

Introduction to Reinforcement Learning: Part IV

Dr. Mayank Shekhar JHA
Slides contribution: Soha KANSO

Maitre de Conférences (Associate Professor),
CRAN,
UMR 7039, CNRS,
Polytech Nancy,
Office: C 225,
Université de Lorraine
France.



Table of Contents

- 1 Temporal Difference
- 2 Function Approximation
- 3 RL: Nonlinear Discrete Time
Forward in time Learning
Neural network based approximation

Introduction

- learn online the solutions to optimal control problems without knowing the full system dynamics.
- leads to true online reinforcement learning
- control actions are improved in real time based on estimating their value functions by observing data measured along the system trajectories.
- based on the Bellman equation and solve Policy Evaluation equation by using data observed along a single trajectory of the system.
- TD learning is applicable for feedback control applications

Learning Along System Trajectories

- Temporal difference reinforcement learning methods are based on the Bellman equation and solve equations of PI, without using systems dynamics knowledge, but using data observed along a single trajectory of the system.
- Temporal difference updates the value at each time step as observations of data are made along a trajectory.
- Periodically, the new value is used to update the policy. Temporal difference methods are related to adaptive control in that they adjust values and actions online in real time along system trajectories.

Temporal Difference Equation

Temporal difference reinforcement learning uses one sample path, namely the current system trajectory, to update the value. Value update can be replaced as:

$$V^\pi(x_k) = r_k + \gamma V^\pi(x_{k+1})$$

which holds for each observed data experience set (x_k, x_{k+1}, r_k) at each time stage k . This data set consists of the current state x_k , the observed cost incurred r_k , and the next state x_{k+1} .

TD Equation

The temporal difference error is defined as

$$e_k = -V^\pi(x_k) + r_k + \gamma V^\pi(x_{k+1})$$

and the value estimate is updated to make the temporal difference error small.

- $V^\pi(x_k)$ may be considered as a predicted performance or value,
- r_k as the observed one-step reward,
- $\gamma V^\pi(x_{k+1})$ as a current estimate of future.

TD Equation: Observations

- The Bellman equation can be interpreted as a consistency equation that holds if the current estimate for the predicted value $V^\pi(x_k)$ is correct.
- Temporal difference methods update the predicted value estimate $\hat{V}^\pi(x_k)$ to make the temporal difference error small.
- Idea: TD application on Bellman's equation repeatedly in policy iteration or value iteration, then on average these algorithms converge toward the solution of the stochastic Bellman equation.

Value Function approximation

For nonlinear systems \rightarrow the value function contains higher order nonlinearities. Then, according to the Weierstrass higher order approximation theorem, there exists a dense basis set $\{\varphi_i(x)\}$ such that

$$\begin{aligned} V(x) &= \sum_{i=1}^{\infty} w_i \varphi_i(x) \\ &= \sum_{i=1}^L w_i \varphi_i(x) + \sum_{i=L+1}^{\infty} w_i \varphi_i(x) \equiv W^T \phi(x) + \varepsilon_L(x), \end{aligned}$$

where basis vector $\phi(x) = [\varphi_1(x) \varphi_2(x) \cdots \varphi_L(x)] : R^n \rightarrow R^L$ and $\varepsilon_L(x)$ converges uniformly to zero as the number of terms retained $L \rightarrow \infty$.

NN based approximation of control policy: ACTOR neural Network

Control policy approximation: Actor neural Network

Introducing a second neural network for the control policy, known as the actor neural network. Consider a parametric approximator structure for the control action

$$u_k = \pi(x_k) = U^T \sigma(x_k),$$

with $\sigma(x) : R^n \rightarrow R^M$ being a vector of M activation functions and $U \in R^{M \times m}$ being a matrix of weights or unknown parameters. In the LQR, the optimal state feedback is linear in the states so that the basis set $\sigma(x)$ can be taken as the state vector.

RL: Discrete time optimal control

Assumption: Stabilizable system

System (1) is *stabilizable* on the prescribed set $\Omega \in \mathbb{R}^n$.

\Rightarrow There is a control policy $u_k^1 = \pi(x)$ such that closed loop system $x_{k+1} = f(x_k) + g(x_k)u_k^1$ is asymptotically stable over Ω i.e. $u_k^1 = (u^1(x_k), u^1(x_{k+1}), u^1(x_{k+2}), \dots, u^1(x_\infty))$ exists that

- that stabilises the system (1)
- associated cost $J(x_k, u_k^1)$ is finite.

U denotes the set of all admissible control inputs.

RL: Discrete time optimal control

For a given admissible prescribed policy $\pi(x)$,
the cost associated is called as its *value* denoted as

$$V_{\pi}(x_k) = J(x_k, \pi(x))$$

RL: Discrete time optimal control

Objective: Optimal Cost

To find a control policy $\pi^*(x_k)$ that minimizes the infinite horizon cost function,

$$V^*(x_k) = \min_{u_k \in U} \sum_{n=k}^{\infty} r(x_n, u_n), \forall x_k \quad (4)$$

or, $V^*(x_k) = \min_{\pi(\cdot)} \sum_{n=k}^{\infty} r(x_n, \pi(x_n)), \forall x_k$

Optimal policy

Optimal control policy: $\pi^*(x_k) = \arg \min_{\pi(\cdot)} \sum_{n=k}^{\infty} r(x_n, \pi(x_n)), \forall x_k$

RL: Discrete time optimal control

Cost (given a prescribed policy $u_k = \pi(x_k)$)

$$V_\pi(x_k) = \sum_{n=k}^{\infty} r(x_n, u_n), \forall x_k$$

$$V_\pi(x_k) = r(x_k, u_k) + V_\pi(x_{k+1})$$

Bellman Eq/ Nonlinear

Lyapunov Eq (Recursive):

$$H(x_k, u_k, V_\pi) = r(x_k, u_k) + V_\pi(x_{k+1}) - V_\pi(x_k)$$

Hamiltonian:

Optimal Cost:

$$V^*(x_k) = \min_{u_k \in U} (r(x_k, u_k) + V_\pi(x_{k+1}))$$

RL: Discrete time optimal control

Bellman Principle Bellman, 1957

“An optimal policy has the property that no matter what the previous decisions (i.e. controls) have been, the remaining decisions must constitute an optimal policy with regard to the state resulting from those previous decisions

RL: Discrete time optimal control

Cost (given a prescribed policy $u_k = \pi(x_k)$)

$$V_\pi(x_k) = \sum_{n=k}^{\infty} r(x_n, u_n), \forall x_k$$

$$V_\pi(x_k) = r(x_k, u_k) + V_\pi(x_{k+1})$$

Bellman Eq/ Nonlinear
Lyapunov Eq (Recursive):
Hamiltonian:

$$H(x_k, u_k, V_\pi) = r(x_k, u_k) + V_\pi(x_{k+1}) - V_\pi(x_k)$$

Optimal Cost:

$$V^*(x_k) = \min_{u_k \in U} (r(x_k, u_k) + V_\pi(x_{k+1}))$$

Bellman principle:

$$V^*(x_k) = \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1}))$$

Backwards in Time!!

Optimal control (policy):

$$\pi^*(x_k) = \arg \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1}))$$



UNIVERSITÉ DE LORRAINE



RL: Discrete time optimal control

Cost (given a prescribed policy $u_k = \pi(x_k)$)

$$V_\pi(x_k) = \sum_{n=k}^{\infty} r(x_n, u_n), \forall x_k$$

$$V_\pi(x_k) = r(x_k, u_k) + V_\pi(x_{k+1})$$

Bellman Eq/ Nonlinear
Lyapunov Eq (Recursive):
Hamiltonian:

$$H(x_k, u_k, V_\pi) = r(x_k, u_k) + V_\pi(x_{k+1}) - V_\pi(x_k)$$

Optimal Cost:

$$V^*(x_k) = \min_{u_k \in U} (r(x_k, u_k) + V_\pi(x_{k+1}))$$

Bellman principle:

$$V^*(x_k) = \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1}))$$

Optimal control (policy):

$$\pi^*(x_k) = \arg \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1}))$$

Only data required!!



RL: Discrete time optimal control

Bellman principle:
(DT Hamilton-
Jacobi-Bellman
Equation)

$$\begin{aligned} V^*(x_k) &= \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1})) \\ &= \min_{u_k \in U} (x_k^T Q x_k + u_k^T R u_k + V^*(x_{k+1})) \\ &= \min_{u_k \in U} (x_k^T Q x_k + u_k^T R u_k + V^*(f(x_k) + g(x_k)u_k)) \end{aligned}$$

Optimal control
(policy):

$$\begin{aligned} \pi^*(x_k) &= \arg \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1})) \\ \pi^*(x_k) &= u_k^* = (-1/2)R^{-1}g^T(x_k) \frac{\partial V^*(x_{k+1})}{\partial x_{k+1}} \end{aligned}$$

RL: Discrete time optimal control

Cost (given a prescribed policy $u_k = \pi(x_k)$)

$$V_\pi(x_k) = \sum_{n=k}^{\infty} r(x_n, u_n), \forall x_k$$

Bellman Eq/ Nonlinear
Lyapunov Eq (Recursive):
Optimal Cost:

$$V_\pi(x_k) = r(x_k, u_k) + V_\pi(x_{k+1})$$

$$V^*(x_k) = \min_{u_k \in U} (r(x_k, u_k) + V_\pi(x_{k+1}))$$

Bellman principle:

$$V^*(x_k) = \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1}))$$

Optimal control (policy):

$$\pi^*(x_k) = \arg \min_{u_k \in U} (r(x_k, u_k) + V^*(x_{k+1}))$$

DT Policy Iteration

Initialization

Select any *stabilizing* /admissible control policy: $\pi_j(x_k)$

Policy Evaluation

Determine the *Value* under the current policy using Bellman Equation/Nonlinear Lyapunov Eq.

$$V_{j+1}(x_k) = r(x_k, \pi_j(x_k)) + V_{j+1}(x_{k+1}) ; V_{j+1}(0) = 0$$

Policy Improvement

Determine an improved policy

$$\pi_{j+1}(x_k) = \arg \min_{u_k \in U} (r(x_k, u_k) + V_{j+1}(x_{k+1}))$$

DT Policy Iteration

Initialization

Select any *stabilizing* /admissible control policy: $\pi_j(x_k)$

Policy Evaluation

Determine the *Value* under the current policy using Bellman Equation/Nonlinear Lyapunov Eq.

$$V_{j+1}(x_k) = r(x_k, \pi_j(x_k)) + V_{j+1}(x_{k+1}) ; V_{j+1}(0) = 0$$

Policy Improvement

Determine an improved policy

$$\pi_{j+1}(x_k) = \arg \min_{u_k \in U} (r(x_k, u_k) + V_{j+1}(x_{k+1}))$$

DT Policy Iteration

Initialization

Select any *stabilizing* /admissible control policy: $\pi_j(x_k)$

Policy Evaluation

Determine the *Value* under the current policy using Bellman Equation/Nonlinear Lyapunov Eq.

$$V_{j+1}(x_k) = r(x_k, \pi_j(x_k)) + V_{j+1}(x_{k+1}) ; V_{j+1}(0) = 0$$

Policy Improvement

Determine an improved policy

$$\pi_{j+1}(x_k) = \arg \min_{u_k \in U} (r(x_k, u_k) + V_{j+1}(x_{k+1}))$$

DT Policy Iteration

Initialization

$$\pi_j(x_k)$$

Policy Evaluation

$$V_{j+1}(x_k) = r(x_k, \pi_j(x_k)) + V_{j+1}(x_{k+1})$$

Policy Improvement

$$\pi_{j+1}(x_k) = \arg \min_{u_k \in U} (r(x_k, u_k) + V_{j+1}(x_{k+1}))$$

When $r(x_k, u_k) = x_k^T Q x_k + u_k^T R u_k$,

$$\pi_{j+1}(x_k) = (-1/2)R^{-1}g^T(x_k) \frac{\partial V_{j+1}(x_{k+1})}{\partial x_{k+1}}$$

DT Policy Iteration: Observations

- Initial policy must be stabilizing.
- Policy Iteration (Howard, 1960; Leake and Liu, 1967) \Rightarrow
 - $V_{j+2}(x_k) \leq V_{j+1}(x_k)$
- As $j \rightarrow \infty$:
 - $V_j(x_k) \rightarrow V^*(x_k)$
 - $\pi_j \rightarrow \pi^*$
- Convergence to optimal cost and thus, optimal control policy.

DT Policy Iteration: Observations

- $V_{j+1}(x_k) = r(x_k, \pi_j(x_k)) + V_{j+1}(x_{k+1}); \forall x_k \in \Omega$
 - value of using a given policy starting in all current states possible.
 - Several states \Rightarrow Significant computations!
- Called *full backup* (Sutton and Barto, 2018) \Rightarrow Massive computational load
- Bellman Eq \rightarrow *fixed point equation*
 - Given admissible policy π_j ,
 $V^{i+1}(x_k) = r(x_k, \pi_j(x_k)) + V^i(x_{k+1})$ is a *contraction map*
 - Upon iterated starting from $V^0(x_k)$, $V^i(x_k) \rightarrow V_{j+1}(x_k)$ as $i \rightarrow \infty$.

DT Policy Iteration: Observations

- $V_{j+1}(x_k) = r(x_k, \pi_j(x_k)) + V_{j+1}(x_{k+1}); \forall x_k \in \mathcal{X}$
 - value of using a given policy starting in all current states possible.
 - Several states \Rightarrow Significant computations!
- Called *full backup* (Sutton and Barto, 2018) \Rightarrow Massive computational load
- Bellman Eq \rightarrow *fixed point equation*
 - Given admissible policy π_j ,
 $V^{i+1}(x_k) = r(x_k, \pi_j(x_k)) + V^i(x_{k+1})$ is a *contraction map*
 - Upon iterated starting from $V^0(x_k)$, $V^i(x_k) \rightarrow V_{j+1}(x_k)$ as $i \rightarrow \infty$.

Issues

- This strategy \Rightarrow *backward in time procedure*
- Good for:
 - Off-line planning, Offline optimization, Offline control synthesis.
 - NOT online learning (optimal control synthesis using real time data measured along system trajectories).
- Exact solutions: very difficult
 - Large state space
 - Highly nonlinear dynamics

Issues

- This strategy \Rightarrow *backward in time procedure*
- Good for:
 - Off-line planning, Offline optimization, Offline control synthesis.
 - NOT online learning (optimal control synthesis using real time data measured along system trajectories).
- Exact solutions: very difficult
 - Large state space
 - Highly nonlinear dynamics

Issues

- This strategy \Rightarrow *backward in time procedure*
- Good for:
 - Off-line planning, Offline optimization, Offline control synthesis.
 - NOT online learning (optimal control synthesis using real time data measured along system trajectories.
- Exact solutions: very difficult
 - Large state space
 - Highly nonlinear dynamics

Temporal Difference (TD) or forward in time learning

Value Function approximation (VFA): Neural Networks

Forward-in-time Learning

Temporal Difference Error (TD error):

$$e_k = r(x_k, \pi_{x_k}) + V_{\pi}(x_{k+1}) - V_{\pi}(x_k)$$

- RHS is DT Hamiltonian
- If Bellman Eq holds, e_k is zero.
- Linear in x .
- Thus, given a policy $\pi(x)$, Least Square based solution at each time k for $e_k = 0$.

NN based approximation

Value Function approximation (VFA): Neural Networks

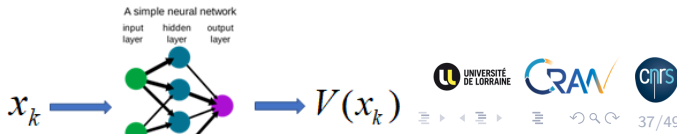
- Value function is sufficiently smooth over compact space
- Consider dense basis set $\{\phi_i(x)\}$ with basis vector (Weierstrass Theorem):

$$\phi(x) = [\varphi_1(x)\varphi_2(x)\dots\varphi_L(x)] : \mathbb{R}^n \rightarrow \mathbb{R}^L$$

$$V_\pi(x) = \sum_{i=1}^L w_i \varphi_i(x) = W^T \phi(x)$$

Substituting in Bellman TD equation:

$$e_k = r(x_k, \pi_{x_k}) + W^T \phi(x_{k+1}) - W^T \phi(x_k)$$



Online DT Policy Iteration: Observations

- At $k + 1$: observe $x_k, u_k = \pi_j(x_k), x_{k+1}$
- Calculate $r(x_k, u_k)$ One scalar Equation in
 $r(x_k, \pi_{x_k}) = W_{j+1}^T (\phi(x_k) - \phi(x_{k+1}))$
- Use same policy $u_k = \pi_j(x_k)$, collect L data $\Rightarrow L$ equations
(!! $\phi(x) = \mathbb{R}^n \rightarrow \mathbb{R}^L$).
- Determine LS based solution \hat{W}_{j+1}
- Repeat till $\hat{W}_{j+1} \equiv \hat{W}_{j+2} \rightarrow W^*$ Apply Improved control

Batch Least Square Approach

- Can be implemented online by standard system identification techniques.
- Note that weight update is a scalar equation, whereas the unknown parameter vector $W_{j+1} \in R^L$ has L elements. Therefore, data from multiple time steps are needed for its solution.
- At time $k + 1$ we measure the previous state x_k , the control $u_k = \pi_j(x_k)$, the next state x_{k+1} , and compute the resulting utility $r(x_k, \pi_j(x_k))$.
- These data result in one scalar equation.
- This procedure is repeated for subsequent times using the same policy $\pi_j(\cdot)$ until at least L equations are obtained, at which point the least-squares solution W_{j+1} can be determined. Batch leastsquares can be used for this procedure.

Recursive Least square approach

- Express Critic Weight update as:

$$W_{j+1}^T \Phi(k) \equiv W_{j+1}^T (\phi(x_k) - \gamma \phi(x_{k+1})) = r(x_k, \pi_j(x_k))$$

with $\Phi(k) \equiv (\phi(x_k) - \gamma \phi(x_{k+1}))$ being a regression vector.

- At step j of the policy iteration algorithm, the control policy is fixed at $u = \pi_j(x)$.
- Then, at each time k the data set $(x_k, x_{k+1}, r(x_k, \pi_j(x_k)))$ is measured.
- One step of RLS is then performed.
- This procedure is repeated for subsequent times until convergence to the parameters corresponding to the value $V_{j+1}(x) = W_{j+1}^T \phi(x)$.
- For RLS to converge, the regression vector $\Phi(k) \equiv (\phi(x_k) - \gamma \phi(x_{k+1}))$ must be persistently exciting.

Execution: Adaptive Critic Structures

RL: Model free approach → Q-function

System $x_{k+1} = f(x_k) + g(x_k)u_k$

$$V_h(0) = 0$$

$$Q_h(x_k, \underline{u}_k) = r(x_k, \underline{u}_k) + \gamma V_h(x_{k+1})$$

$$Q_h(x_k, u_k) = r(x_k, u_k) + \gamma Q_h(x_{k+1}, h(x_{k+1}))$$

$$V^*(x_k) = \min_{u_k} (Q^*(x_k, u_k))$$

$$h^*(x_k) = \arg \min_{u_k} (Q^*(x_k, u_k))$$


~~$$u^*(x_k) = -\frac{1}{2} R^{-1} g(x_k)^T \frac{dV^*(x_{k+1})}{dx_{k+1}}$$~~



References I

Lewis, F. L., Vrabie, D., & Vamvoudakis, K. G. (2012). Reinforcement learning and feedback control: Using natural decision methods to design optimal adaptive controllers. *IEEE Control Systems Magazine*, 32(6), 76-105.

 Bellman, R. (1957). A markovian decision process. *Journal of mathematics and mechanics*, 679–684.

 Howard, R. A. (1960). Dynamic programming and markov processes..

 Leake, R., & Liu, R.-W. (1967). Construction of suboptimal control sequences. *SIAM Journal on Control*, 5(1), 54–63.

 Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.