

Reinforcement Learning

Lecture 7: Value Function Approximation

Lecturer: Haim Permuter

Scribe: Nave Algarici

I. INTRODUCTION

II. VALUE FUNCTION APPROXIMATION

In previous lectures, we have represented the value function using a lookup table, where every state s has an entry $V(s)$, or alternatively, every state-action pair s, a has an entry $Q(s, a)$. The problem comes when we start dealing with large MDPs. In such cases, there are too many states and/or actions to store in memory, and the process of learning each state's value individually takes too long. The solution in this case is to estimate the value function with function approximation:

$$\hat{v}(s, w) \approx v_{\pi}(s). \quad (1)$$

where w is a trainable parameter vector. Using this approximation we can generalize from seen states to unseen states, and update parameter vector w to improve the accuracy of the function approximation.

The function approximator can be one of many options, including: Linear combinations of features, Neural Network, Decision Tree, Nearest Neighbour, etc. We will consider differentiable function approximators: Linear combinations of features and Neural Networks.

A. Value Function Approximation By Stochastic Gradient Descent

Reminder 1 (Gradient Descent) Let $J(w)$ be a differentiable function of parameter vector w . Define the gradient of $J(w)$ to be

$$\nabla_w J(w) = \begin{pmatrix} \frac{\partial J(w)}{\partial w_1} \\ \vdots \\ \frac{\partial J(w)}{\partial w_n} \end{pmatrix}. \quad (2)$$

To find a local minimum of $J(w)$, adjust w in the opposite direction of the gradient

$$\Delta w = -\frac{1}{2}\alpha\nabla_w J(w). \quad (3)$$

Our goal is to find a parameter vector w that minimizes the mean-squared error between the approximate value function $\hat{v}(s, w)$ and the true value function $v_\pi(s)$. The objective function:

$$J(w) = \mathbb{E}_\pi [(v_\pi(S) - \hat{v}(S, w))^2]. \quad (4)$$

Gradient descent finds a local minimum

$$\Delta w = -\frac{1}{2}\alpha\nabla_w J(w) \quad (5)$$

$$= \alpha\mathbb{E}_\pi [(v_\pi(S) - \hat{v}(S, w))\nabla_w \hat{v}(S, w)], \quad (6)$$

and Stochastic gradient descent (SGD) samples the gradient

$$\Delta w = \alpha(v_\pi(S) - \hat{v}(S, w))\nabla_w \hat{v}(S, w). \quad (7)$$

The expected update of SGD is equal to the full gradient update.

Each state will be represented by a feature vector

$$x(S) = \begin{pmatrix} x_1(S) \\ \vdots \\ x_n(S) \end{pmatrix}. \quad (8)$$

This feature vector can many things, such as: distance of a robot from landmarks, Chess pieces configuration on a board, etc.

Example 1 (Linear Value function Approximation) In this case, the value function is represented by a linear combination of features

$$\hat{v}(S, w) = x(S)^T w = \sum_{j=1}^n x_j(S)w_j, \quad (9)$$

and the objective function is quadratic in parameters w

$$J(w) = \mathbb{E}_\pi [(v_\pi(S) - x(S)^T w)^2]. \quad (10)$$

By deriving the objective function by w

$$\nabla_w \hat{v}(S, w) = x(S), \quad (11)$$

we get the update rule

$$\Delta w = \alpha(v_\pi(S) - x(S)^T w)x(S). \quad (12)$$

Note 1 Table lookup is a special case of linear value function approximation, where the table lookup features are

$$x^{table}(S) = \begin{pmatrix} \mathbf{1}(S = s_1) \\ \vdots \\ \mathbf{1}(S = s_n) \end{pmatrix}, \quad (13)$$

and the parameter vector w is the value of each individual state

$$\hat{v}(S, w) = \left(\mathbf{1}(S = s_1), \dots, \mathbf{1}(S = s_n) \right) \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix}. \quad (14)$$

B. Incremental Prediction Algorithms

Up until this point, we have assumed the true value function $v_\pi(s)$ is given, but in real world problems, we observe only rewards. In practise, we substitute a *target* for $v_\pi(s)$ using different methods. We will discuss four of them.

- For Monte Carlo (MC), the *target* is the return G_t

$$\Delta w = \alpha(G_t - \hat{v}(S_t, w))\nabla_w \hat{v}(S_t, w). \quad (15)$$

- For TD(0), the *target* is the TD target $R_{t+1} + \gamma\hat{v}(S_{t+1}, w)$

$$\Delta w = \alpha(R_{t+1} + \gamma\hat{v}(S_{t+1}, w) - \hat{v}(S_t, w))\nabla_w \hat{v}(S_t, w). \quad (16)$$

- For forward-view TD(λ), the *target* is the *lambda*-return G_t^λ

$$\Delta w = \alpha(G_t^\lambda - \hat{v}(S_t, w))\nabla_w \hat{v}(S_t, w). \quad (17)$$

- For backward-view TD(λ). the equivalent update is

$$\delta_t = R_{t+1} + \gamma\hat{v}(S_{t+1}, w) - \hat{v}(S_t, w), \quad (18)$$

$$E_t = \gamma\lambda E_{t-1} + \nabla_w \hat{v}(S_t, w), \quad (19)$$

$$\Delta w = \alpha\delta_t E_t. \quad (20)$$

III. ACTION-VALUE FUNCTION APPROXIMATION

In lecture 5, we learned controlling the policy by repeating two steps: policy evaluation and policy improvement. Action-value function approximation is used to perform an approximate policy evaluation, for the same reasons that were mentioned in the introduction. After that, an ϵ -greedy policy improvement can be applied using this approximation. The approximation process of the Action-value function is similar to the value function. The Action-value function approximation:

$$\hat{q}(s, a, w) \approx q_\pi(s, a). \quad (21)$$

The objective function:

$$J(w) = \mathbb{E}_\pi [(q_\pi(S, A) - \hat{q}(S, A, w))^2]. \quad (22)$$

Using SGD we can find a local minimum

$$\Delta w = -\frac{1}{2}\alpha \nabla_w J(w) \quad (23)$$

$$= \alpha(q_\pi(S, A) - \hat{q}(S, A, w))\nabla_w \hat{q}(S, A, w). \quad (24)$$

Each state-action pair will be represented by a feature vector

$$x(S) = \begin{pmatrix} x_1(S, A) \\ \vdots \\ x_n(S, A) \end{pmatrix}. \quad (25)$$

Example 2 (Linear Action-Value function Approximation) In this case, the action-value function is represented by a linear combination of features

$$\hat{q}(S, A, w) = x(S, A)^T w = \sum_{j=1}^n x_j(S, A)w_j, \quad (26)$$

The SGD update rule:

$$\nabla_w \hat{q}(S, A, w) = x(S, A), \quad (27)$$

$$\Delta w = \alpha(q_\pi(S, A) - x(S, A)^T w)x(S, A). \quad (28)$$

A. Incremental Control Algorithms

As in prediction, we will substitute a *target* for $q_\pi(s, a)$.

- For Monte Carlo (MC), the *target* is the return G_t

$$\Delta w = \alpha(G_t - \hat{q}(S_t, A_t, w)) \nabla_w \hat{q}(S_t, A_t, w). \quad (29)$$

- For TD(0), the *target* is the TD target $R_{t+1} + \gamma \hat{v}(S_{t+1}, A_{t+1}, w)$

$$\Delta w = \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, w) - \hat{q}(S_t, A_t, w)) \nabla_w \hat{q}(S_t, A_t, w). \quad (30)$$

- For forward-view TD(λ), the *target* is the *lambda*-return G_t^λ

$$\Delta w = \alpha(G_t^\lambda - \hat{q}(S_t, A_t, w)) \nabla_w \hat{q}(S_t, A_t, w). \quad (31)$$

- For backward-view TD(λ), the equivalent update is

$$\delta_t = R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, w) - \hat{q}(S_t, A_t, w), \quad (32)$$

$$E_t = \gamma \lambda E_{t-1} + \nabla_w \hat{q}(S_t, A_t, w), \quad (33)$$

$$\Delta w = \alpha \delta_t E_t. \quad (34)$$

IV. BATCH REINFORCEMENT LEARNING

This method addresses the objective of best fitting the value function to the whole training data. Instead of using only the last data point, or *event*, to improve the accuracy of the value function, it can use all of the agent's experience as training data, and improve using it. This method gives us more stability and is more data efficient.

A. Stochastic Gradient Descent with Experience Replay

In this algorithm we sample randomly a state-value pair from the complete experience, and use it for training the value function approximation.

Algorithm 1 *SGD Experience Replay*

Given $\mathcal{D} = \{\langle s_1, v_1^\pi \rangle, \langle s_2, v_2^\pi \rangle, \dots, \langle s_T, v_T^\pi \rangle\}$ **repeat**

Sample state, value from experience

$$\langle s, v^\pi \rangle \sim \mathcal{D}$$

Apply SGD update

$$\Delta w = \alpha(v_\pi - \hat{v}(s, w)) \nabla_w \hat{v}(s, w)$$

REFERENCES

- [1] UCL course on RL by David Silver, Lecture 6, Value Function Approximation.
http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching_files/FA.pdf.