

Supervisory Meeting

Meeting Methods/Venue:

Wechat

Items for discussion:

1. The physical structure of the robot
2. Members in the Lab who I can work with
3. Time for me to go to the lab in Shanghai
4. Useful materials in the field.

Record of discussion:

1. Professor Rosendo currently has a robot with two legs that is currently under construction.
2. Since I come from an EE background it would be nice if I work with it, as at the moment no one is controlling it
3. I am asked that it would be possible to go to Shanghai before the Chinese new year. Since it is just announced today that my last final exam is on Jan. 02, I would go there around that
4. I could do some self-learning before that whatever is relative to the project.

Action list:

1. Finish the visiting student application form and submit it to the ShanghaiTech University.
2. Find some learning material online.

Experiment: Simulation of Walking Robot Based on MATLAB

I try to do some exercises according to the tutorials "MATLAB and Simulink Robotics Arena" (<https://www.youtube.com/playlist?list=PLn8PRpmu08oLufaYWEucuzRq7q4o7D>). First, I begin to with a walking robot model based on Simulink, given by the tutorial. I explore the model, run the simulation and try to understand each part of the model, as shown in the Figures followed. In addition, plots of angles of each joint (ankle roll, ankle pitch, knee, hip pitch, hip roll, and hip yaw) for right and left legs are also attached. (Total time costs: around 3 hours, including the set up of the env.)

8 pm

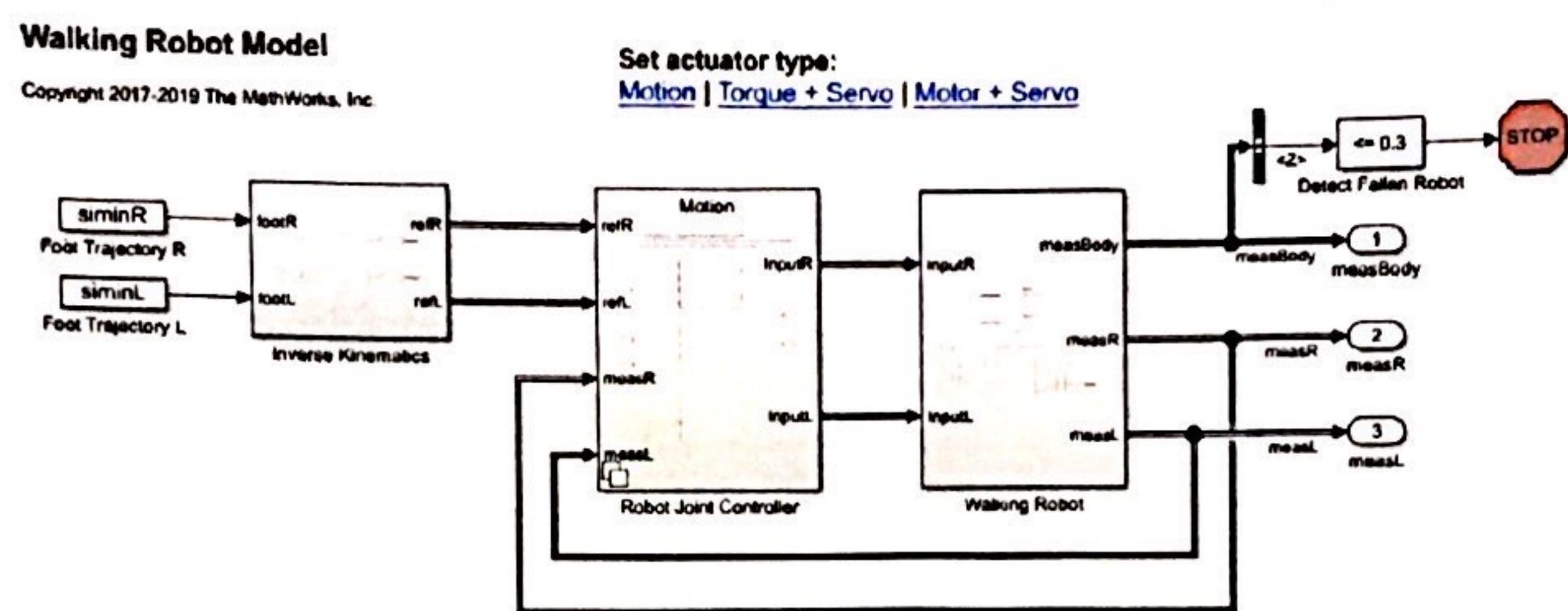


Figure 1. Result of training a walking robot with DDPG on Python (average reward per episode).



Figure 2 Simulation result of the model.

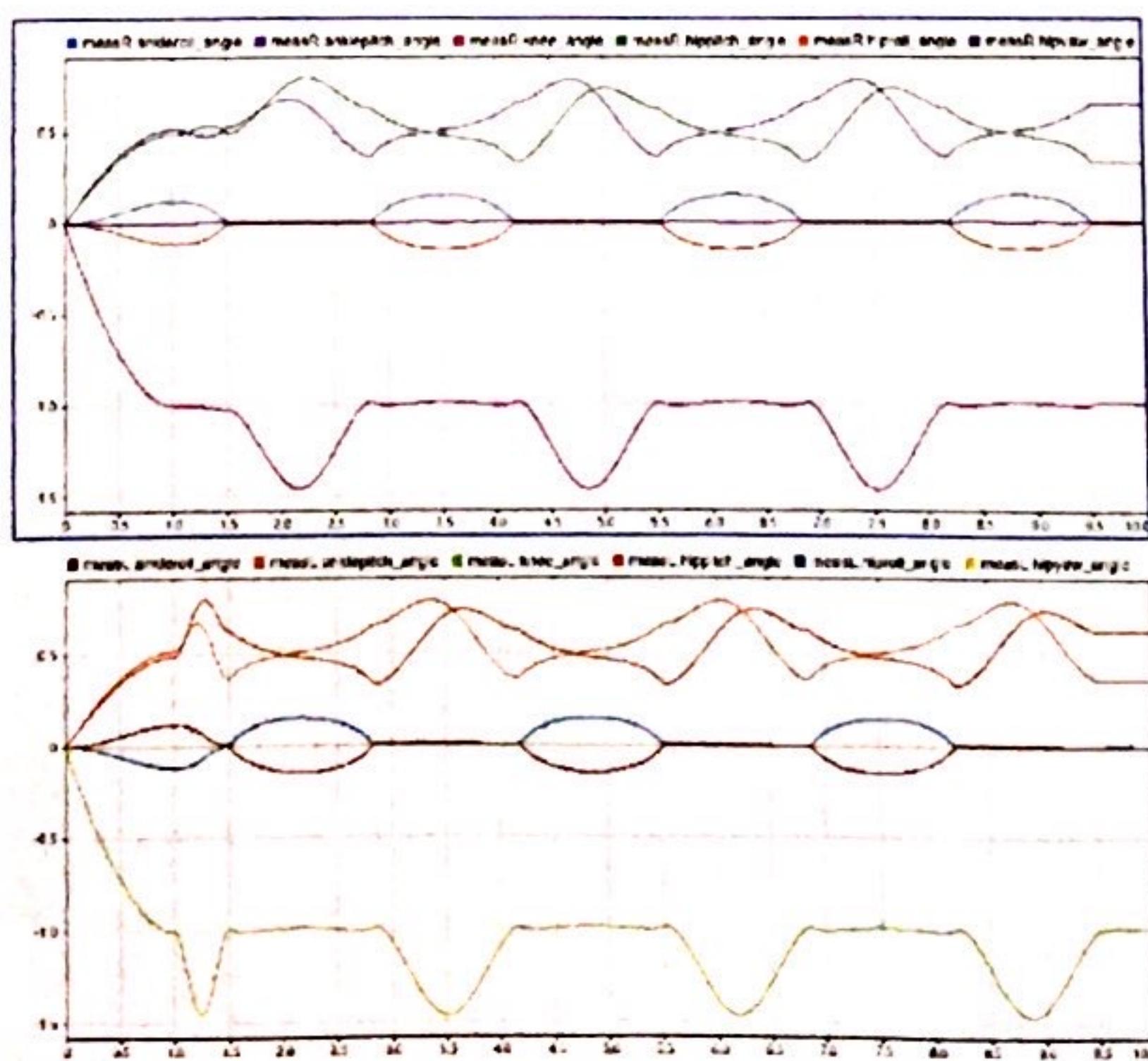


Figure 3. Angles for right leg (upper) and left leg (lower).

Notes on Introduction to Reinforcement Learning

This is the note for the first lecture of the course Introduction to Reinforcement Learning online by David Silver from UCL.

(<https://www.youtube.com/playlist?list=PLqYmG7hTtaZDM-OYHwgPebj2MtCFzFObQ>)

ps: The courses cover most of the basic concepts of RL, which is necessary and meaningful to refresh again and again while diving into more advanced topic.

{ History: the sequence of observation, actions, rewards $H_t = A_1, O_1, R_1, \dots, A_t, O_t, R_t$

State: the information used to determine what happens next, $S_t = f(H_t)$ (summary of history)

↪ Environment state

Agent state: the agent's internal representation

• A state S_t is Markov if and only if $P(S_{t+1} | S_t) = P(S_{t+1} | S_1, \dots, S_t)$

{ Full observability: Agent state = environment state = information state \Rightarrow Markov decision process

Partial observability: Agent indirectly observes environment \Rightarrow Partially observable Markov decision process

\Rightarrow Complete history: $S_t^a = H_t$

Beliefs of environment state: $S_t^a = \{P[S_t^a = s'], \dots, P[S_t^a = s^n]\}$

Recurrent neural network: $S_t^a = \sigma(\sum_w w_s + o_a w_o)$

Policy is the agent's behaviour { Deterministic policy: $a = \pi(s)$

Stochastic policy: $\pi(a|s) = P[A=a | S=s]$

Value function is a prediction of future reward

Model predicts what the environment will do next

{ Transitions: predicts the next state $P_{S|A}^a = P[s'=s' | S=s, A=a]$

Rewards: predicts the next reward $R_{S|A}^a = E[R | S=s, A=a]$

* Prediction: evaluate the policy

Control: find the best policy

Experiment: Training a Walking Robot on MATLAB (GA)

Equipment: PC (require MATLAB 2016b or later)

Files: createSmoothTrajectory.m

doSpeedupTasks.m

optimizeRobotMotion.m main optimization script

simulateWalkingRobot.m

walkingRobotOptim.d

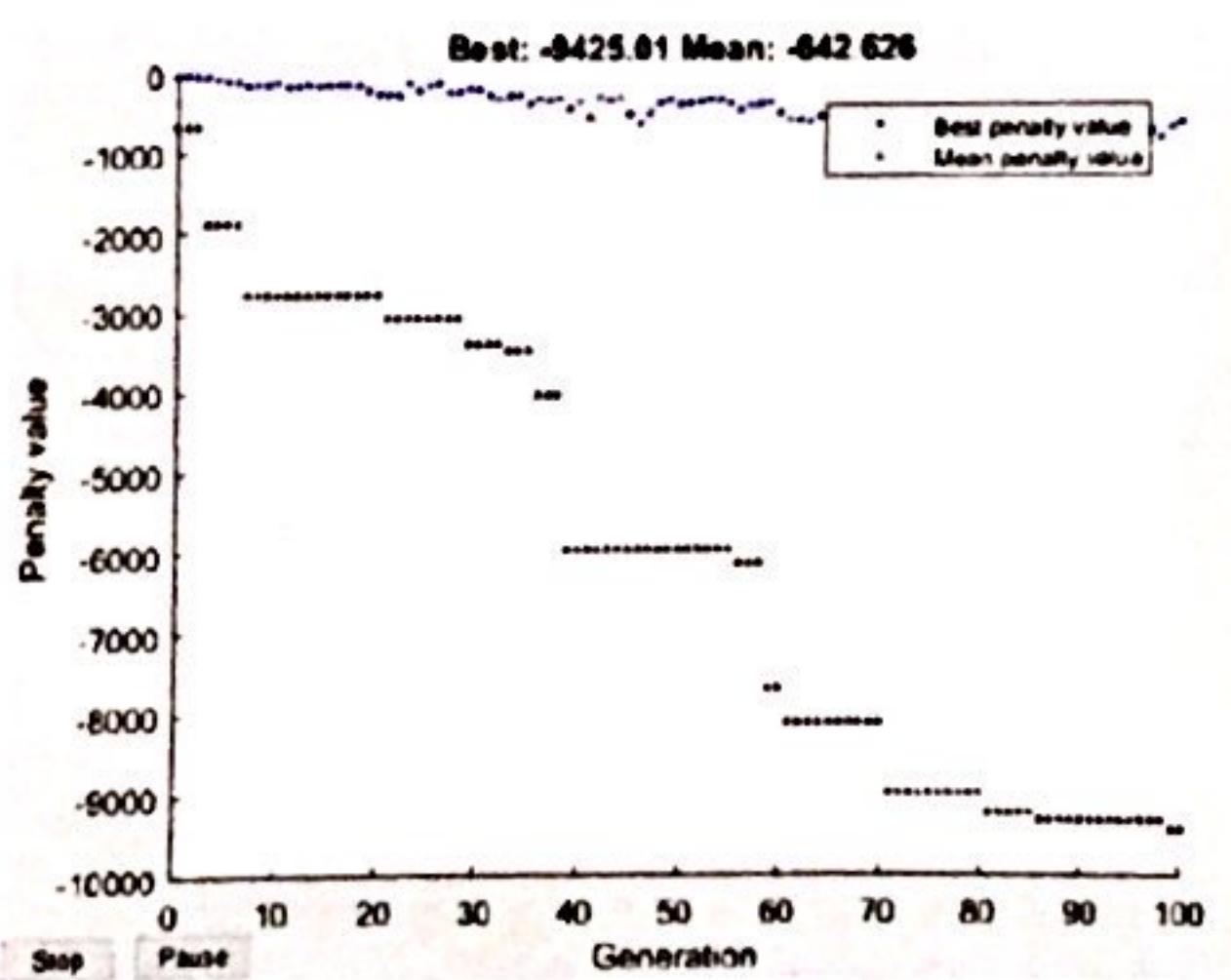
Process:

- ① Tuning the parameter (6:34pm) and run to train.

The screenshot shows the MATLAB environment with the following details:

- Script Editor:** The file `optimizeRobotMotion.m` is open, containing code for setting up a genetic algorithm to optimize a walking robot's motion. It includes parameters like population size (100), number of points (6), and gait period (1.5s).
- Toolbox Window:** A 'Running Optimization Population: 100' window displays the progress of the genetic algorithm. It shows generations from 1 to 16, function counts, and best and mean penalty values. The best value starts at -2.44 and decreases to -27.7.
- Plot Window:** Below the toolbox, a plot titled 'Best and Mean Penalty Value vs Generation' shows the convergence of the algorithm. The y-axis is labeled 'Penalty value' and ranges from 0 to -10000. The x-axis is labeled 'Generation' and ranges from 0 to 100. Two lines are plotted: 'Best penalty value' (blue line with circles) and 'Mean penalty value' (red line with squares). Both lines show a downward trend, indicating improvement over time.

- ② plot the learning results (6:58pm)



Notes on Markov Decision Process

This is the notes for the second lecture of the course Introduction to Reinforcement Learning by David Silver.
 (Almost all RL problems can be formalised as MDPs)

Markov Processes: State transition probability: $P_{ss'} = \Pr(S_{t+1} = s' | S_t = s, A_t = a) \Rightarrow$ transition matrix P

A Markov Process (Markov Chain) is a tuple $\langle \vec{s}, \vec{P} \rangle$

Markov Reward Process \Rightarrow a Markov chain with values \Rightarrow a tuple $\langle S, P, R, \gamma \rangle$

(Immediate) Reward $R_s = \mathbb{E}(R_{t+1} | S_t = s)$

Return $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ (in one sample)

Value Function $v(s) = \mathbb{E}(G_t | S_t = s)$

Bellman Equation for MRP, $v(s) = \mathbb{E}(R_{t+1} + \gamma v(S_{t+1}) | S_t = s) = R_s + \gamma \sum_{s'} P_{ss'} v(s')$

Markov Decision Process \Rightarrow a Markov Reward Process with decision \Rightarrow a tuple $\langle S, A, P, R, \gamma \rangle$

Policy $\pi(a|s) = \Pr(A_t = a | S_t = s)$ - A policy fully defines the behaviour of an agent

- MDP policies depend on the current state (not the history)

$$\text{MDP} \rightarrow \text{MRP } S_1, R_2, S_2, R_3, \dots \leftarrow \langle S, P^\pi, R^\pi, \gamma \right\rangle \Leftarrow \begin{cases} P_{ss'}^\pi = \sum_{a \in A} \pi(a|s) P_{s,a,s'} \\ R_s^\pi = \sum_{a \in A} \pi(a|s) R_s^a \end{cases}$$

} State-value function $v_\pi(s) = \mathbb{E}_\pi(G_t | S_t = s) \rightarrow$ how good is it in state s when following policy π

} Action-value function $q_\pi(s, a) = \mathbb{E}_\pi(G_t | S_t = s, A_t = a) \rightarrow$ Starting from state s , taking action a , and then following policy π

$$\text{Bellman Expectation Equation} \quad \begin{cases} v_\pi(s) = \mathbb{E}_\pi(R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s) = \sum_{a \in A} \pi(a|s) (R_s^a + \gamma \sum_{s' \in S} P_{s,a,s'}^\pi v_\pi(s')) \\ q_\pi(s, a) = \mathbb{E}_\pi(R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a) = R_s^a + \gamma \sum_{s' \in S} \sum_{a' \in A} \pi(a'|s) q_\pi(s', a') \end{cases}$$

Optimal Value Function $v^*(s) = \max_\pi v_\pi(s)$ $q^*(s, a) = \max_\pi q_\pi(s, a)$

Optimal Policy $\pi^* \geq \pi'$ if $v_{\pi^*}(s) \geq v_{\pi'}(s)$, $\forall s$

There is always a deterministic optimal policy for any MDP $\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \arg\max_a q^*(s, a) \\ 0 & \text{otherwise} \end{cases}$

$$\text{Bellman Optimality Equation } v^* = \max_a q^*(s, a) \quad \left\} v^*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{s,a,s'}^\pi v^*(s') \right. \\ q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{s,a,s'}^\pi v^*(s')$$

Notes on Dynamic Programming

This is the notes for the third lecture of the course Introduction to Reinforcement Learning by David Silver.

Problems with properties } Optimal substructure

} Overlapping subproblems \Leftarrow Bellman equation gives recursive decomposition

} Prediction MDP (S, A, P, R, γ) and policy $\pi \rightarrow$ value function V_π

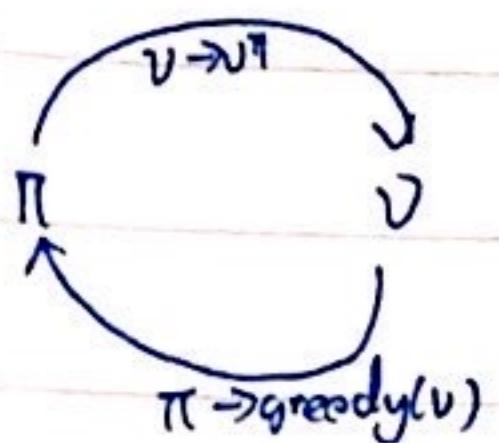
Control MDP $(S, A, P, R, \gamma) \rightarrow$ optimal value function V^* , optimal policy π^*

Policy Evaluation: Evaluate a given policy $\pi \quad V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V_n$

Synchronous backups: each iteration, all states $V_{k+1}(s) = \sum_{a \in A} \pi(a|s) (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s')) \quad \vec{V}^{k+1} = \vec{R}^k + \gamma \vec{P}^k \vec{V}^k$

Policy Iteration: Given a policy $\pi \rightarrow$ Evaluate the policy $V_\pi(s) = \mathbb{E}_{\pi}[R_{t+1} + \gamma R_{t+2} + \dots | s_t = s]$

Improve the policy by action greedily with respect to $V_\pi \quad \pi' = \text{greedy}(V_\pi)$



$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} q_\pi(s, a)$$

Improvement step $\rightarrow V_{\pi'}(s) = \max_{a \in A} q_\pi(s, a)$ the Bellman optimality equation has been satisfied

Value Iteration: Principle of Optimality: A policy $\pi(a|s)$ achieves the optimal value from state s , $V_\pi(s) = V^*(s)$

if and only if } for any state s' reachable from s

π achieves the optimal value from state s' , $V_\pi(s') = V^*(s')$

\Rightarrow if we know the solution to subproblem $V^*(s')$ one-step lookahead $\rightarrow V^*(s) = \max_{a \in A} R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')$

Iterative application of Bellman optimality backup \rightarrow find optimal policy π

$V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V^* \quad \text{update } V_{k+1}(s) \text{ from } V_k(s)$

[C http://www.cs.ubc.ca/~pode/demos/mdp/Vi.html](http://www.cs.ubc.ca/~pode/demos/mdp/Vi.html)

Summary: Prediction: Iterative Policy Evaluation (Bellman Expectation Equation)

Control } Policy Iteration (Bellman Expectation + Greedy Policy Improvement)

} Value Iteration (Bellman Optimality Equation)

Algorithms based on } state-value function $V_\pi(s)$ or $V^*(s)$: Complexity $O(mn^2)$ per iteration, m actions, n states

} action-value function $q_\pi(s, a)$ or $q^*(s, a)$: Complexity $O(m^2n^2)$ per iteration

Organisation of Literatures

This note is for the Interim Report.

"Perception and navigation in unknown environments: the DAPPA robotics challenge" E.J.Molinos [Molinos 2013]

→ The paper presents different techniques to achieve the tasks proposed in DAPPA.

(Not very useful because of the application limitation.)

"A small humanoid robot sdr-4x for entertainment applications" M. Fujita

→ It describes a new small humanoid type robot, SDR-4X, for entertainment purpose.

(It could be used as a reference when dealing with similar robots.)

"Bipedal Robots: Modeling, Design and Walking Synthesis" C. Chevallereau

→ A book published in 2007 introducing bipedal robot in terms of modeling, design and walking synthesis

(A very comprehensive introduction, but just limits state-of-the-art methods)

"Zero-moment point thirty-five years of its life" M. Vukobratovic

→ It is devoted to the permanence of the concept of Zero-Moment Point.

(It is a good material to learn ZMP)

"Modifiable walking pattern of a humanoid robot by using allowable ZMP variation" B.J. Lee

→ It proposes a novel algorithm that can modify a walking period and a step length in both sagittal and lateral planes.

(An application of ZMP)

"Strategies for adjusting the ZMP reference trajectory for maintaining balance in humanoid walking" K. Nishiwaki

→ It addresses strategies of changing the reference trajectories of the future ZMP that are used for online repetitive

(walking pattern generation).

"The development of honda humanoid robot" K. Hirai 1998

→ It presents the mechanism, system configuration, basic control algorithm and integrated function of the Honda humanoid robot.

(a really classic robot design)

"Optimizing foot centers of pressure through force distribution in a humanoid robot" P.M. Wensing 2013

→ It presents two formulation of the FDP for humans in double support, and propose objective functions within a general framework to address the variety of competing requirements for the realization of balance.

"Multi-contact bipedal robotic locomotion" H. Zhao 2019

→ It presents a formal framework for achieving multi-contact bipedal robotic walking.

(Traditional method to control AMBER2 and ATRIAS)

"Toward an expressive bipedal robot: variable gait synthesis and validation in a planar model" U. Huraita 2018

→ It presents a framework for stylized gait generation in a compass-like under-actuated planar biped model.

(Stylish control)

" Real time motion generation and control for biped robot - 4th report: integrated balance control "

→ It proposes a design technique for feedback gains to stabilize the upper body position under varying vertical ground reaction force.

" Robust feedback control of ZMP-based gait for the humanoid robot Nao "

→ It proposes an approach combines the ZMP stability criterion with angular momentum suppression and step timing control.

C! it points out the importance of working in physical platform!)

" Towards natural bipedal walking: virtual gravity compensation and capture point control "

→ It proposes a use of virtual gravity compensation (VGC), essentially a dynamic controller that outputs joint torque, to generate the force and moment as desired.

" Motion and walking stabilization of humanoids using sensory reflex control " J. Kim 2016

→ It employs and elaborates on sensory reflex control to stabilize standing motion and biped walking using basic sensors such as an inertial measurement unit (IMU) and a force-sensing resistor (FSR).

" Observer-based postural balance control for humanoid robots " M. Shahbazi 2013

→ It presents an observer-based balancing control strategy that is robust to persistent perturbations on the ground slope which is based on the estimation of the disturbance without the need for an inertial measurement unit.

[6] " Machine learning algorithms in bipedal robot control " S. Wang 2012

→ A 2012-published paper reviewing the state-of-the-art learning algorithms and their applications to bipedal robot control.

C A really useful survey and the structure of the paper is inspiring!)

* " Combining model-based policy search with online model learning for control of physical humanoids " J. Mordatch 2016

→ It presents an automatic method for interactive control of physical humanoid robots based on high-level tasks.

(Intro) An really interesting and inspiring method that trained on both simulated and real robots in different phases.

" Real time implementation of CT-RNN and BPTT algorithm to learn on-line biped robot balance: Experiments on the standing posture "

→ It uses CT-RNN and BPTT for online learning control laws.

ANN, balance

" Biologically inspired control system for 3D locomotion of a humanoid biped robot " A.A. Sapatra 2016

→ It proposes a control system for 3-D locomotion of a humanoid biped robot based on a biological approach.

RNN, stabilization system

" Rejection of an external force in the sagittal plane applied on a biped robot using a neuro-fuzzy controller " J.P. Ferreira 2009

→ It describes the control of an autonomous biped robot capable to be subjected to external forces applied in the sagittal plane.

(Continues)

- "A new hybrid intelligent control algorithm for a seven-link biped walking robot" F. Farzadpour 2014
 - It proposes the implementation of a new hybrid intelligent control approach for a seven-link biped walking robot.
 - (A method between classic control method and learning method.)
- * "Gait balance and acceleration of a biped robot based on Q-Learning" K.S. Huang 2016
 - It presents a method for the biped dynamic walking and balance control using reinforcement learning.
 - (Train on both simulation and reality environment !!)
- "Motion segmentation and balancing for a biped robot's imitation learning" K.S. Huang 2017
 - It proposes an imitating algorithm called posture-based imitation with balance learning (Post-BL).
 - (It works in physical robot)
- "Posture self-stabilizer of a bipedal robot based on training platform and reinforcement learning" W. Wu 2017
 - An automatic abstraction method for state space is proposed by using the basis function and inner evaluation index to speed up the learning process. (only in simulation)
- * "Survey of model-based reinforcement learning: application on robotics" A.-S. Polydoros 2017
 - A comprehensive survey of model-based RL that applied on robotics
 - (It's a really nice summary of RL algorithm, which could be used as a index of RL literature search. [4, 17, 59, 61, 63])
 - "PILCO : a model-based and data-efficient approach to policy search" M.P. Deisenroth 2011
 - Introduction of a new RL algorithm
 - (A 2011 paper that first introduces this method.)
 - "Efficient reinforcement learning for robots using informative simulated priors" M. Cutler 2015
 - It presents a novel method for transferring data from a simulator to a robot, using simulated data as a prior for real-world learning. (Improved PILCO)
 - "Model-based imitation learning by probabilistic trajectory matching" P. Englot 2013
 - It presents an imitation-learning approach to efficiently learning a task from expert demonstrates.
 - (An application of PILCO)
 - "Improving PILCO with Bayesian Neural Network Dynamics Models" Y. Gal 2017
 - It extends PILCO's framework to use Bayesian deep dynamics models with approximate variational inference, allowing PILCO to scale linearly with number of trials and observation space dimensionality*
 - "Integral control of humanoid balance" B. Stephens 2007
 - It presents a balance controller that allows a humanoid to recover from large disturbance and still maintain an upright posture.

* "Reinforcement learning in robotics: a survey" J. Kober 2013

- It attempts to strengthen the links between the two research communities (Robotics, RL) by providing a survey of work in reinforcement learning for behaviour generation in robots.

"Multi-task policy search for robotics" [Deisenroth 2014]

- An application of PILCO

(It contains a brief and useful introduction of RL)

"Learning to control a low-cost manipulator using data-efficient reinforcement learning" [Durrant-Whyte 2012]

- Another application case of PILCO

(A more detail description of the last application in fact)

"Toward fast policy search for learning legged locomotion" [Deisenroth 2012]

- It presents a data-efficient approach to learning feedback controllers for legged locomotive systems.

(The outcome is really close to the objective of my project!!)

"Model-based imitation learning by probabilistic trajectory matching" [Englot 2013]

- It presents an imitation-learning approach to efficiently learn a task from expert demonstrations.

(It introduces a new concept "imitation learning" but the application case is similar as before)

"Policy gradient based reinforcement learning for real autonomous underwater cable tracking"

- It proposes a field application of a high-level RL control system for solving the action selection problem of an autonomous robot in cable tracking task.

"Comparison of model-free and model-based methods for time optimal hit control of a badminton robot"

- A model-based method is compared to model-free reinforcement learning methods.

(The developed controllers are illustrated on a simulation model of robot.)

"Survey of behaviour learning applications in robotics - state of the art and perspectives."

- It analyzes the current state of machine learning for robotic behaviours.

Not only focus on algorithms but also real world application. (a higher level of view)

"A survey of dynamic robot legged locomotion"

- It proposes a survey focuses on legged dynamic locomotion of robots.

Main part of the paper is focus on history and introduction.

#

Draft for literature organisation

Walking { six legs }

Kirchner (1997)

quadrupedal { Siedewell and Livingston (2007)

Kohl and Stone (2004)

Bartek et al. (2016)

biped { Benbrahim and Franklin

Matsubara et al. (2005)

Geng et al. (2006)

Kormushev et al. (2011c)

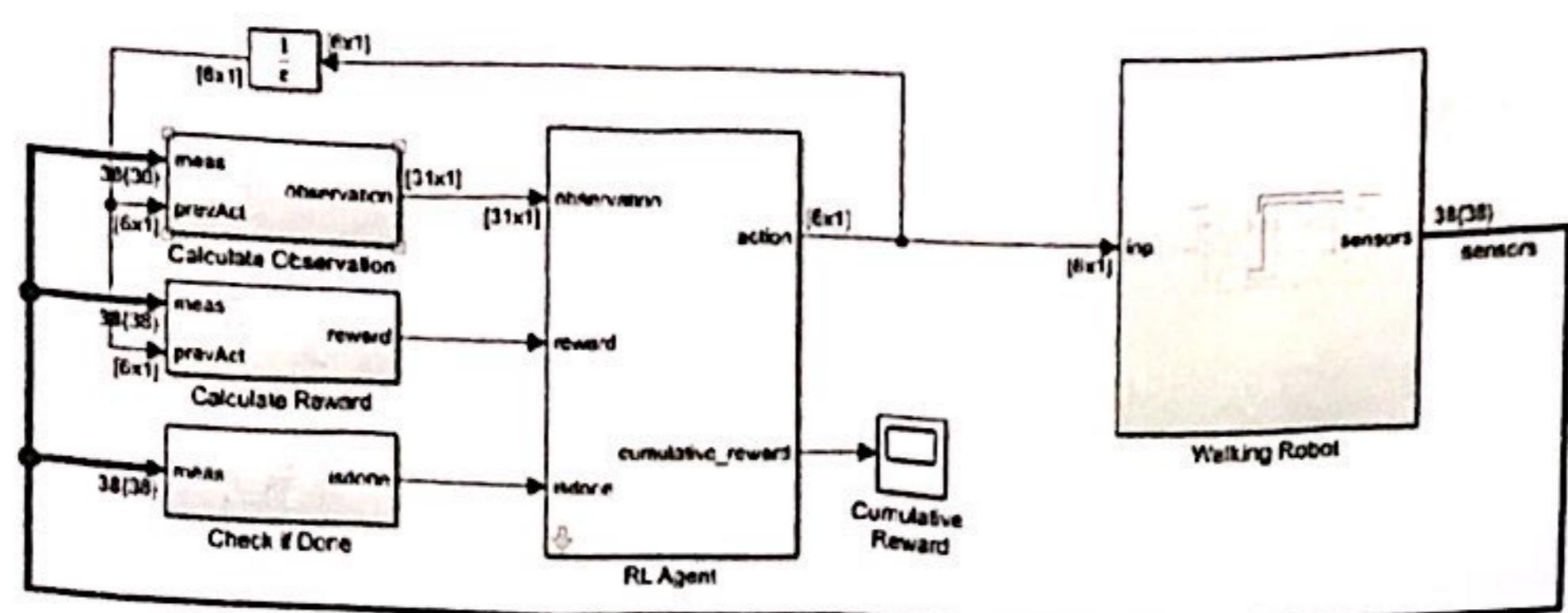
Misra and Behnke (2015)

Experiment: Training a robot on MATLAB using DDPG

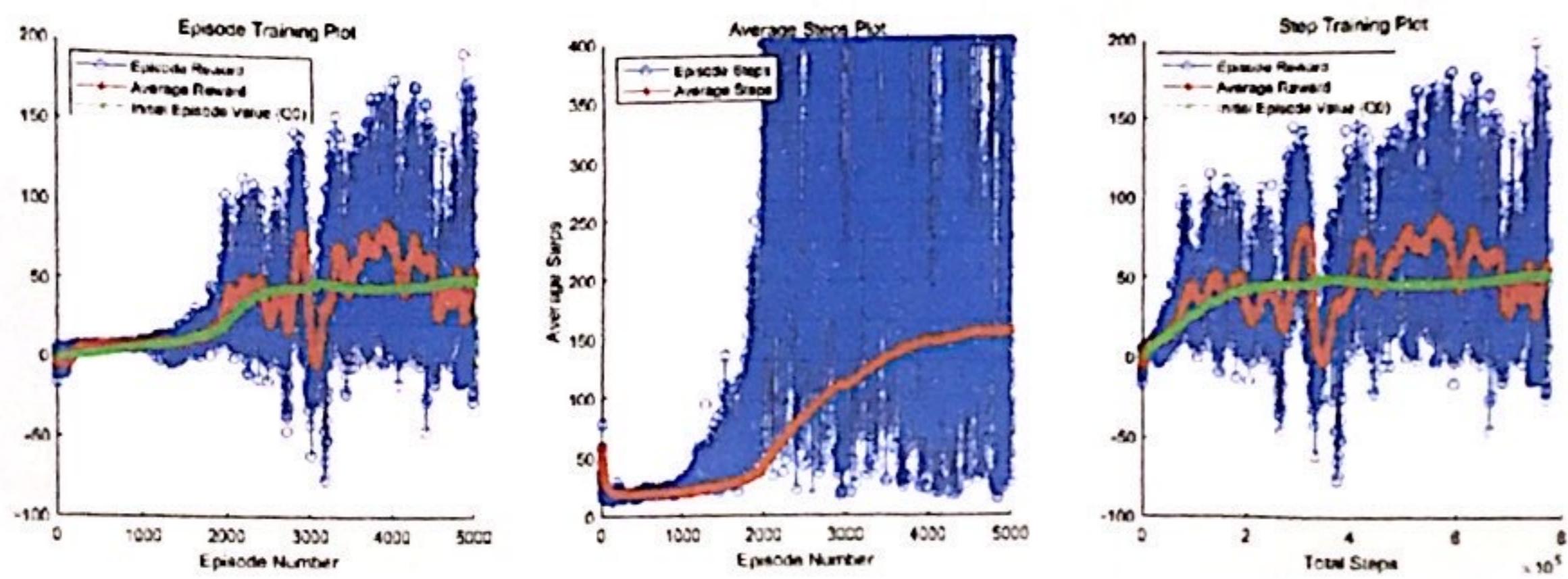
Equipment: PC (Require MATLAB R2019b or later)

Files: createDDPGNetworks.m, createDDPGOption.m, CreateWalkingAgent2D.m, createWalkingAgent3D.m, plotTrainingResults.m, robotParametersRL.m, walkerInkin.m, walkerResetFn.m, walkingRobotRL2D.slx, (walking Robot RL2D.sxc), walkingRobotRL3D.slx.

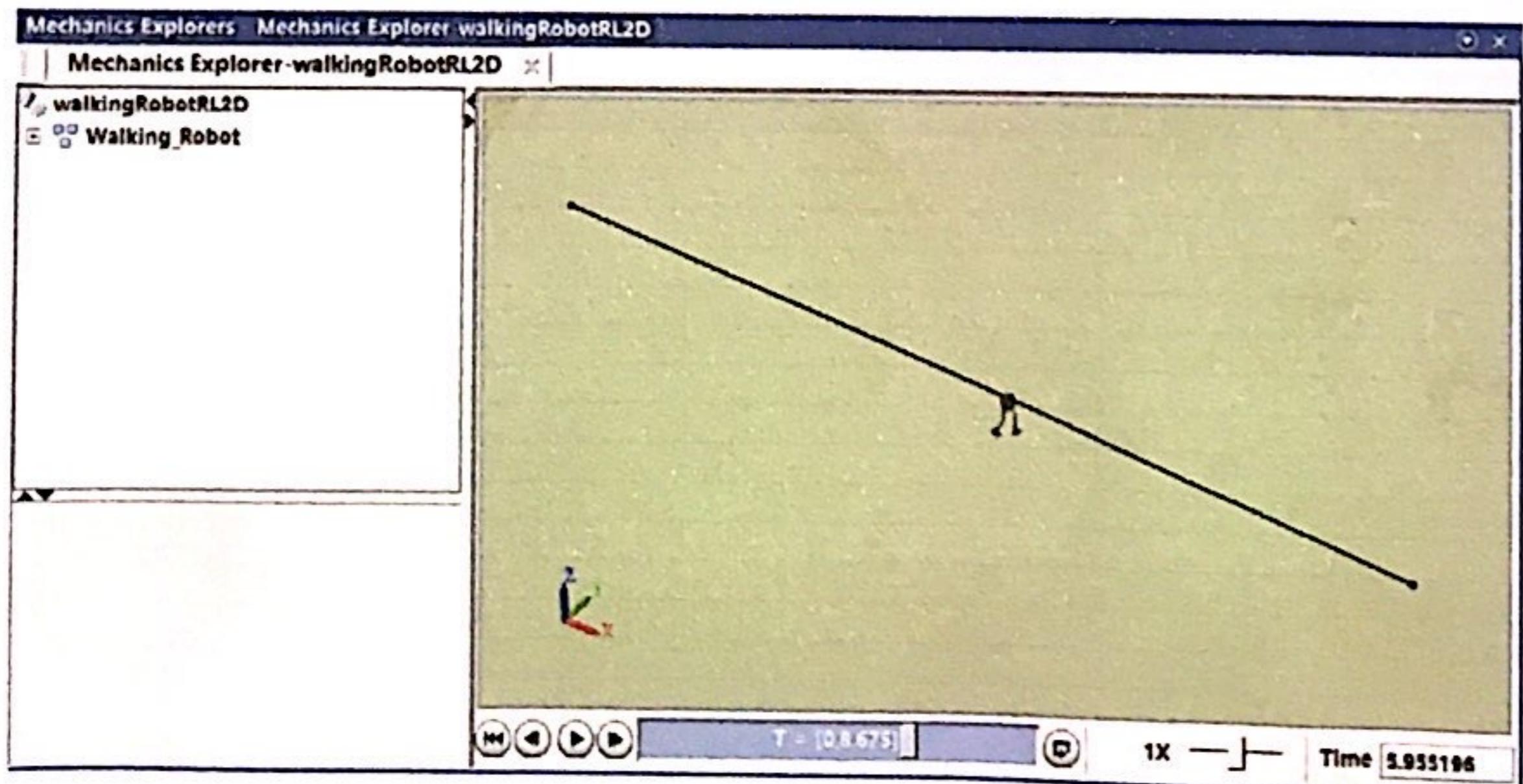
Process: ① Explode the model and run. (1:36pm)



② Got the trained results and plot it. (5:52 pm)



③ Demonstrate the simulated results. (6:12 pm)



Notes on Model-Free Prediction

This is a lecture note for the fourth lecture of the course Introduction to Reinforcement Learning by David Silver

Model-free: no one tell the environment (Estimate the value function of an unknown MPP)

Monte-Carlo learning (Learn directly from episodes of experience)

Value = mean return \Rightarrow can only apply MC to episodic MPPs

Goal: learn v_π from episodes of experience under policy π $S_1, A_1, R_1, \dots, S_k \sim \pi$

Recall of return: $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$

Value function $v_\pi(s) = \mathbb{E}_\pi(G_t | S_t=s)$ expected return \xrightarrow{MC} empirical mean

To evaluate state s : The first time-step t that state s is visited in an episode (First-Visit MC Policy Evaluation)

Every time step t that state s is visited in an episode (Every-Visit MC Policy Evaluation)

- Increment counter $N(s) \leftarrow N(s) + 1$

- Increment total return $S(s) \leftarrow S(s) + G_t$

- Value is estimated by mean return $v_\pi(s) = S(s)/N(s)$

- By law of large numbers, $v_\pi(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

To evaluate state s : Every time step t that state s is visited in an episode

For each state s_t with return G_t : $N(s_t) \leftarrow N(s_t) + 1$

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)}(G_t - V(s_t))$$

Temporal-Difference Learning TD learns from incomplete episodes, by bootstrapping (Update a guess toward a guess)

Goal: learn v_π online from experience under policy π

{ Incremental every-visit MC: $V(s_t) \leftarrow V(s_t) + \alpha(G_t - V(s_t))$

$$TD(\alpha): V(s_t) \leftarrow V(s_t) + \alpha(R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

{ MC converges to solution with minimum mean-squared error * MC does not exploit Markov property

TD(α) converges to solution of max likelihood Markov model * TD exploits Markov property