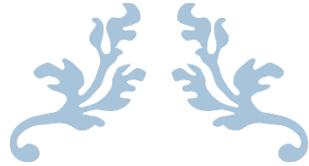


S. Eben Li. <Reinforcement Learning for Decision-making and Control>



REINFORCEMENT LEARNING FOR DECISION-MAKING AND CONTROL

Full Chapters



Shengbo Eben Li
(李升波)

Email: lishbo@tsinghua.edu.cn
Intelligent Driving Laboratory (iDLab)
www.idlab.tsinghua.edu.cn

TSINGHUA UNIVERSITY
Beijing, 100084, China

DECEMBER 25, 2021

About the Author

Dr. Li is the full professor at Tsinghua University. He is now leading Intelligent Driving Laboratory (iDLab) at School of Vehicle and Mobility. His active research interests include intelligent vehicles and driver assistance, reinforcement learning, optimal control and estimation, etc.

He is the author of over 130 peer-reviewed journal/conference papers, and the co-inventor of over 30 patents. Dr. Li is the recipient of best paper (student) awards of IEEE ITSC 2020/2021, ICCAS 2020 IEEE ICUS 2020, CCCC 2018/2019, ITSAPF 2015, IEEE ITSC 2014, etc. His academic services include Member of Board of Governors of IEEE ITS Society, Senior Editor of IEEE OJ ITS, AEs of IEEE ITSM, IEEE Trans ITS, JICV, and Automotive Innovation, etc.

Contact information:

Shengbo Eben Li (李升波)
Ph.D., Tenured Professor
Rm 643, Lee Shau Kee Sci&Tech Building,
Tsinghua University, 100084, Beijing
Office: 86-10-62796150
Email: lishbo@tsinghua.edu.cn
Website: www.idlab.tsinghua.com

Preface/前言

自 2000 年以来，人工智能的快速崛起正重塑人类社会的各个角落，有望引导工业文明进入第四次革命浪潮。以道路交通为例，汽车的智能化变革促使整个行业发生了翻天覆地的变化，包括驾驶辅助、自动驾驶、云控协同等一系列新技术如雨后春笋般涌现，它们在提升地面车辆行驶性能的同时，也为解决交通事故、排放污染、城市拥堵等问题提供了一条可行的途径。近年随着机器学习和自动控制的融合发展，以模仿人类大脑学习机制为原理的强化学习（RL, Reinforcement Learning）技术迅速进入人们的视野，它为大规模复杂动态系统的高性能决策与高实时控制提供了一套极具前景的解决方案。一个引人注目的成功案例是以 Alpha Go 为代表的围棋智能：它利用深度强化学习算法实现围棋智能的自我进化，自我超越，以超乎想象的速度进化出打败人类专业棋手的能力，引发学术界和工业界的热切关注。

尽管强化学习具有优异的潜在优势，但是该方法的工程应用尚属于起步阶段。一个重要的原因是该方法既具有理论学习的复杂度，又具有工程实践的挑战性。该方法隶属于统计学习、最优控制、最优化三者的交叉结合部，涉及的数理基础较深，知识涵盖面较广，难以学习更难工程应用。入门者不易短期内掌握关联的理论体系，若对原理不够熟悉，难以对代码进行针对性调整，不能发挥算法应有的性能。为了应对上述挑战，依托笔者在清华大学开设的研究生课程《强化学习与控制》，撰写了这一本教学参考书，主要面向工程应用领域的科研工作者和技术人员，按照原理剖析、主流算法、典型示例的架构，介绍用于动态系统的决策与控制的强化学习方法。所涉及的知识点包括马尔科夫决策、蒙特卡洛学习、时序差分学习、函数近似学习、策略梯度学习、近似动态规划、深度强化学习等。希望本书为该领域的行业同仁，包括研究生和本科生，提供一本体系较为全面，且适合研究者学习和应用的参考书籍。

全书总共包括 11 章。第 1 章介绍 RL 概况，包括发展历史、知名学者、典型应用以及主要挑战等。第 2 章介绍 RL 的基础知识，包括定义概念、自洽条件、最优化原理与问题架构等。第 3 章介绍免模型学习的蒙特卡洛法，包括 Monte Carlo 估计，On-policy/Off-policy，重要性采样等。第 4 章介绍免模型学习的时序差分法，包括它衍生的 Sarsa, Q-learning, Expected Sarsa 等算法。第 5 章介绍带模型学习的动态规划法，包括策略迭代、值迭代、通用迭代架构与收敛性证明等。第 6 章介绍间接型 RL 的函数近似方法，包括常用近似函数，值函数近似，策略函数近似以及所衍生的 Actor-critic 架构等。第 7 章介绍直接型 RL 的策略梯度法，包括 On-policy gradient, Off-policy gradient，以及它们的代价函数、优化算法等。第 8 章介绍带模型的强化学习，即近似动态规划（ADP），包括离散时间系统的 ADP，连续时间系统的 ADP，以及 ADP 与 MPC 的联系与区别等。第 9 章介绍有限时域的近似动态规划，同时探讨了状态约束的处理手段以及它与求解可行性之间的关系。第 10 章介绍深度强化学习，即以神经网络为载体的 RL，包括神经网络的原理与训练，深度化挑战以及 DQN、DDPG、TD3、TRPO、DSAC 等典型深度化算法。第 11 章介绍 RL 的各类拾遗，包括鲁棒性、POMDP、多智能体、元学习、逆强化学习、离线强化学习、训练框架与平台等。

借此机会，诚挚感谢清华大学智能驾驶课题组的老师和同学们，他们为本书的撰写付出了极大的心血与努力，包括文献查阅、文字编排、公式推导、示例代码等

一系列工作，贡献良多。还有很多师长、同仁和朋友，也为本书的撰写和完善提供了大量宝贵地建议，受篇幅所限不能一一赘述，借此机会深表感谢。

受笔者的水平所限，书中难免存在疏漏和不足，包括图片版权冲突、公式推导错误，文字写作粗陋等不一而足。借此机会，恳请国内外同仁给予批评和指正。如有发现任何错误，烦请将意见发送给笔者邮箱：lishbo@tsinghua.edu.cn。每一位读者的反馈将是本书进一步完善的重要保障，提前谢谢大家！

Preface/前言

From the beginning of the 21st century, artificial intelligence (AI) is reshaping almost all areas of human society, which is promising to spark the fourth industrial revolution. Noticeable examples can be found in the sector of road transportation, where AI has drastically changed automobile design and traffic management. Lots of new technologies, such as driver assistance, autonomous driving, and cloud-based cooperation, are rising in an unbelievable speed. These technologies have the potential to significantly improve driving ability, reduce traffic accidents, and relieve urban congestion.

As one of the most important AI branches, reinforcement learning (RL) is attracting increasing attention in the past decades. RL is an interdisciplinary field of trial-and-error learning and optimal control, which provides a promising solution for decision-making and control of large-scale and complex dynamic processes. One of its most eye-catching success is AlphaZero from Google DeepMind, which beats the most professional human player. The key technology behind is called deep reinforcement learning, which equips AlphaGo with an amazing self-evolution ability.

Despite a few success, the application of RL is still in its infancy stage. This is because most RL algorithms are difficult to comprehend. In one hand, RL deeply connects with statistic learning and convex optimization, and involves a wide range of concepts and theories. On the other hand, it is a tedious and long learning process for a beginner to become an RL master. Without fully understanding those principles, it is very difficult for users to make necessary adjustments to achieve the best performance. This book aims to provide a systematic introduction of fundamental RL theories, mainstream RL algorithms and typical RL applications to fellow researchers and engineers. The topics mainly include Markov Decision Processes, Monte Carlo learner, Temporal Difference learner, RL with function approximation, policy gradient method, approximate dynamic programming, deep reinforcement learning, etc.

The book contains 11 chapters. Chapter 1 provides an overview of RL, including its history, famous scholars, successful examples and up-to-date challenges. Chapter 2 briefs the basis of RL, including its concepts, optimality conditions, and problem formulation. Chapter 3 introduces Monte Carlo methods for model-free RL, including on-policy/off-policy and importance sampling technique. Chapter 4 introduces temporal-difference methods for model-free RL, including Sarsa, Q-learning, expected Sarsa, etc. Chapter 5 introduces stochastic dynamic programming, i.e., model-free RL with tabular representation, including value iteration, policy iteration and their convergence mechanism. Chapter 6 introduces how to approximate policy and value function in indirect RLs, as well as its associated actor-critic architecture. Chapter 7 derives different kinds of direct policy gradients, including likelihood ratio gradient, natural policy gradient and a few variants. Chapter 8 introduces infinite horizon ADP, and its connection with model predictive control. Chapter 9 introduces finite-horizon ADP, and puts great emphasis on how to handle state constraints. Chapter 10 devotes to deep reinforcement learning, including artificial neural networks and typical deep RL algorithms like DQN, DDPG, TD3, TRPO, DSAC, etc. Finally, Chapter 11 provides various RL relics, including robust RL, POMDP, multi-agent RL, meta-RL, inverse RL, offline RL, major RL libraries and platforms, etc.

In closing, I wish to offer my sincere gratitude to all the faculties and students in Intelligent Driving Laboratory (iDLab) for their great contribution to this book. I also express my deep appreciation to those friends and colleagues, who support writing and

polishing this book. They have provides numerous priceless suggestions. Any comments and corrections from readers would be much appreciated. I look forward to seeing your email at lishbo@tsinghua.edu.cn. Thanks a lot in advance!

Contents

0	Summary of Notations	15
0.1	Symbols	15
0.2	Abbreviation.....	21
1	Introduction of Reinforcement Learning	25
1.1	History of RL	25
1.1.1	Dynamic Programming.....	26
1.1.2	Trial-and-Error Learning.....	28
1.2	Examples of Using RL.....	30
1.2.1	Tic-Tac-Toe.....	30
1.2.2	Go Game	31
1.2.3	Autonomous Vehicles	32
1.3	Key Challenges in Today's RL.....	33
1.3.1	Exploration-Exploitation Dilemma	33
1.3.2	Uncertainty and Partial Observability.....	34
1.3.3	Temporally Delayed Reward	34
1.3.4	Infeasibility from Safety Constraint	34
1.3.5	Entangled Stability and Convergence	34
1.3.6	Non-stationary Environment	34
1.3.7	Lack of Generalizability	35
1.4	References [In chronological order].....	35
2	Principles of RL Problems.....	37
2.1	Four Elements of RL Problems	38
2.1.1	Environment Model	38
2.1.2	State-Action Sample.....	39
2.1.3	Policy.....	39
2.1.4	Reward Signal.....	40
2.2	Classification of RL Methods	44
2.2.1	Definition of RL Problems	44
2.2.2	Bellman's Principle of Optimality.....	46
2.2.3	Indirect RL Methods.....	48
2.2.4	Direct RL Methods	49
2.3	Revisiting RL in A Broad View	50

2.3.1	Influence of Initial State Distribution	50
2.3.2	Difference between RL and MPC	51
2.3.3	Combination of Four Elements	53
2.4	Measures of Learning Performance	54
2.4.1	Policy Performance	54
2.4.2	Learning Accuracy	55
2.4.3	Learning Speed.....	55
2.4.4	Sample Efficiency	56
2.4.5	Approximate Accuracy	56
2.5	Two Examples of Markov Decision Processes	56
2.5.1	Example: Indoor Cleaning Robot	57
2.5.2	Example: Autonomous Driving System.....	58
2.6	References [In chronological order].....	60
3	Model-Free Indirect RL: Monte Carlo	61
3.1	MC Policy Evaluation	61
3.2	MC Policy Improvement.....	63
3.2.1	Greedy Policy	64
3.2.2	Policy Improvement Theorem	64
3.2.3	MC Policy Selection.....	66
3.3	On-Policy Strategy vs. Off-Policy Strategy.....	66
3.3.1	On-Policy Strategy.....	67
3.3.2	Off-Policy Strategy	69
3.4	Understanding Monte Carlo RL from a Broad Viewpoint.....	73
3.4.1	On-Policy MC.....	74
3.4.2	Off-Policy MC	75
3.4.3	Incremental Estimation of Value Function	76
3.5	Example of Monte Carlo RL.....	78
3.5.1	Cleaning Robot in a Grid Room	79
3.5.2	MC with Action-Value Function	80
3.5.3	Influences of Key Parameters	81
3.6	References [In chronological order]	84
4	Model-Free Indirect RL: Temporal Difference	85
4.1	TD Policy Evaluation	86
4.2	TD Policy Improvement	87

4.2.1	On-Policy Strategy.....	87
4.2.2	Off-Policy Strategy	88
4.3	Typical TD Learning Algorithms.....	89
4.3.1	On-Policy TD: Sarsa	89
4.3.2	Off-Policy TD: Q-Learning.....	91
4.3.3	Off-Policy TD: Expected Sarsa	94
4.3.4	Recursive Value Initialization	95
4.4	Unified View of TD and MC	95
4.4.1	N-Step TD Policy Evaluation.....	96
4.4.2	TD-Lambda Policy Evaluation.....	97
4.5	Examples for Temporal Difference.....	98
4.5.1	Results of Sarsa	99
4.5.2	Results of Q-Learning.....	101
4.5.3	Comparison of MC, Sarsa, and Q-Learning	102
4.6	References [In chronological order]	103
5	Model-Based Indirect RL: Dynamic Programming	105
5.1	Stochastic Sequential Decision.....	106
5.1.1	Models for Stochastic Environments	106
5.1.2	Discounted Cost vs Average Cost.....	107
5.1.3	Policy Iteration vs Value Iteration.....	111
5.2	Policy Iteration Algorithms.....	112
5.2.1	Policy Evaluation (PEV)	112
5.2.2	Policy Improvement (PIM)	114
5.2.3	Proof of Convergence	115
5.2.4	Explanation with Newton-Raphson Mechanism.....	117
5.3	Value Iteration Algorithms	120
5.3.1	Explanation with Fixed-point iteration Mechanism.....	121
5.3.2	Convergence of DP Value Iteration.....	121
5.3.3	Value Iteration for Problems with Average Costs.....	122
5.4	Stochastic Linear Quadratic Control.....	125
5.4.1	Average Cost LQ Control.....	125
5.4.2	Discounted Cost LQ Control.....	126
5.4.3	Performance Comparison with Simulations	127
5.5	More Viewpoints about DP	129

5.5.1	Unification of Policy Iteration and Value Iteration	130
5.5.2	Unification of Model-Based and Model-Free RLs	131
5.5.3	PEV with Other Fixed-point Iterations	133
5.5.4	Value Iteration with Other Fixed-Point Iterations	134
5.5.5	Finite-Horizon DP with Backward Recursion.....	136
5.6	More Definitions of Better Policy and Their Convergence.....	138
5.6.1	Penalize Greedy Search with Policy Entropy.....	138
5.6.2	Expectation-Based Definition for Better Policy.....	139
5.6.3	Another Form of Expectation-Based Condition	142
5.7	Example: Autonomous Car on Grid Road	144
5.7.1	Results of DP	145
5.7.2	Influences of Key Parameters on DP	146
5.7.3	Influences of Key Parameters on Sarsa and Q-learning	148
5.8	Appendix: Fixed-point iteration Theory	150
5.8.1	Procedures of Fixed-Point Iteration.....	150
5.8.2	Fixed-Point Iteration for Linear Equations.....	151
5.9	References.....	152
6	Indirect RL with Function Approximation	153
6.1	Linear Approximation and Basis Functions	154
6.1.1	Binary Basis Function	154
6.1.2	Polynomial Basis Function	155
6.1.3	Fourier Basis Function.....	155
6.1.4	Radial Basis Function	155
6.2	Parameterization of Value Function and Policy	156
6.2.1	Parameterized Value Function.....	156
6.2.2	Parameterized Policy	157
6.2.3	Choice of Tabular Function or Approximation Function	159
6.3	Value Function Approximation.....	160
6.3.1	On-Policy Value Approximation	161
6.3.2	Off-Policy Value Approximation.....	164
6.3.3	Deadly Triad Issue	167
6.4	Policy Approximation	169
6.4.1	Indirect On-Policy Gradient.....	170
6.4.2	Indirect Off-Policy Gradient	171

6.4.3	Revisit Better Policy for More Designs.....	172
6.5	Actor-Critic Architecture from Indirect RL	174
6.5.1	Generic Actor-Critic Algorithms	175
6.5.2	On-policy AC with Action-Value	177
6.5.3	On-policy AC with State-Value	180
6.6	Example: Autonomous Car on a Circular Road.....	182
6.6.1	Autonomous Vehicle Control Problem	182
6.6.2	Design of On-policy Actor-Critic Algorithms	184
6.6.3	Training Results and Performance Comparison.....	185
6.7	References [In chronological order]	189
7	Direct RL with Policy Gradient	190
7.1	Overall Objective Function for Direct RL.....	191
7.1.1	Stationary Properties of Markov Process	192
7.1.2	Discounted Objective Function	194
7.1.3	Average Objective Function	196
7.2	Likelihood Ratio Gradient in General	196
7.2.1	Gradient Descent Algorithm	196
7.2.2	Method I: Gradient Derivation from Trajectory Concept	197
7.2.3	Method II: Gradient Derivation from Cascading Concept.....	200
7.3	On-Policy Gradient	201
7.3.1	On-Policy Gradient Theorem	202
7.3.2	Extension of the On-Policy Gradient.....	203
7.3.3	How a Baseline Works	206
7.4	Off-Policy Gradient.....	209
7.4.1	Off-Policy True Gradient	210
7.4.2	Off-Policy Quasi-Gradient	210
7.5	Actor-Critic Architecture from Direct RL	212
7.5.1	Off-policy AC with Advantage Function	213
7.5.2	Off-policy Deterministic Actor-Critic.....	215
7.5.3	State-of-The-Art of AC Algorithms	216
7.6	Miscellaneous Topics in Direct RL Algorithms.....	219
7.6.1	Tricks in Stochastic Derivative.....	219
7.6.2	Types of Numerical Optimization	221
7.6.3	Some Variants of Objective Function.....	224

7.7	References [In chronological order]	226
8	Infinite Horizon Approximate Dynamic Programming	228
8.1	Discrete-time Problems with Infinite Horizon	229
8.1.1	Bellman Equation.....	230
8.1.2	Discrete-time ADP Framework	233
8.1.3	Convergence and Stability	234
8.1.4	Convergence of Inexact Policy Iteration	238
8.1.5	Discrete-time ADP with Parameterized Function	239
8.2	Continuous-time Problems with Infinite Horizon	242
8.2.1	Hamilton-Jacobin-Bellman (HJB) Equation	243
8.2.2	Continuous-time ADP Framework	244
8.2.3	Convergence and Stability	245
8.2.4	Continuous-time ADP with Parameterized Function	247
8.2.5	Linear Quadratic Control Problem	248
8.3	How to Run RL and ADP Algorithms.....	249
8.3.1	Two Kinds of Environment Models	250
8.3.2	Implementation of RL/ADP Algorithms.....	251
8.3.3	Training Efficiency Comparison: RL VS ADP	252
8.4	ADP for Tracking Problem and its Policy Structure	254
8.4.1	Two Kinds of Time Domain in RHC.....	254
8.4.2	Traditional Real-time Definition of Tracking ADP	255
8.4.3	Reformulate Tracking ADP in Virtual-time Domain	256
8.4.4	Quantitative Analysis with Linear Quadratic Control.....	258
8.5	Example: Lane Keeping Control on Curved Road	261
8.5.1	Lane Keeping Problem Description	261
8.5.2	Design of ADP Algorithm with Neural Network	262
8.5.3	Training Results of ADP and MPC	263
8.6	References [In chronological order]	267
9	Finite Horizon ADP and State Constraints.....	269
9.1	Finite Horizon ADP in Discrete-time Domain	270
9.1.1	Optimal Regulator.....	272
9.1.2	Optimal Tracker with Multi-stage Policy.....	275
9.1.3	Optimal Tracker with Recurrent Policy	278
9.2	Input Constraints.....	281

9.2.1	Saturated Policy Function	281
9.2.2	Penalized Utility Function	282
9.3	State Constraints and Feasibility	283
9.3.1	Understand State Constraints in Two Domains	284
9.3.2	Definition of Infeasibility.....	286
9.3.3	Types of State Constraints and Feasibility Analysis	288
9.3.4	Classification of Constrained OCP Methods.....	292
9.3.5	Penalty Function Methods	296
9.3.6	Methods of Lagrange Multipliers.....	297
9.4	Feasible Descent Direction Methods	301
9.4.1	Feasible Direction and Descent Direction	301
9.4.2	Constrained LP Optimization	304
9.4.3	Constrained QP Optimization	306
9.5	Actor-Critic-Scenery Architecture	308
9.5.1	Two Kinds of ACS Architectures	308
9.5.2	Reachability-based Scenery Update	309
9.5.3	Solvability-based Scenery Update.....	313
9.6	Safety Consideration in RL/ADP	316
9.6.1	Two Basic Modes to Train Safe Policy.....	317
9.6.2	Safety Guarantee during Online Training	319
9.6.3	Safety Shield Mechanism in Implementation	324
9.7	References [In chronological order]	325
10	Deep Reinforcement Learning	330
10.1	Introduction to Artificial Neural Networks.....	331
10.1.1	Neurons	331
10.1.2	Layers.....	333
10.1.3	Typical Neural Networks.....	336
10.2	Training of ANNs	337
10.2.1	Loss Function	338
10.2.2	Training Algorithms	339
10.3	Challenges of Deep Reinforcement Learning	341
10.3.1	Challenge: Non-iid Sequential Data	342
10.3.2	Challenge: Easy Divergence	343
10.3.3	Challenge: Overestimation	345

10.3.4	Challenge: Sample Inefficiency	349
10.4	Deep RL Algorithms	351
10.4.1	Deep Q-Network.....	352
10.4.2	Double DQN.....	353
10.4.3	Trust Region Policy Optimization	354
10.4.4	Proximal Policy Optimization.....	354
10.4.5	Asynchronous Advantage Actor-Critic.....	355
10.4.6	Deep Deterministic Policy Gradient	356
10.4.7	Twin Delayed DDPG.....	357
10.4.8	Soft Actor-Critic	357
10.4.9	Distributional SAC.....	359
10.5	References [In chronological order]	360
11	Miscellaneous RL Topics	362
11.1	Robust RL with Bounded Uncertainty	363
11.1.1	H-infinity Control and Zero-Sum Game	364
11.1.2	Linear Version of Robust RL.....	366
11.1.3	Nonlinear Version of Robust RL.....	368
11.2	Partially Observable MDP.....	370
11.2.1	Problem Description of POMDP	370
11.2.2	Linear Quadratic Gaussian Control.....	372
11.2.3	Belief State and Separation Principle	375
11.3	Meta Reinforcement Learning	378
11.3.1	Transferable Experience	379
11.3.2	Transferable Policy	380
11.3.3	Transferable Loss.....	382
11.4	Multi-agent Reinforcement Learning	382
11.4.1	Stochastic Multi-Agent Games	384
11.4.2	Fully Cooperative RL	385
11.4.3	Fully Competitive RL	386
11.4.4	Multi-agent RL with Hybrid Rewards.....	387
11.5	Inverse Reinforcement Learning	388
11.5.1	Essence of Inverse RL.....	389
11.5.2	Max Margin Inverse RL	389
11.5.3	Max Entropy Inverse RL.....	390