

RL: Lecture 21— Approximation

Harvey Mudd College

April 15, 2020

Neil Rhodes

Overview

- Why approximation?
- Value-function Approximation
- Prediction Objective
- Linear methods
- Stochastic Gradient Descent (SGD) methods

Why approximation?

Value-function Approximation

weight vector $\mathbf{w} \in \mathbb{R}^d$

approximate value function: $\hat{v}(\mathbf{s}, \mathbf{w}) \approx v_\pi(\mathbf{s})$

$$\mathbf{s} \mapsto u$$

\mathbf{s} : state being updated

u : update target \mathbf{s} 's value is being shifted toward

Prediction Objective

How much we care about the error at each state: state distribution $\mu(\mathbf{s}) \geq 0, \sum_{\mathbf{s}} \mu(\mathbf{s}) = 1$

Mean Squared Value Error: $\overline{VE}(\mathbf{w}) = \sum_{\mathbf{s} \in \mathcal{S}} \mu(\mathbf{s}) [v_{\pi}(\mathbf{s}) - \hat{v}(\mathbf{s}, \mathbf{w})]^2$

On-policy distribution in episodic tasks

$\mu(\mathbf{s})$: fraction of the time spent in \mathbf{s} under policy π .

Probability starting in state \mathbf{s} : $h(\mathbf{s})$

time spent in \mathbf{s} : $\eta(\mathbf{s}) = h(\mathbf{s}) + \sum_{\bar{\mathbf{s}}} \eta(\bar{\mathbf{s}}) \sum_a \pi(a|\bar{\mathbf{s}}) p(\mathbf{s}|\bar{\mathbf{s}}, a)$

Normalized to sum to one: $\mu(\mathbf{s}) = \frac{\eta(\mathbf{s})}{\sum_{\mathbf{s}'} \eta(\mathbf{s}')}$

Goal for \overline{VE}

Find a global optimum:

$$\boldsymbol{w}^* \text{ such that } \forall \boldsymbol{w} : \overline{VE}(\boldsymbol{w}^*) \leq \overline{VE}(\boldsymbol{w})$$

May have to settle for local optimum:

$$\boldsymbol{w}^* \text{ such that } \forall \boldsymbol{w} \text{ in neighborhood of } \boldsymbol{w}^* : \overline{VE}(\boldsymbol{w}^*) \leq \overline{VE}(\boldsymbol{w})$$

Gradient

Gradient: operation that takes a function of a vector and creates a vector of partial derivatives:

$$\nabla f(\mathbf{w}) = \left[\frac{\partial f(\mathbf{w})}{\partial w_1}, \frac{\partial f(\mathbf{w})}{\partial w_2}, \dots, \frac{\partial f(\mathbf{w})}{\partial w_d} \right]^T$$

SGD

Given loss function $\overline{VE}(\mathbf{w})$ which consists of a sum of individual losses: $\sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$

1. Pick an $s \in \mathcal{S}$ according to μ
2. Find loss for that s : $[v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$
3. Find gradient of individual loss w.r.t. \mathbf{w} :
 $\nabla [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 = 2(v_\pi(s) - \hat{v}(s, \mathbf{w})) \nabla \hat{v}(s, \mathbf{w})$
4. Reduce individual loss by adjusting \mathbf{w} in opposite direction of gradient:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{1}{2} \alpha \nabla [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 = \\ &= \mathbf{w}_t + \alpha (\hat{v}(s, \mathbf{w}) - v_\pi(s)) \nabla \hat{v}(s, \mathbf{w}) \end{aligned}$$

Linear Methods

$\boldsymbol{x}(s)$: feature vector of s

$$\hat{v}(\boldsymbol{w}) = \boldsymbol{w}^T \boldsymbol{x}(s) = \sum_{i=1}^d w_i x_i(s)$$

SGD—Linear methods

For any weight, w_i , $\frac{\partial \hat{v}}{\partial w_i} = x_i$

Gradient of \hat{v} with respect to \mathbf{w} :

$$\nabla \hat{v}(s, \mathbf{w}) = \left[\frac{\partial \hat{v}(s, \mathbf{w})}{\partial w_1}, \frac{\partial \hat{v}(s, \mathbf{w})}{\partial w_2}, \dots, \frac{\partial \hat{v}(s, \mathbf{w})}{\partial w_d} \right] = \mathbf{x}(s)^T$$

To reduce $\overline{VE}(\mathbf{w})$, tweak \mathbf{w} in a direction such that $\overline{VE}(\mathbf{w})$ is reduced

SGD—Linear methods

Given loss function $\overline{VE}(\mathbf{w})$ which consists of a sum of individual losses: $\sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$

1. Pick an $s \in \mathcal{S}$ according to μ

2. Find loss for that s : $[v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$

3. Find gradient of individual loss w.r.t. \mathbf{w} :

$$\nabla [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 = 2(v_\pi(s) - \hat{v}(s, \mathbf{w})) \nabla \hat{v}(s, \mathbf{w})$$

4. Reduce individual loss by adjusting \mathbf{w} in opposite direction of gradient:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2} \alpha \nabla [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 =$$

$$\mathbf{w}_t + \alpha (\hat{v}(s, \mathbf{w}) - v_\pi(s)) \nabla \hat{v}(s, \mathbf{w}) = \mathbf{w}_t + \alpha (\hat{v}(s, \mathbf{w}) - v_\pi(s)) \mathbf{x}(s)$$