go

AlphaGo Zero learns to play completely on its own, without human knowledge (Oct 2017, Nature)



Knowledge





MuZero



Domains

Known rules

Knowledge

AlphaZero masters three perfect information games using a single algorithm for all games (Dec 2018, Science)



MuZero learns the rules of the game, allowing it to also master environments with unknown dynamics. (Dec 2020, Nature)







Deep Reinforcement Learning Summary

Foundations

- Agents acting in environment
- State-action pairs → maximize future rewards
- Discounting



Q-Learning

- Q function: expected total reward given **s**, **a**
- Policy determined by selecting action that maximizes Q function



Policy Gradients

- Learn and optimize the policy directly
- Applicable to continuous action spaces





Reinforcement Learning Concepts

- Value-Based:
 - Learn value function
 - Implicit policy (e.g. greedy selection)
 - Example: Deep Q Networks (DQN)
- Policy-Based:
 - No value function
 - Learn explicit (stochastic) policy
 - Example: Stochastic Policy Gradients
- Actor-Critic:
 - Learn value function
 - Learn policy using value function
 - Example: Asynchronous Advantage Actor Critic (A3C)





Multi-Agent Reinforcement Learning



Motivation

- In multi-agent systems it may be beneficial to learn
 - Environment dynamics
 - Reward functions
 - Other agents' strategies
- Can be useful for cooperative as well as competitive scenarios
- Finding solutions by hand may be difficult and time-consuming
 - Learning may help
 - Also lets the agent adapt to changes in the environment



Challenges

- Non-stationary environment when multiple agents learn
 - Simultaneous learning and teaching for individual agent
- Lack of observability (e.g., actions and rewards of other agents)



Model-based and Model-free Learning

- Model-based algorithms
 - Try to model the behavior of other agents
 - More specific models may improve performance, but at a loss of generality
- Model-free algorithms
 - Do not use models
 - May perform worse or take longer to converge



Markov Models

	No Agents	Single Agent	Multiple Agents
State Known	Markov Chain	Markov Decision Process (MDP)	Markov Game (a.k.a. Stochastic Game)
State Observed Indirectly	Hidden Markov Model (HMM)	Partially-Observable Markov Decision Process (POMDP)	Partially-Observable Stochastic Game (POSG)



MDP, POMDP, and Dec-POMDP



Figure: (a) Markov decision process (MDP) (b) Partially observable Markov decision process (POMDP) (c) Decentralized partially observable Markov decision process with two agents (Dec-POMDP)



Multi-Agent Reinforcement Learning





Multi-Agent Reinforcement Learning





Decentralized Multi-Agent Deep Reinforcement Learning

- Though no theoretical guarantees exist, single-agent algorithms may produce interesting results in multi-agent systems
- Ways to stabilize the learning process
 - Clever design of reward systems
 - Training populations of agents
 - Can allow agents to generalize



Centralized Learning, Decentralized Execution

- Extra information is used for guidance during learning, e.g., actor-critic setup or value function decomposition
- At execution time agents act based on local observations
- Examples of algorithms
 - QMIX
 - COMA
 - MADDPG





2023-11-24

69



Research in Multi-Agent RL

Large Problems	Approximate Solution Methods	Approximate Solution Methods
Small Problems	Tabular Solution Methods	Tabular Solution Methods
	Single Agent	Multiple (e.g. 2) Agents



Axes of Multi-Agent RL

- Centralized:
 - One brain/algorithm deployed across many agents
- Prescriptive:
 - Suggests how agents should behave
- Cooperative: Agents cooperate to achieve a goal
 - Shared team reward
- Numbers of agents
 - One (single-agent)
 - Two (very common)
 - Finite
 - Infinite

- Decentralized:
 - All agents learn individually
 - Communication limitations defined by environment
- Descriptive:
 - Forecast how agent will behave
- Competitive: Agents compete against each other
 - Zero-sum games
 - Individual opposing rewards
- Neither: Agents maximize their utility which may require cooperating and/or competing
 - General-sum games



Foundations of MARL




First MARL Algorithm: Minimax-Q (Littman '94)

Q-values are over joint actions: Q(s, a, o)

■ s = state

- a = your action
- o = action of the opponent

Instead of playing action with highest Q(s, a, o), play MaxMin

$$Q(s, a, o) = (1 - \alpha)Q(s, a, o) + \alpha(r + \gamma V(s'))$$
$$V(s) = \max_{\pi_s} \min_o \sum_a Q(s, a, o)\pi_s(a)$$



MARL Formulation

- The agents choose actions according to their policies.
- For agent *j*, the corresponding policy is defined as $\pi^j : S \to \Omega(A^j)$, where $\Omega(A^j)$ is the collection of probability distributions over agent *j*'s action space A^j .
- Let $\boldsymbol{\pi} = [\pi^1, \cdots, \pi^N]$ is the joint policy of all agents, then

$$v^{j}_{\pi}(s) = v^{j}(s;\pi) = \sum_{t=0}^{\infty} \gamma^{t} \mathbb{E}_{\pi,p}[r^{j}_{t}|s_{0}=s,\pi]$$

• Q-function such that the Q-function $Q_{\pi}^{j}: S \times A^{1} \times \cdots \times A^{N} \to \mathbb{R}$ of agent j under the joint policy π :

$$Q^j_{\pi}(s, \boldsymbol{a}) = r^j(s, \boldsymbol{a}) + \gamma \mathbb{E}_{s' \sim p}[v^j_{\pi}(s')]$$



Nash Q-Learning

- In MARL, the objective of each agent is to learn an optimal policy to maximize its value function
- Optimizing the v_{π}^{j} for agent j depends on the joint policy π of all agents
- A Nash equilibrium is a joint policy π such that no player has incentive to deviate unilaterally. It is represented by a particular joint policy

$$\boldsymbol{\pi_*} = [\pi^1_*, \cdots, \pi^N_*]$$

such that for all $s \in S, j \in \{1, \cdots, N\}$ it satisfies:

$$v^j(s; \pi_*) = v^j(s; \pi^j_*, \pi^{-j}_*) \geq v^j(s; \pi^j, \pi^{-j}_*)$$

Here π_*^{-j} is the joint policy of all agents except j as

$$\boldsymbol{\pi}_{*}^{-j} = [\pi_{*}^{1}, \cdots, \pi_{*}^{j-1}, \pi_{*}^{j+1}, \cdots, \pi_{*}^{N}]$$



Nash Q-Learning (cont)

- In a Nash equilibrium, each agent acts with the best response π_*^j to others, provided that all other agents follow the policy π_*^{-j}
- For a N-agent stochastic game, there is at least one Nash equilibrium with stationary policies, assuming players are rational
- \blacksquare Given Nash policy π_* , the Nash value function

$$oldsymbol{v}^{\mathsf{Nash}} = [v_{\pi_*}^1(s), \cdots, v_{\pi_*}^N(s)]$$

 $oldsymbol{Q}(s, oldsymbol{a})^{\mathsf{Nash}} = \mathbb{E}_{s' \sim p}[oldsymbol{r}(s, oldsymbol{a}) + \gamma oldsymbol{v}^{\mathsf{Nash}}(s')]$

where $\boldsymbol{r}(s, \boldsymbol{a}) = [r^1(s, \boldsymbol{a}), \cdots, r^N(s, \boldsymbol{a})]$



MARL Policies



(a) Single-agent



(b) Multiple logical entities, single "super-agent"





MARL Policies





- *n* pursuit-evasion a set of agents (the pursuers) are attempting to chase another set of agents (the evaders)
- The agents in the problem are self-interested (or heterogeneous), i.e. they have different objectives
- The two pursuers are attempting to catch the two evaders





Challenge: defining the problem in such a way that an arbitrary number of agents can be represented without changing the architecture of the deep Q-Network.

Solution (under some assumptions):

- The image tensor is of size 4× W × H, where W and H are the height and width of our two dimensional domain and four is the number of channels in the image.
- Channels:
 - **Background Channel**: contains information about any obstacles in the environment
 - Opponent Channel: contains information about all the opponents
 - Ally Channel: contains information about all the allies
 - **Self Channel:** contains information about the agent making the decision







- Train one agent at a time, and fix policies of all the other agents
- After a number of iterations distribute the policy learned by the training agent to all the other agents of its type







MADQN: Dealing with Ambiguity

- **Challenge**: When two ally agents are occupying the same position in the environment, the image-like state representation for each agent will be identical, so their policies will be exactly the same.
- **Solution**: To break this symmetry use a stochastic policy for agents. The actions taken by the agent are drawn from a distribution derived by taking a softmax over the Q-values of the neural network. This allows allies to take different actions if they occupy the same state and break the ambiguity.



Multi-Agent Deep Deterministic Policy Gradient

- **Multi-agent DDPG** (**MADDPG**) (Lowe et al., 2017) extends DDPG to an environment where multiple agents are coordinating to complete tasks with only local information.
- In the viewpoint of one agent, the environment is nonstationary as policies of other agents are quickly upgraded and remain unknown.
- MADDPG is an actor-critic model redesigned particularly for handling such a changing environment and interactions between agents.
 - Centralized critic + decentralized actors, changes non-stationary problems to stationary problems
 - Actors can use estimated policies of other agents for learning
 - Policy ensembles is good for reducing variance





Algorithm 1: Multi-Agent Deep Deterministic Policy Gradient for N agents

for episode = 1 to M do Initialize a random process \mathcal{N} for action exploration Receive initial state x for t = 1 to max-episode-length **do** for each agent *i*, select action $a_i = \mu_{\theta_i}(o_i) + \mathcal{N}_t$ w.r.t. the current policy and exploration Execute actions $a = (a_1, \ldots, a_N)$ and observe reward r and new state \mathbf{x}' Store $(\mathbf{x}, a, r, \mathbf{x}')$ in replay buffer \mathcal{D} $\mathbf{x} \leftarrow \mathbf{x}'$ for agent i = 1 to N do Sample a random minibatch of S samples $(\mathbf{x}^j, a^j, r^j, \mathbf{x}'^j)$ from \mathcal{D} Set $y^{j} = r_{i}^{j} + \gamma Q_{i}^{\mu'}(\mathbf{x}^{\prime j}, a_{1}^{\prime}, \dots, a_{N}^{\prime})|_{a_{L}^{\prime} = \mu_{L}^{\prime}(o_{L}^{j})}$ Update critic by minimizing the loss $\mathcal{L}(\theta_i) = \frac{1}{S} \sum_j \left(y^j - Q_i^{\mu}(\mathbf{x}^j, a_1^j, \dots, a_N^j) \right)^2$ Update actor using the sampled policy gradient: $\nabla_{\theta_i} J \approx \frac{1}{S} \sum_i \nabla_{\theta_i} \boldsymbol{\mu}_i(o_i^j) \nabla_{a_i} Q_i^{\boldsymbol{\mu}}(\mathbf{x}^j, a_1^j, \dots, a_i, \dots, a_N^j) \big|_{a_i = \boldsymbol{\mu}_i(o_i^j)}$

end for

Update target network parameters for each agent i:

$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau) \theta_i'$$

end for end for



Challenges in Multi-Agent Learning

- Computational complexity
 - AlphaGo Zero (per agent):
 - 64 GPUs & 19 CPUs
 - OpenAI Dota Five
 - 256 GPUs & 128000 CPUs
- Lack of good benchmarks
- Reproducability



AlexNet to AlphaGo Zero: A 300,000x Increase in Compute



TDDE13 MAS LE3 HT2023: Reinforcement learning Deep reinforcement learning Multi-agent learning

www.ida.liu.se/~TDDE13

