

Lecture 15 - Value Function Approximation

Guiliang Liu

The Chinese University of Hong Kong, Shenzhen

DDA4230: Reinforcement Learning
Course Page: [\[Click\]](#)

Value Function Approximation

Motivations:

- So far we represented the value function by a **lookup table** where each state has an entry, $V(s)$, or each state-action pair has an entry, $Q(s, a)$.
- However, this approach might not generalize well to problems with large state and action spaces. A popular solution is via **value function approximation (VFA)**

$$V^\pi(s) \approx \hat{V}(s, w) \quad \text{or} \quad Q^\pi(s, a) \approx \hat{Q}(s, a, w).$$



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Value Function Approximation

$$V^\pi(s) \approx \hat{V}(s, w) \quad \text{or} \quad Q^\pi(s, a) \approx \hat{Q}(s, a, w).$$

In the approximation, w is usually referred to as **the parameter or weights of our function approximator**. Some choices for function approximators are listed below.

- Linear combinations of features
- Neural networks
- Decision trees
- Nearest neighbors
- Fourier and wavelet basis



Linear feature representations

In **linear function representations**, we use a feature vector to represent a state

$$x(s) = (x_1(s), x_2(s), \dots, x_d(s))^T,$$

where d is the **dimensionality of the feature space**. We then approximate our value functions using a linear combination of features as

$$\hat{V}(s, \mathbf{w}) = x(s)^T \mathbf{w} = \sum_{j=1}^d x_j(s) w_j.$$



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Linear feature representations

The error of the approximation is defined on **the measure space of the occupancy measure**, which denotes the cumulative probability that a state is visited under π

$$\rho^\pi(s) = \lim_{T \rightarrow \infty} \frac{\sum_{t=0}^T \gamma^t \mathbb{P}(s_t = s \mid \pi)}{\sum_{t=0}^T \gamma^t}.$$

The **quadratic objective function** (also known as the loss function) of the approximation error is then defined as

$$J(\mathbf{w}) = \mathbb{E}_{s \sim \rho^\pi(s)} \left[(V^\pi(s) - \hat{V}(s, \mathbf{w}))^2 \right].$$

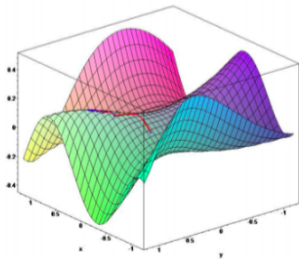
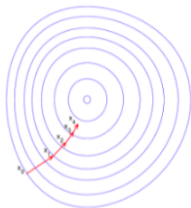


香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Gradient Descent

A common technique to minimize the above objective function is gradient descent.



- Start at some particular spot x_0 , corresponding to some initial value of our parameter w .
- Evaluate the gradient at x_0 , which is the direction of the steepest increase of objective.
- Take a step along the negative direction of the gradient vector and arrive at x_1 .
- This process is repeated until convergence.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Gradient Descent

Mathematically, this can be summarized as

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \left(\frac{\partial J(\mathbf{w})}{\partial w_1}, \frac{\partial J(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_n} \right),$$
 compute the gradient

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

compute an update step using gradient descent

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w},$$

take a step towards the local minimum

where α is the learning rate.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Gradient Descent

Stochastic Gradient Descent

- In practice, gradient descent is not a sample-efficient optimizer. so we use stochastic gradient descent (SGD).
- In minibatch SGD, we sample a minibatch of past experiences, compute our objective function on that minibatch, and update our parameters using gradient descent on the minibatch.



Monte-Carlo policy evaluation with linear VFA

Algorithm 1: Monte-Carlo policy evaluation with linear VFA

Initialize $\mathbf{w} = 0$, $R(s) = 0 \forall s$, $k = 1$

while *true* **do**

 Sample k -th episode $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \dots, s_{k,H_k})$ given π

for $t = 1, \dots, H_k$ **do**

if *first visit to s in episode k* **then**

 Append $\sum_{j=t}^{H_k} r_{k,j}$ to $R(s_t)$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha(\text{avg}(R(s_t)) - \hat{V}(s_t, \mathbf{w}))x(s_t)$

$k = k + 1$



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Monte-Carlo policy evaluation with linear VFA

This algorithm is a modification of first-visit Monte-Carlo policy evaluation, while we replace our value function with our linear VFA. Recall that the mean squared error of a linear VFA for a particular policy π relative to the true value is:

$$J(\mathbf{w}) = \sum_s \rho^\pi(s) (V^\pi(s) - \hat{V}^\pi(s, \mathbf{w}))^2.$$

Lemma

Monte-Carlo policy evaluation with linear VFA converges to the weights \mathbf{w}_{MC} with minimum mean squared error.

$$J(\mathbf{w}_{MC}) = \min_{\mathbf{w}} \sum_s \rho^\pi(s) (V^\pi(s) - \hat{V}^\pi(s, \mathbf{w}))^2.$$



Temporal-difference methods with linear VFA

Recall that in the tabular setting, we approximate V^π via bootstrapping and sampling and update $V^\pi(s)$ by

$$V^\pi(s) \leftarrow V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s)),$$

where $r + \gamma V^\pi(s')$ represents our TD target. Using linear VFA, we replace V^π with \hat{V}^π and our update equation becomes

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha(r + \gamma \hat{V}^\pi(s', \mathbf{w}) - \hat{V}^\pi(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{V}^\pi(s, \mathbf{w}) \\ &= \mathbf{w} + \alpha(r + \gamma \hat{V}^\pi(s', \mathbf{w}) - \hat{V}^\pi(s, \mathbf{w})) \mathbf{x}(s). \end{aligned}$$



Temporal-difference methods with linear VFA

In value function approximation, although our target is a biased and approximated estimate of the true value $V^\pi(s)$, linear TD(0) will still converge to some global approximate optimum.

Lemma

TD(0) policy evaluation with VFA converges to the weights w_{TD} which is optimum up to $1/(1-\gamma)$ of the minimum mean squared error.

$$J(w_{TD}) \leq \frac{1}{1-\gamma} \min_w \sum_s \rho^\pi(s) (V^\pi(s) - \hat{V}^\pi(s, w))^2.$$



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Control using VFA

Use function approximators for **action-value functions** by $\hat{Q}(s, a, w) \approx Q^\pi(s, a)$. We may then interleave policy evaluation, by approximating using $\hat{Q}(s, a, w)$, and policy improvement, by ϵ -greedy policy improvement. To be more concrete, we define our objective function $J(w)$ as

$$J(w) = \mathbb{E}_\pi \left[(Q^\pi(s, a) - \hat{Q}^\pi(s, a, w))^2 \right].$$

Similar to what we did earlier in policy evaluation, we may then use either gradient descent or stochastic gradient descent to minimize the objective function.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Control using VFA

For example, for a linear action value function approximator, this can be summarized as

$$x(s, a) = (x_1(s, a), x_2(s, a), \dots, x_n(s, a))^T,$$

$$\hat{Q}(s, a, w) = x(s, a)^T w,$$

$$J(w) = \mathbb{E}_\pi \left[(Q^\pi(s, a) - \hat{Q}^\pi(s, a, w))^2 \right],$$

$$\begin{aligned} -\frac{1}{2} \nabla_w J(w) &= \mathbb{E}_\pi \left[(Q^\pi(s, a) - \hat{Q}^\pi(s, a, w)) \nabla_w \hat{Q}^\pi(s, a, w) \right] \\ &= \mathbb{E}_\pi \left[(Q^\pi(s, a) - \hat{Q}^\pi(s, a, w)) x(s, a) \right], \end{aligned}$$

$$\begin{aligned} \Delta w &= -\frac{1}{2} \alpha \nabla_w J(w) \\ &= \alpha (Q^\pi(s, a) - \hat{Q}^\pi(s, a, w)) x(s, a), \end{aligned}$$

$$w \leftarrow w + \Delta w.$$

action value features

action value linear in features

objective function

compute the gradient

 香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen

compute the update

take a step of gradient descent 

Control using VFA

For Monte Carlo methods, we substitute our target $Q^\pi(s, a)$ with a return G_t .

$$\Delta \mathbf{w} = \alpha(G_t - \hat{Q}(s, a, \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a, \mathbf{w}).$$

For SARSA, we substitute our target with a TD target

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

For Q-learning, we substitute our target with a maximum TD target

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



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Neural networks

Motivations:

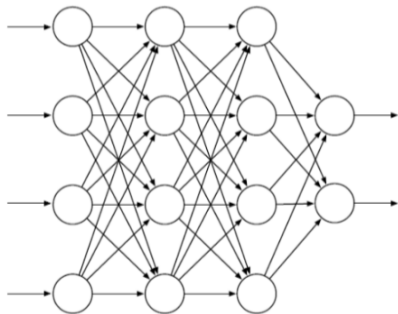
- Although linear VFAs often work well given the right set of features, it can also be difficult to hand-craft such feature set.
- Neural networks provide a much richer function approximation class that is able to directly go from states without requiring an explicit specification of features.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Neural networks



- Neural networks with a single hidden layer can have the “universal approximation” property, which has been demonstrated both empirically and theoretically.
- Complicated functions can be approximated with a hierarchical composition of multiple hidden layers.



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen

Question and Answering (Q&A)



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen