

# TD and Friends

## RL Context

- $\langle s, a, r \rangle^* \rightarrow \text{RL algo} \rightarrow \pi$
- model-based (more supervised, less direct learning):
  - $\langle s, a, r \rangle^* \leftrightarrow \text{model learner} \rightarrow T, R \rightarrow \text{MDP solve} \rightarrow Q^* \rightarrow \text{argmax} \rightarrow \pi$
- value-function based (model-free):
  - $\langle s, a, r \rangle^* \leftrightarrow Q \rightarrow \text{argmax} \rightarrow \pi$
- policy search (more supervised, less direct learning):
  - $\langle s, a, r \rangle^* \rightarrow \text{policy update} \leftrightarrow \pi$

## TD( $\lambda$ ) basics

- Computing estimates incrementally
  - $V_T(s) = \frac{(T-1)V_{T-1}(s) + R_T(s)}{T} = \frac{T-1}{T}V_{T-1} + \frac{1}{T}R_T(s) = V_T(s) + \alpha_T(R_T(s) - V_{T-1}(s))$ , where  $\alpha_T = \frac{1}{T}$
  - interpretation: value is updated by the difference b/w the return at current time and value at previous time
  - Isbell: looks like perceptron update (the diff is error)
- properties of learning rates
  - if  $\sum_T \alpha_T = \infty$  and  $\sum_T \alpha_T^2 < \infty$ , then  $\lim_{T \rightarrow \infty} V_T(s) = V(s)$
  - interpretation: sum of alphas must be big enough to move to true value, but not so big they hide noise

## TD(1)

- let  $e(s)$  be the eligibility of  $s$  (kind of like a filter for application of the rule that reflects age of states)
- $V_T(s) = V_T(s) + \alpha_T(r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))e(s)$
- Big idea: telescoping series on updating values
- behaves like an infinite-step ahead estimator
- update rule:

# TD(1) Rule

Episode T

For all  $s$ ,  $e(s) = 0$  at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{r_t} s_t$  : (step t)

$$e(s_{t-1}) = e(s_{t-1}) + 1$$

$$\text{For all } s, \\ V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1})) e(s) \\ e(s) = \gamma e(s)$$

- $\Delta V_T(s_n) = \alpha(r_n + \gamma r_{n+1} + \gamma^2 r_{n+2} + \dots + \gamma^{N-1} V_{T-1}(s_N) - V_{T-1}(s_n)) =$  "what you predicted for val of s vs. what value of s was last episode"
- TD(1) is the same outcome-based updates (if no repeated states)
  - TD(1) allows for extra learning when repeating a state
- maximum likelihood estimate can be better than outcome-based TD(1) since it uses more data (see [here](#) for more)

## TD(0)

- equivalent to maximum likelihood estimate?
- removes eligibility vector
- behaves like a 1-step ahead estimator

## TD( $\lambda$ )

- updates based on differences b/w temporally successive predictions
- rules comparison:

# TD( $\lambda$ ) Rule

Episode T

For all  $s$ ,  $e(s) = 0$  at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{r_t} s_t$ : (step t)

$$e(s_{t-1}) = e(s_{t-1}) + 1$$

For all  $s$ ,

$$V_T(s) = V_T(s) + \alpha_r (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1})) e(s)$$

$$e(s) = \gamma e(s)$$

✓

TD(1)

TD(0) ✓

Episode T

For all  $s$ ,  $e(s) = 0$  at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{r_t} s_t$ : (step t)

$$e(s_{t-1}) = e(s_{t-1}) + 1$$

For all  $s$ ,

$$V_T(s) = V_T(s) + \alpha_r (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1})) e(s)$$

$$e(s) = \gamma e(s)$$

Episode T

For all  $s$ , at start of episode,  $V_T(s) = V_{T-1}(s)$

After  $s_{t-1} \xrightarrow{r_t} s_t$ : (step t)

For all  $s = s_{t-1}$

$$V_T(s) = V_T(s) + \alpha_r (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))$$

- if  $\lambda$  is not 0 or 1, TD( $\lambda$ ) is a weighted combination of all k-step estimators