Introduction of Reinforcement Learning with Related Applications

Minsung Yoon Ph.D. student at SGVR Lab.





Korea Advanced Institute of Science and Technology

School of Computing

Contents

Overview of Reinforcement Learning (RL) Technology

- Applications for Robot Arms
 - Learning-based Initialization of Trajectory Optimization for Path-following Problems of Redundant Manipulators, ICRA 2023
 - Towards Safe Remote Manipulation: User Command Adjustment based on Risk Prediction for Dynamic Obstacles, ICRA 2023
- Applications for Quadruped Robots
 - Enhancing Navigation Efficiency of Quadruped Robots via Leveraging Personal Transportation Platforms, ICRA 2025
 - Learning-based Adaptive Control of Quadruped Robots for Active Stabilization on Moving Platforms, IROS 2024





Branches of Machine Learning







Agents in Reinforcement Learning







Language model (software)



Reinforcement Learning

 Reinforcement learning is a type of machine learning where agents learn <u>optimal</u> <u>behaviors</u> through <u>trial and error interactions with an environment</u>.









Agent and Environment

- At each time step t, the agent:
 - **Executes** action A_t
 - **Receives** observation O_t (sensing)
 - **Receives** scalar reward R_t
- The environment:
 - **Receives** action A_t
 - Emits observation O_{t+1}
 - **Emits** scalar reward R_{t+1}







Markov Decision Process (MDP)

- MDP is a formal mathematical modeling of RL problem, defined by a 5-tuple.
 - **States (S):** All information describing the agent's situation.
 - Actions (A): Available choices the agent can make in each state.
 - **Rewards (R):** Immediate (positive or negative) feedback for an action taken in a state.
 - **Transitions (T):** Probabilities of moving between states after actions.
 - **Discount Factor (γ):** Reduces the value of future rewards compared to immediate ones.





Optimization objective of RL

• The goal of RL is to find an optimal "policy (agent) π^* " which produces the action maximizing the discounted expected future rewards at each state.

$$\pi^* = \operatorname{argmax}_{\pi} J(\pi)$$

$$J(\pi) = E_{a \sim \pi(s), s' \sim T(s, a)} \left[\sum_{t=0}^{\infty} \gamma^{t} R(s, a, s') \right]$$



Arrows are induced optimal actions at each state (-1 rewards at each step, except when reaching the goal state)



Partially Observable Markov Decision Process (POMDP)

- POMDP is an extension of MDP, where an agent cannot fully observe the entire state of the environment.
 - States (S):
 - All environmental states.
 - Observations (O):
 - Partial information the agent receives (e.g., RGB-D images).
 - Actions (A)
 - Rewards (R)
 - Transitions (T)
 - Discount Factor (γ)



To convert POMDP into MDP,
estimate hidden states (i.e., belief states of HMM) from a history of observations.





Taxonomy of RL algorithms



Taxonomy of RL methods

- **Model-free RL**: agent learns optimal policies directly from interactions with the environment.
 - Simplicity, Flexibility to changes in the environment
 - Poor sample efficiency, Responding to unexperienced states, Lack of planning
- Model-based RL (Planning): utilizes an environment model where the agent lives (the model can be given naturally, or it can be learned from data)
 - Sample efficient learning, Planning capability
 - **8** Risk of model accuracy and bias



10









Characteristics of Reinforcement Learning

- What makes RL different from other machine learning paradigms?
 - No supervisor, only a reward signal
 - Feedback is delayed, not instantaneous
 - Time matters
 - agent's actions affect the subsequent data it will receive
 - collected data is sequential \rightarrow not i.i.d. (independent and identically distributed)





Application of RL for Robot Arm Tasks (1)





Learning-based Initialization of Trajectory Optimization for Path-following Problems of Redundant Manipulators

Min-Sung Yoon, Min-Cheul Kang, Dae-Hyung Park, Sung-Eui Yoon





Korea Advanced Institute of Science and Technology

School of Computing

Problem Statement of Path-following Problems

 Generate a joint trajectory precisely following a given 6-dimensional Cartesian path (i.e., target path) with an end-effector.



13

KAIST



Relationship with IK Problem of Path-following Problems



- End-effector matching
- Collision avoidance

...





Redundant Manipulators

 Redundant manipulators have an infinite number of Inverse kinematics (IK) solutions given a single endeffector pose.





[Inverse kinematics (IK) solutions given a pose] (Fetch robot manipulator has **7 degree of freedom**)





Approaches for Path-following Problems



Optimization-based approaches

[Kang et al. 2020, Schulman et al. 2013, Zucker et al. 2013, ...]

- Variable: a whole joint trajectory, $\xi = \{q_0, \dots, q_{N-1}\}$
- Objective function:

$$\mathcal{U}(\xi) = \mathcal{F}_{pose}(\xi) + \lambda_1 \mathcal{F}_{obs}(\xi) + \lambda_2 \mathcal{F}_{smooth}(\xi),$$



Still susceptible to poor local minima with a practical time constraint ..

KAIST



Overview

The robot must follow the target path given a fixed end-effector orientation.



Solution: Learning-based Initial Trajectory Generator







Solution: RL-based Initial Trajectory Generator (*RL-ITG*)

• Training Pipeline



Randomly generated path-following problems





Solution: RL-based Initial Trajectory Generator (*RL-ITG*)



Solution: Initial Trajectories ξ_{init} generated by *RL-ITG*





KAIST

Results: Real-world Experiment



- Total execution time: **00:46.76**
 - Average pose error: **2.85 x 10**-3
 - with RL-ITG (ours)

about 220% faster than Greedy

01:43.16 (min: sec) 4.28 x 10⁻³ with Greedy

ΚΔΙST

Benchmark name: 'Random #64' | Trajectory optimizer: TORM [Kang et al. 2020]





Results: Real-world Experiment



- Total execution time: **00:15.23**
 - Trajectory jerkiness: 576.729

with RL-ITG (ours)

about **350%** smoother than Greedy

about 187% faster than Greedy

01:28.01 (min: sec) 2023.25 (^{rad}/_{sec³}) with Greedy

ΚΔΙΣΤ

Benchmark name: 'S' | Trajectory optimizer: TORM [Kang et al. 2020]



Take-Home Message

Hybrid frameworks integrating learning and planning is an important strategy that works in a complementary manner.
→ Improves accuracy and efficiency by combining the two approaches.

Learning-based methods

- \rightarrow may not guarantee optimality
- → but offer a good starting point for optimization quickly.



Optimization-based methods

- → may struggle with highdimensional and non-convex problems
- → but find optimal solutions around starting point by iterative refinement.





Application of RL for Robot Arm Tasks (2)



Towards Safe Remote Manipulation: User Command Adjustment based on Risk Prediction for Dynamic Obstacles

Min-Cheul Kang, Min-Sung Yoon, Sung-Eui Yoon





Korea Advanced Institute of Science and Technology

School of Computing

Remote manipulation

- Performs sophisticated or hazardous tasks on behalf of humans
- Expands to irregular environments around us



Convenience store







Motivation

- A robot accident can be a significant threat.
- A user observes a restricted environment through a camera.
 - \rightarrow A user may not be aware of obstacles.



Example of a robot accident



Restricted environment information

KAIST



https://www.youtube.com/watch?v=mJ6I9thm-8s

Problem

- Avoiding dynamic obstacles depends on a user's judgment.
- We need a method to avoid the risk of dynamic obstacles.
 - The method should minimize the delay of remote manipulation.



Non-risky situation



Risky situation





System flow

 Δx_u : user command Δx_{π} : obstacle avoidance command $\hat{\rho}$: predicted risk for dynamic obstacles Δx_a : adjusted command x_a : end-effector pose for Δx_a







Experimental setup

•We constructed a system for performing a remote manipulation task.



Real environment



Predicted risk 0.963000 User command Adjusted command User's monitor screen

User's remote environment





Real robot experiment







Real robot experiment

S^V_RG



ΚΔΙΣΤ

Enhancing Navigation Efficiency of Quadruped Robots via Leveraging Personal Transportation Platforms



Min-Sung Yoon, Sung-Eui Yoon

Learning-based Adaptive Control of Quadruped Robots for Active Stabilization on Moving Platforms

Min-Sung Yoon, Heechan Shin, Jeil Jeong, Sung-Eui Yoon







Korea Advanced Institute of Science and Technology

School of Computing

Recent Progress of Quadruped Robots



Traversing Challenging Terrains

Agile Locomotion





Battery Limitation of Legged Robots





Unitree B2's Battery Spec.

- Battery Life: **4** ~ **6** hours
- Battery Capacity: 45Ah (2250 Wh)
- Standard Voltage: 50.4V

Unitree Go1's Battery Spec.

- Battery Life: ~ 2.5 hours
- Battery Capacity, 6Ah (133.2Wh)
- Standard Voltage: 22.2V

Actual running time is way more below...





Motivation: Human Mobility Augmentation

• Humans Use Transporters to Move Farther and Faster.



Beyond Footsteps: Transporter-Riding Skills

 Legged animals, like dogs, instinctively use transporters to improve mobility and reduce energy use.







Research Goal

 We aim to ensure that quadruped robots adeptly utilize transportation platforms, also known as transporters, for efficient long-range navigation.







Main Contribution

- We introduce RL-ATR (Reinforcement Learning-based Active Transporter Riding method).
 - Built a simulation with transporter dynamics for reinforcement learning.
 - Trained a transporter riding policy.
 - Added state estimators for stability in non-inertial frames (moving platforms).







Experimental Result





Experimental Result

Distribution of Cost of Transport (CoT) values measured during the path tracking tasks

Legged Locomotion (Mean: 0.3050) Wheel-Legged Locomotion (Mean: 0.1577) 25 Type-1 Transporter Riding (Mean: 0.0376) Type-2 Transporter Riding (Mean: 0.0428) 20 Density 12 10 5 -0 0.2 0.5 0.0 0.1 0.3 0.4 Cost of Transport (CoT) $_{42}$

100m sine path (5m amp, 50m wave)





Summary

- Covered basics of RL and its applications to robot arms and quadrupeds.
- RL enables diverse real-world tasks, showing broad applicability.
- Reward engineering is often needed to guide desired behavior.
- Despite these challenges, RL offers unique and powerful capabilities.





Thanks

Any Questions?





Korea Advanced Institute of Science and Technology

School of Computing