

Reinforcement Learning

An Introductory Note

Jingye Wang

✉ wangjy5@shanghaitech.edu.cn

Spring 2020

Contents

1	Introduction	3
2	Review of Basic Probability	5
2.1	Interpretation of Probability	5
2.2	Transformations	5
2.3	Limit Theorem	5
2.4	Sampling & Monte Carlo Methods	6
2.5	Basic Inequalities	8
2.6	Concentration Inequalities	10
2.7	Conditional Expectation	12
3	Bandit Algorithms	14
3.1	Bandit Models	14
3.2	Stochastic Bandits	14
3.3	Greedy Algorithms	15
3.4	UCB Algorithms	16
3.5	Thompson Sampling Algorithms	17
3.6	Gradient Bandit Algorithms	18
4	Markov Chains	20
4.1	Markov Model	20
4.2	Basic Computations	20
4.3	Classifications	21

CONTENTS	2
4.4 Stationary Distribution	22
4.5 Reversibility	22
4.6 Markov Chain Monte Carlo	23
5 Markov Decision Process	25
5.1 Markov Reward Process	25
5.2 Markov Decision Process	26
5.3 Dynamic Programming	28
6 Model-Free Prediction	33
6.1 Monte-Carlo Policy Evaluation	33
6.2 Temporal-Difference Learning	35
7 Model-Free Control	37
7.1 On Policy Monte-Carlo Control	37
7.2 On Policy Temporal-Difference Control: Sarsa	39
7.3 Off-Policy Temporal-Difference Control: Q-Learning	40
8 Value Function Approximation	41
8.1 Semi-gradient Method	41
8.2 Deep Q-Learning	43
9 Policy Optimization	46
9.1 Policy Optimization Theorem	46
9.2 REINFORCE: Monte-Carlo Policy Gradient	49
9.3 Actor-Critic Policy Gradient	51
9.4 Extension of Policy Gradient	52

8 Value Function Approximation

For some problems, the state spaces may be huge or even infinity, e.g.:

- Backgammon: 10^{20} states;
- Chess: 10^{47} states;
- Game of Go: 10^{170} states;
- Robot Arm and Helicopter: continuous state space;

How can we scale up the model-free methods for these problems?

A solution for these is *function approximation*, which is an instance of *supervised learning*. Concretely, it follows that:

- $\hat{v}(s; \mathbf{w}) \approx v^\pi(s)$,
- $\hat{q}(s, a; \mathbf{w}) \approx q^\pi(s, a)$,
- $\hat{\pi}(s, a; \mathbf{w}) \approx \pi(a|s)$,

where \mathbf{w} is the parameter to be learned. With such approximation, we can generalize from seen states to unseen states. Obviously, we have many available function approximators:

- Linear combinations of features,
- Neural networks,
- Decision trees,
- Nearest neighbors.

Among those, we will focus on the first two as they are differentiable and thus easy to be optimized. It is worth to notice that though such generalization makes the learning potentially more powerful (it may even be applicable to partially observable problems), it is potentially more difficult to manage and understand.

8.1 Semi-gradient Method

Now we consider a naive case where we have an oracle for knowing the true value $v^\pi(s)$ for any given state s under policy π . Then the object is to find the estimator $\hat{v}(s, \mathbf{w})$ approximating $v^\pi(s)$. We can define the loss function by the mean squared error:

$$J(\mathbf{w}) = \mathbb{E}_\pi[(v^\pi(s) - \hat{v}(s, \mathbf{w}))^2],$$

and the gradient descent for the loss function is:

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w}),$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \Delta \mathbf{w}.$$

Following the gradient descent, we can find a local minimum. Further, if the value function is represented by a linear combination of features, then such method can converge to the global optimum.

8.1.1 Semi-gradient Method for Prediction

However, in practice, no access to the oracle of the true value $v^\pi(s)$. What we have is only the reward, or, the bootstrapping estimation. Thus we can substitute the target for $v^\pi(s)$:

- For MC, the target is the actual return G_t and hence

$$\Delta \mathbf{w} = \alpha(G_t - \hat{v}(s_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s_t, \mathbf{w});$$

- For TD(0), the target is the TD target $R_{t+1} + \gamma \hat{v}(s_{t+1}, \mathbf{w})$ and hence

$$\Delta \mathbf{w} = \alpha(r_{t+1} + \gamma \hat{v}(s_{t+1}, \mathbf{w}) - \hat{v}(s_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s_t, \mathbf{w})$$

For TD(0), the gradient is also called as semi-gradient, as the target $r_{t+1} + \gamma \hat{v}(s_{t+1}, \mathbf{w})$ depends on the current value of the weight vector \mathbf{w} and we just ignore such effect.

Algorithm 16 Gradient Monte Carlo Policy Evaluation

```

1: initialize the differentiable function  $\hat{v}(s, \mathbf{w}) \in \mathbb{R}$  arbitrarily for all  $s \in \mathcal{S}$ ;
2: input the policy  $\pi$  to be evaluated; the step size  $\alpha$ ;
3: for true do:
    # variants for this alg. can be start with  $s_0 \in \mathcal{S}, a_0 \in \mathcal{A}(s_0)$  randomly,  $\varepsilon$ -greedy, etc.
4:   Generate a complete episode  $\tau = (s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T)$ ;
5:   for  $t = 0, 1, \dots, T-1$  do:
6:      $\mathbf{w} \leftarrow \mathbf{w} + \alpha[G_t - \hat{v}(s_t, \mathbf{w})] \nabla \hat{v}(s_t, \mathbf{w})$ ;
7:   end for
8: end for

```

Algorithm 17 Semi-gradient TD(0) Policy Evaluation

```

1: initialize the differentiable function  $\hat{v}(s, \mathbf{w}) \in \mathbb{R}$  arbitrarily for all  $s \in \mathcal{S}$ ;
2: input the policy  $\pi$  to be evaluated; the step size  $\alpha$ ;
3: for true do:
4:   Generate a start state  $s$ ;
5:   while  $s$  is not terminal do:
6:     Choose action  $a \leftarrow \pi(s)$ ;
7:      $r, s' \leftarrow \text{environment}(s, a)$ ;
8:      $\mathbf{w} \leftarrow \mathbf{w} + \alpha[r + \gamma \hat{v}(s', \mathbf{w}) - \hat{v}(s, \mathbf{w})] \nabla \hat{v}(s, \mathbf{w})$ ;
9:      $s \leftarrow s'$ ;
10:  end while
11: end for

```

8.1.2 Semi-gradient Method for Control

The extension of the semi-gradient prediction methods to state-action values is straightforward. In this case it is the approximate action-value function, $\hat{q} \approx q^\pi$, that is represented as a parameterized functional

form with weight vector \mathbf{w} . The control version is derived from the *generalized policy iteration*. It quite like the way we make prediction:

- For MC, the target is return G_t

$$\Delta \mathbf{w} = \alpha (G_t - \hat{q}(s_t, a_t, \mathbf{w})) \nabla \hat{q}(s_t, a_t, \mathbf{w});$$

- For Sarsa, the target is TD target $r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, \mathbf{w})$:

$$\Delta \mathbf{w} = \alpha (r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, \mathbf{w}) - \hat{q}(s_t, a_t, \mathbf{w})) \nabla \hat{q}(s_t, a_t, \mathbf{w});$$

- For Q-learning, the target is TD target $r_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a, \mathbf{w})$:

$$\Delta \mathbf{w} = \alpha \left(r_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a, \mathbf{w}) - \hat{q}(s_t, a_t, \mathbf{w}) \right) \nabla \hat{q}(s_t, a_t, \mathbf{w}).$$

One can get the corresponding control algorithm by modifying the update rule in the prediction version and then extracting the optimal policy according to \hat{q} , which is $\pi(s) = \arg \max_a \hat{q}(s, a, \mathbf{w})$.

Though TD control is powerful and practical in most case, there is no guarantees in the convergence of TD. Firstly, TD with VFA follows the *semi-gradient* rather than the true gradient of any objective function. Secondly, the updates involve doing an approximate Bellman backup (model-free) followed by fitting the underlying value function (approximation). Further, for off-policy, behavior policy and target policy are not identical, thus value function approximation may diverge.

The Deadly Triad for *Danger of Instability and Divergence*:

- Function approximation: A scalable way of generalizing from a state space much larger than the memory and computational resources;
- Bootstrapping: Update target that includes existing estimates (as in dynamic programming or TD methods) rather than relying exclusively on actual rewards and complete returns (as in MC methods);
- Off-policy learning: training on a distribution of transitions other than that produced by the target policy.

8.2 Deep Q-Learning

Deep Q-learning leverages nonlinear function approximator, deep neural networks, to avoid manual designing of features. Such method is also called DQN, and it is a well-known case of *deep reinforcement learning*.

Deep Reinforcement Learning:

- Frontier in machine learning and artificial intelligence;
- Deep neural networks can be used to represent value function, policy function, and even model;
- Optimize loss function by stochastic gradient descent.

Algorithm 18 Semi-gradient TD(0) Policy Evaluation

```

1: initialize the differentiable function  $\hat{v}(s, \mathbf{w}) \in \mathbb{R}$  arbitrarily for all  $s \in \mathcal{S}$ ;
2: input the policy  $\pi$  to be evaluated; the step size  $\alpha$ ;
3: for true do:
4:   Generate a start state  $s$ ;
5:   while  $s$  is not terminal do:
6:     Choose action  $a \leftarrow \pi(s)$ ;
7:      $r, s' \leftarrow \text{environment}(s, a)$ ;
8:      $\mathbf{w} \leftarrow \mathbf{w} + \alpha[r + \gamma \hat{v}(s', \mathbf{w}) - \hat{v}(s, \mathbf{w})] \nabla \hat{v}(s, \mathbf{w})$ ;
9:      $s \leftarrow s'$ ;
10:  end while
11: end for

```

The success of DQN lies in two innovative techniques that greatly improve and stabilize the training procedure.

Experience Replay

- To reduce the correlations among samples, DQN stores transition (s_t, a_t, r_t, s_{t+1}) in a replay memory \mathcal{D} ;
- To perform experience replay, DQN repeats the following
 - samples an experience tuple from the dataset:

$$(s, a, r, s') \sim \mathcal{D};$$

- computes the target value with the sampled tuple: $r + \gamma \max_{a'} \hat{q}(s', a', \mathbf{w})$;
- uses stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha(r + \gamma \max_{a'} \hat{q}(s', a', \mathbf{w}) - q(s, a, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(s, a, \mathbf{w}).$$

Fixed Target

- To help improve stability, DQN fixes the weights used in the target calculation for multiple updates;
- Let a different set of parameter \mathbf{w}^- be the set of weights used in the target, and \mathbf{w} be the weights to be updated in each step;
- To perform experience replay with fixed target, DQN repeats the following
 - samples an experience tuple from the dataset: $(s, a, r, s') \sim \mathcal{D}$;
 - computes the target value for the sampled tuple:

$$r + \gamma \max_{a'} \hat{q}(s', a', \mathbf{w}^-);$$

- uses stochastic gradient descent to update the network weights:

$$\Delta \mathbf{w} = \alpha(r + \gamma \max_{a'} \hat{q}(s', a', \mathbf{w}^-) - q(s, a, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(s, a, \mathbf{w});$$

- updates $\mathbf{w}^- \leftarrow \mathbf{w}$ periodically.

Algorithm 19 Deep-Q Network

- 1: initialize the replay memory \mathcal{D} with capacity N ; initialize the state-action value function q with random weight \mathbf{w} ; initialize the target state-action value function \hat{q} with random weight $\mathbf{w}^- = \mathbf{w}$;
 - 2: **input** the period C ;
 - 3: **for** true **do**:
 - 4: Generate a start state s ;
 - 5: **while** s is not terminal **do**:
 - 6: Choose action $a = \begin{cases} \arg \max_{a' \in \mathcal{A}} q(s, a'), & \text{with probability } \varepsilon; \\ \text{a random action,} & \text{otherwise;} \end{cases}$
 - 7: $r, s' \leftarrow \text{environment}(s, a)$;
 - 8: Store transition (s, a, r, s') in \mathcal{D} ;
 - 9: Sample random minibatch of transitions (s_j, a_j, r_j, s_{j+1}) from \mathcal{D} ;
 - 10: Set $y_j = \begin{cases} r_j, & \text{if } s_{j+1} \text{ is the terminal state;} \\ r_j + \gamma \max_{a'} \hat{q}(s_{j+1}, a', \mathbf{w}^-), & \text{otherwise;} \end{cases}$
 - 11: Perform a gradient descent step on $(y_j - q(s_j, a_j, \mathbf{w}))^2$ with respect to the parameter \mathbf{w} ;
 - 12: Reset $\hat{q} = q$ every C steps;
 - 13: **end while**
 - 14: **end for**
-