

## Recitation 8

Lecturer: Yishay Mansour

TA: Lee Cohen

## Value Function Approximation

**Definition 1.** We would like to reduce our reinforcement learning problem to a supervised learning problem. To do so we need:

1. Loss function: we will use **Mean Squared Error (MSE)** by defining

$$J(w) = \sum_s \mu(s) (V^\pi(s) - \hat{V}^\pi(s; w))^2$$

, where  $\mu(s)$  is the steady state distribution of  $\pi$ .

2. Update algorithm - in this case we will **Stochastic Gradient Descent (SGD)**.
3. Build samples for training by replacing  $V^\pi(s_t)$  with an estimated value (Monte-Carlo, TD(0), TD( $\lambda$ )).

Once we implement the conditions above, we can approximate the value function. Depending on the approximation used, the gradients updates accordingly:

1. **Monte-Carlo**

$$\Delta w_t = \alpha (R_t(s) - \hat{V}^\pi(s; w_t)) \nabla_w \hat{V}^\pi(s; w)$$

$$R_T(s_j) = \sum_{t=j}^T r_t$$

2. **TD(0)**

$$\Delta w(s) = \alpha [r_t + \gamma \hat{V}(s_{t+1}; w) - \hat{V}(s_t)] \nabla_w \hat{V}(s_t; w)$$

Note: This is a **semi-gradient**, since  $w$  influences the target  $R_t(s_t) = r_t + \gamma \hat{V}(s_{t+1}; w)$ .

3. **TD( $\lambda$ )** - forward view

$$\Delta w = \alpha [R_t^\lambda - \hat{V}(s_t; w)] \nabla_w \hat{V}(s_t; w)$$

4. **TD**( $\lambda$ ) - backward view

$$\begin{aligned}e_t &= \gamma \lambda e_{t-1} + \nabla_w \hat{V}(s_t; w) \\ \Delta_t &= r_t + \gamma \hat{V}(s_{t+1}; w) - \hat{V}(s_t; w) \\ \Delta w &= \alpha \Delta_t e_t\end{aligned}$$

In the case of linear function approximation the formulas are reduced to:

1. **Monte-Carlo**

$$\Delta w = \alpha [R_t - \hat{V}(s_t; w_t)] x(s_t)$$

2. **TD**(0)

$$\Delta w = \alpha [r_t + \gamma \hat{V}(s_{t+1}; w_t) - \hat{V}(s_t; w_t)] x(s_t)$$

3. **TD**( $\lambda$ ) - forward view

$$\Delta w = \alpha [R_t^\lambda - \hat{V}(s_t; w)] x(s_t)$$

4. **TD**( $\lambda$ ) - backward view

$$\begin{aligned}e_t &= \gamma \lambda e_{t-1} + x(s_t) \\ \Delta_t &= r_t + \gamma \hat{V}(s_{t+1}; w_t) - \hat{V}(s_t; w_t) \\ \Delta w &= \alpha \Delta_t e_t\end{aligned}$$



We now calculate the gradient

$$\nabla Q_t = \left( \frac{dQ_t}{dw_{hat}} \quad \frac{dQ_t}{dw_{hbt}} \quad \frac{dQ_t}{dw_{hct}} \right)$$
$$\nabla Q_0 = (c_{a0} \quad c_{b0} \quad c_{c0}) = (4 \quad 7 \quad 1)$$

Now we can calculate the new weights:

$$w_{h1} = (3 \quad 2 \quad 1) + 0.01[3 - 27] (4 \quad 7 \quad 1) = (2.04 \quad 0.32 \quad 0.76)$$

Repeating for sample 2:

$$Q(s_1, harvest_1) = 10 * 2.04 + 6 * 0.32 + 0 * 0.76 = 22.32$$
$$w_{h2} = (2.04 \quad 0.32 \quad 0.76) + 0.01[-15 - 22.32] (10 \quad 6 \quad 0) = (-1.69 \quad -0.72 \quad 0.76)$$

Repeating for sample 3:

$$Q(s_2, harvest_2) = 20 * -1.69 + 1 * -0.72 + 15 * 0.76 = -23.12$$
$$w_{h3} = (-1.69 \quad -0.72 \quad 0.76) + 0.01[5 + 23.12] (20 \quad 1 \quad 15) = (3.93 \quad 0.44 \quad 4.98)$$

Repeating for sample 4:

$$Q(s_3, harvest_3) = 4 * 3.93 + 19 * -0.44 + 3 * 4.98 = 22.3$$
$$w_{h4} = (3.93 \quad -0.44 \quad 4.98) + 0.01[21 - 22.3] (4 \quad 19 \quad 3) = (3.88 \quad -0.69 \quad 4.94)$$

**Exercise 2.** The orchard also has a record from its own last harvest. Continue to approxi-

A (ppm)	B (ppm)	C (ppm)	Action	Profit/Reward
6	7	2	Wait	-1
3	8	4	Harvest	19

mate  $Q(\text{state}, \text{harvest})$  and start approximating  $Q(\text{state}, \text{wait})$ . Run through this episode, performing gradient descent and using the TD(0) goal. Instead of choosing your actions in an  $\epsilon$ -greedy way, assume the  $\epsilon$ -greedy algorithm chose the actions in the table above.

**Solution 2.** For TD(0):

$$T_t = R_t + \gamma Q(s_{t+1}, a_{t+1})$$

Where  $a_{t+1}$  is chosen  $\epsilon$ -greedily. For this question we assume it always picks the actions in the episode provided. We can set  $\gamma = 1$  as there is an absorbing state but we could also select a lower value.

Need initial **wait** weight values:  $w_{w4} = (3 \ 2 \ 1)$

Day 1: Wait

$$Q(s_4, \text{wait}_4) = 6 * 3 + 7 * 2 + 2 * 1 = 34$$

We also need to calculate:

$$Q(s_5, \text{harvest}_5) = 3 * 3.88 + 8 * -0.69 + 4 * 4.94 = 25.88$$

$$w_{w5} = (3 \ 2 \ 1) + 0.01[-1 + 25.88 - 34] (6 \ 7 \ 2) = (2.45 \ 1.36 \ 0.82)$$

Day 2: Harvest

$$Q(s_5, \text{harvest}_5) = 25.88$$

The state after harvest is an absorbing state so  $Q(s_6, *) = 0$ .

We also need to calculate:

$$w_{h6} = (3.88 \ -0.69 \ 4.94) + 0.01[19 + 0 - 25.88] (3 \ 8 \ 4) = (3.67 \ 1.24 \ 4.6648)$$

**Exercise 3.** The orchard enters a new harvest season. Over the next three days the concentrations will be:

Day	A (ppm)	B (ppm)	C (ppm)
1	20	6	1
2	10	7	2
3	5	8	4

Run through an episode using these values. Choose your actions greedily and don't update  $Q$ . If you wait on the third day then the episode is over. If your algorithm harvested, did it harvest at its maximum possible  $Q(s, \text{harvest})$  or would it have had a higher value on one of the other two days? If your algorithm waited on the third day, which day was it closest to harvesting?

**Solution 3.** We need to find  $Q$  values for each day and choose our actions greedily. Using the final weight vectors so far:

$$w_{w5} = (2.45 \quad 1.36 \quad 0.82)$$

$$w_{h6} = (3.67 \quad 1.24 \quad 4.66)$$

Day 1:

$$Q(s_6, \text{wait}_6) = 20 * 2.45 + 6 * 1.36 + 1 * 0.82 = 57.98$$

$$Q(s_6, \text{harvest}_6) = 20 * 3.67 + 6 * -1.24 + 1 * 4.66 = 70.62$$

So the episode immediately ends in a harvest. For the other days:

Day 2:

$$Q(s_7, \text{harvest}_7) = 10 * 3.67 + 7 * -1.24 + 2 * 4.66 = 37.34$$

Day 3:

$$Q(s_8, \text{harvest}_8) = 5 * 3.67 + 8 * -1.24 + 4 * 4.66 = 27.07$$

So the orchard harvested on the highest value of  $Q(s, \text{harvest})$  that it would have had in those 3 days. If the orchard had waited on the third day then the closest day to harvesting would have been the one with the smallest value of

$$Q(s_t, \text{wait}_t) - Q(s_t, \text{harvest}_t)$$

#### **Exercise 4.**

1. What is the overall effect of increasing the learning rate? What happens when you set it too high? What happens when you set it too low?
2. Write a short RL problem where it would be better to use feature vectors to represent a state. Write one where it would be better to use discrete states.

#### **Solution 4.**

1. Increasing the learning rate overall increases the effect each sample has on determining the weights. If you set it too high, then the weights will never converge as each small error will create too large a change in the weights. If you set the learning rate too low, then the weights will take a long time to converge and will be vulnerable to getting stuck in local maxima.
2. For the inverted pendulum problem, a robot is trying to balance a rod on a hinge on top of it. As the rod can be at any angle the problem would be better modelled continuously. For a chess AI every possible state of the board can be a state in the model and so a discrete approach would be better.