

## 10.0 : Intro

In RL, we talk a lot about prediction and control.  
Prediction refers to the prediction of how well a fixed policy performs, and control refers to the search of a better policy.

- In ch 9, we talked about on-policy prediction, Now we talk about the control setting:

Function approximation of the action-value function

$$\hat{q}(s, a, w) \approx q_*(s, a).$$

- we will see semi-gradient SARSA which is an extension of the semi-gradient TD(0) from last chapter, but with action values.
- Episodic extension is straightforward but continuous is trickier.

## 10.1 : Episodic Semi-Gradient Control

Instead of learning state values  $V(s)$ , we're now learning action-values  $Q(s, a)$  using function approximation,

The approximate action-value function  $\hat{q} \approx q_T$  is the parameterized function w/ weight vector  $w$ . the general gradient-descent update for action-value prediction is :

$$w_{t+1} = w_t + \alpha [u_t - \hat{q}(s_t, a_t, w_t)] \nabla \hat{q}(s_t, a_t, w_t)$$

- It looks a lot like the state value version in Ch 9.
- $u_t$  is the target, for example,  $r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t)$  for one-step SARSA, but it could be the full Monte-Carlo return or any of the n-step SARSA returns.
- With the one-step SARSA target, the method is called Episodic Semi-gradient one-step SARSA.

This forms the action-value prediction part, and for control we need to couple that to policy improvement and action selection techniques;

- action selection :  $A_t^* = \arg \max_a \hat{q}(s_t, a, w_t)$
- policy improvement :  $\epsilon$ -greedy.

# Algo : 1-step SARSA using function approximation

## Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step size  $\alpha > 0$ , small  $\varepsilon > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop for each episode:

$S, A \leftarrow$  initial state and action of episode (e.g.,  $\varepsilon$ -greedy)

Loop for each step of episode:

Take action  $A$ , observe  $R, S'$

If  $S'$  is terminal:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

Go to next episode

Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedy)

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \underbrace{\gamma \hat{q}(S', A', \mathbf{w})}_{\text{1-step}} - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

$$S \leftarrow S'$$

1-step

$$A \leftarrow A'$$

// semi-gradient bc we're treating the target as if it were fixed, even though it actually depends on the weight we're updating.

## 10.2: Semi-gradient n-step SARSA

use n-step return as the target ( $U_t$ ) in the usual update.

$$w_{t+n} \doteq w_{t+n-1} + \alpha [G_{t:t+n} - \hat{q}(S_t, A_t, w_{t+n-1})] \nabla \hat{q}(S_t, A_t, w_{t+n-1})$$

where  $G_{t:t+n} \doteq$

$$\begin{cases} R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{q}(S_{t+n}, A_{t+n}, w_{t+n-1}), & \text{if } t+n < T \\ G_t, & \text{if } t+n \geq T. \end{cases}$$

## Algo: n-step SARSA using function approximation

### Episodic semi-gradient n-step Sarsa for estimating $\hat{q} \approx q_*$ or $q_\pi$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Input: a policy  $\pi$  (if estimating  $q_\pi$ )

Algorithm parameters: step size  $\alpha > 0$ , small  $\varepsilon > 0$ , a positive integer  $n$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

All store and access operations ( $S_t$ ,  $A_t$ , and  $R_t$ ) can take their index mod  $n + 1$

Loop for each episode:

    Initialize and store  $S_0 \neq$  terminal

    Select and store an action  $A_0 \sim \pi(\cdot | S_0)$  or  $\varepsilon$ -greedy wrt  $\hat{q}(S_0, \cdot, \mathbf{w})$

$T \leftarrow \infty$

    Loop for  $t = 0, 1, 2, \dots$ :

        If  $t < T$ , then:

            Take action  $A_t$

            Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$

            If  $S_{t+1}$  is terminal, then:

$T \leftarrow t + 1$

            else:

                Select and store  $A_{t+1} \sim \pi(\cdot | S_{t+1})$  or  $\varepsilon$ -greedy wrt  $\hat{q}(S_{t+1}, \cdot, \mathbf{w})$

$\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)

                If  $\tau \geq 0$ :

$G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$

                    If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n \hat{q}(S_{\tau+n}, A_{\tau+n}, \mathbf{w})$

$(G_{\tau:\tau+n})$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [G - \hat{q}(S_\tau, A_\tau, \mathbf{w})] \nabla \hat{q}(S_\tau, A_\tau, \mathbf{w})$

    Until  $\tau = T - 1$

## 10.3: Average rewards : a new problem setting for continuing tasks

We introduce a third classical setting for formulating the goal in MDPs.

- Episodic setting.
- Discounted setting.
- Average reward setting.

It applies to continuing problems (no start or end), there is no discounting. Discounting setting is problematic with function approximation; this setting replaces it.

In the average reward setting, the quality of a policy  $\pi$  is defined as the average rate of reward (short average reward)  $r(\pi)$  while following that policy:

$$r(\pi) = \lim_{h \rightarrow \infty} \underbrace{\frac{1}{h} \sum_{t=1}^h E[R_t | S_0, A_0:t-1 \sim \pi]}_{\substack{\text{add up} \\ \text{all the} \\ \text{expected} \\ \text{rewards} \\ \text{from } t=1 \\ \text{to } t=h.}} // \begin{array}{l} \text{Expected reward} \\ \text{per time step when} \\ \text{following policy } \pi, \\ \text{conditioned on initial} \\ \text{state } S_0. \end{array}$$

*h → ∞ implies this is continuous task.*

$$= \lim_{t \rightarrow \infty} E[R_t | A_0:t-1 \sim \pi] // \text{assumes } \mu_\pi(s) \doteq \lim_{t \rightarrow \infty} \Pr(S_t = s | A_0:t-1 \sim \pi)$$

$$= \sum_s \mu_\pi(s) \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) r,$$

$\mu_\pi(s)$  = the long-run probability of being in state  $s$  (steady-state distribution).

This works under the ergodicity assumption:  $\mu_\pi$  is the steady-state distribution, and it exists for any  $\pi$  independently of  $s_0$  (starting state and early decisions have no effect in the limit).

- All policies that attain the maximal value of  $r(\pi)$  are considered optimal.
- The steady state distribution is the distribution under which when you select an action you end up in the same distribution. Formally,

$$\sum_s \mu_\pi(s) \sum_a \pi(a|s) p(s'|s, a) = \mu_\pi(s')$$

Ergodicity assumptions:

- 1) Irreducibility: There is a non-zero probability to go from any state to any other state.
- 2) Aperiodicity: You can return to any state at irregular intervals.

## Differential Returns

Definition: In the average-reward setting, differential returns = differences btwn rewards and average rewards:

$$G_t = R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + \dots$$

The corresponding value functions are different value functions (remove all  $\gamma$  and replace rewards by the difference):

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) [r - r(\pi) + V_\pi(s')]$$

$$Q_\pi(s) = \sum_{s', r} P(s', r | s, a) [r - r(\pi) + \sum_{a'} \pi(a'|s') Q_\pi(s', a')]$$

$$V_*(s) = \max_a \sum_{s', r} P(s', r | s, a) [r - \max_\pi r(\pi) + V_*(s')]$$

$$Q_*(s) = \sum_{s', r} P(s', r | s, a) [r - \max_\pi r(\pi) + \max_{a'} Q_*(s', a')]$$

## Differential forms of TD error:

For state-values:

$$\delta_t = R_{t+1} - \bar{R}_{t+1} + \hat{V}(S_{t+1}, w_t) - \hat{V}(S_t, w_t)$$

For action-values:

$$\delta_t = R_{t+1} - \bar{R}_{t+1} + \hat{a}(S_{t+1}, A_{t+1}, w_t) - \hat{q}(S_t, A_t, w_t)$$

where  $\bar{R}_t$  is an estimate at time  $t$  of the average reward  $r(\pi)$ .

With these alternative definitions, most of the algorithms already seen carry through the average-reward setting without change. For example, the average reward version of semi-gradient SARSA is defined with the new version of the TD error:

$$w_{t+1} = w_t + \alpha \delta_t \nabla \hat{q}(S_t, A_t, w_t)$$

# Algo: Differential 1-step SARSA.

## Differential semi-gradient Sarsa for estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step sizes  $\alpha, \beta > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Initialize average reward estimate  $\bar{R} \in \mathbb{R}$  arbitrarily (e.g.,  $\bar{R} = 0$ )

Initialize state  $S$ , and action  $A$

Loop for each step:

Take action  $A$ , observe  $R, S'$

Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedy)

$$\delta \leftarrow R - \bar{R} + \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$$

$$\bar{R} \leftarrow \bar{R} + \beta \delta$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \nabla \hat{q}(S, A, \mathbf{w})$$

$$S \leftarrow S'$$

$$A \leftarrow A'$$

Algo has 2 learning processes:

1) Average reward estimation:  $\bar{R} \leftarrow \bar{R} + \beta \delta$

2) Action-value learning:  $\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \nabla \hat{q}(S, A, \mathbf{w})$

Others:

The general form for a semi-gradient temporal difference update is:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [\text{Target} - \text{Estimate}] \nabla (\text{Estimate}).$$

For semi-gradient n-step off-policy TD, the update formula is:

$$\mathbf{w}_{t+n} \doteq \mathbf{w}_{t+n-1} + \alpha p_t p_{t+1} \dots p_{t+n-1} [G_{t:t+n} - \hat{V}(S_t, \mathbf{w}_{t+n-1})] \nabla \hat{V}(S_t, \mathbf{w}_{t+n-1})$$

## 10.4: Depreciating the Discounted Setting

Section explains why discounting becomes problematic in the continuous task setting.

- The discounted setting is useful for the tabular case where the return for each state can be separated and averaged.
  - In continuing tasks, we could measure the discounted return at each time step, but it's useless:
    - Average of discounted rewards is proportional to average reward.
    - Policy ranking remains the same regardless of discount factor.
- ⇒ Discounting does not change which policy is "best".

Proof: Idea: Symmetry argument. Each time step is the same as every other, with discounting each reward appears exactly once in return, the  $t^{\text{th}}$  reward will be undiscounted in the  $t^{\text{st}}$  return, discounted once in the  $t-1^{\text{th}}$  return, and so on. The weight of the  $t^{\text{th}}$  reward is thus  $1 + \gamma + \gamma^2 + \dots = 1/(1-\gamma)$ . Because all states are the same, they are all weighted by this, and thus the average of the returns will be this weight times the average reward,

Proof:

## The Futility of Discounting in Continuing Problems

Perhaps discounting can be saved by choosing an objective that sums discounted values over the distribution with which states occur under the policy:

$$\begin{aligned} J(\pi) &= \sum_s \mu_\pi(s) v_\pi^\gamma(s) && \text{(where } v_\pi^\gamma \text{ is the discounted value function)} \\ &= \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) [r + \gamma v_\pi^\gamma(s')] && \text{(Bellman Eq.)} \\ &= r(\pi) + \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) \gamma v_\pi^\gamma(s') && \text{(from (10.7))} \\ &= r(\pi) + \gamma \sum_{s'} v_\pi^\gamma(s') \sum_s \mu_\pi(s) \sum_a \pi(a|s) p(s'|s, a) && \text{(from (3.4))} \\ &= r(\pi) + \gamma \sum_{s'} v_\pi^\gamma(s') \mu_\pi(s') && \text{(from (10.8))} \\ &= r(\pi) + \gamma J(\pi) \\ &= r(\pi) + \gamma r(\pi) + \gamma^2 J(\pi) \\ &= r(\pi) + \gamma r(\pi) + \gamma^2 r(\pi) + \gamma^3 r(\pi) + \dots \\ &= \frac{1}{1 - \gamma} r(\pi). \end{aligned}$$

The proposed discounted objective orders policies identically to the undiscounted (average reward) objective. The discount rate  $\gamma$  does not influence the ordering!

Root cause of the difficulties w/ discounting is that we lost the policy improvement theorem with function approximation. It is no longer true that if we change the policy to improve the discounted value of one state then we are guaranteed to have improved the overall policy. This lack of theoretical guarantee is an ongoing area of research.

## 10.5: Differential Semi-gradient n-step SARSA

To generalize to n-step bootstrapping, we need an n-step version of the TD error.

### N-step return (differential form)

each reward is differential.

$$G_{t:t+n} = \overbrace{R_{t+1} - \bar{R}_{t+1} + R_{t+2} - \bar{R}_{t+2} + \dots + R_{t+n} - \bar{R}_{t+n}}^{\text{bootstrap at the end: add estimated Q-value after } n\text{-steps.}} + \underbrace{\hat{q}(S_{t+n}, A_{t+n}, w_{t+n-1})}_{}$$

•  $\bar{R}_t$  is an estimate of  $r(\pi)$  at time  $t$ , with  $n \geq 1$  and  $t+n < T$ .

• If  $t+n \geq T$ , then we define  $G_{t:t+n} = G_t$  as usual.

### n-step error:

$$\delta_t = G_{t:t+n} - \hat{q}(S_t, A_t, w)$$

And then we can apply the usual semi-gradient SARSA update.

Algo :

### Differential semi-gradient $n$ -step Sarsa for estimating $\hat{q} \approx q_\pi$ or $q_*$

Input: a differentiable function  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$ , a policy  $\pi$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Initialize average-reward estimate  $\bar{R} \in \mathbb{R}$  arbitrarily (e.g.,  $\bar{R} = 0$ )

Algorithm parameters: step size  $\alpha, \beta > 0$ , a positive integer  $n$

All store and access operations ( $S_t$ ,  $A_t$ , and  $R_t$ ) can take their index mod  $n + 1$

Initialize and store  $S_0$  and  $A_0$

Loop for each step,  $t = 0, 1, 2, \dots$ :

Take action  $A_t$

Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$

Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ , or  $\varepsilon$ -greedy wrt  $\hat{q}(S_{t+1}, \cdot, \mathbf{w})$

$\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)

If  $\tau \geq 0$ :

$$\delta \leftarrow \sum_{i=\tau+1}^{\tau+n} (R_i - \bar{R}) + \hat{q}(S_{\tau+n}, A_{\tau+n}, \mathbf{w}) - \hat{q}(S_\tau, A_\tau, \mathbf{w})$$

$$\bar{R} \leftarrow \bar{R} + \beta \delta$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \nabla \hat{q}(S_\tau, A_\tau, \mathbf{w})$$

## 10.6 : Summary

- Extension of parametrized action-value functions and semi-gradient descent to control,
- Average-reward setting replacing discounting for continuous case.
- This new setting involves a differential version of everything.