

27 Oct

Mercedes-Benz R & D India

## Chapter 6 : Temporal - Difference Learning

No model is assumed (Interplay between dynamic programming and monte-carlo techniques)

Used for problem of prediction - find value function corresponding to policy  $\pi$ .

Look-Up table methods: Estimated values of all states are stored in an array.

Basic Algorithm:  $V(S_t) \leftarrow V(S_t) + \alpha (r_t - V(S_t))$

$\alpha$  - step size parameter.

$V(S_t)$  - non-stationary env

$$V_{t+1}(S_t) = V_t(S_t) + \alpha_t (r_t - V_t(S_t))$$

suppose  $\alpha_t, t \geq 0$  are such that

$$\sum_t \alpha_t = \infty, \quad \sum_t \alpha_t^2 < \infty$$

eg:  $\frac{1}{t+1}$ ,  $\frac{1}{(t+1)^\beta}$   $\beta > 1$ ,  $t \geq 0$

eg:  $\frac{\log(t+2)}{(t+2)}$ ,  $t \geq 0$

$$V_{t+1}(S_t) = V_t(S_t) + \alpha_t (E[r_t - V_t(S_t) | S_t] + M_{t+1})$$

$$\text{where } M_{t+1} = (r_t - V_t(S_t)) - E[(r_t - V_t(S_t)) | S_t]$$

$$E[M_{t+1} | S_0, S_1, \dots, S_t] = 0 \quad \text{a.s.}$$

If the rewards are bounded,  $\sup_t |r_t| \leq B < \infty$

then,  $\sup_t |G_t - V_t(s_t)| < \infty$  as well

$$\text{then } N_t = \sum_{m=0}^t \alpha_m M_{m+1}$$

$\{N_t\}$  - martingale sequence which is convergent

thus  $\alpha_m M_{m+1} \rightarrow 0$  as  $m \rightarrow \infty$

ODE:

$$\dot{V}_t(s_t) = E[G_t - V_t(s_t) | s_t]$$

Suppose  $\{s_t\}$  has a unique stationary distribution

$d^\pi = (d^\pi(s), s \in S)$  where  $\pi$  is the policy

$$V_t(s_t) = \sum_{s \in S} d^\pi(s) \left( \sum_{s', r} p(s', r | s, \pi) (r + \gamma V_t(s') - V_t(s)) \right)$$

Stochastic fixed points of the ODE will correspond to

$$\sum_{s \in S} d^\pi(s) V^*(s) = \sum_{s \in S} d^\pi(s) \left( \sum_{s', r} p(s', r | s, \pi) (r + \gamma V^*(s')) \right)$$

TD(0) for estimating  $V_\pi$

Input: Policy  $\pi$  to be evaluated, step size  $\alpha \in (0, 1]$   
& discount factor  $\gamma$

Initialize  $V(s), \forall s \in S^+$ , except  $V(\text{terminal state}) = 0$   
loop (for each episode)

Initialize  $s$

loop for each step of episode

$A \leftarrow$  action given by  $\pi$  for  $s$

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

$$V(s) \leftarrow V(s) + \alpha (R + \gamma V(s') - V(s))$$

$$s \leftarrow s'$$

Note:

$$V_{\pi}(s) = E_{\pi}[G_t | S_t = s] \text{ — Monte Carlo Estimate}$$

$$= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= E_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s]$$

$\hookrightarrow$  TD estimate

Bootstrapping method.

$$\delta_t \triangleq (R_t + \gamma V(S_{t+1}) - V(S_t))$$

$\uparrow$  TD term error

Note:  $E[\delta_t | S_t] = 0$ .

$$G_t - V(S_t) = R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1})$$

$$= \delta_t + \gamma (G_{t+1} - V(S_{t+1}))$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 (G_{t+2} - V(S_{t+2}))$$

$\vdots$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 \delta_{t+2} + \dots + \gamma^{T-t-1} \delta_{T-1} + \gamma^{T-t} (G_T - V(S_T))$$

$$= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k$$



$TD(x)$

Given a trajectory of states  $i_0, i_1, \dots, i_N$  with  $i_N = 0$  [terminal state]

$S = \{1, 2, \dots, N\}$  &  $0$  is the terminal state.

assume policy is fixed as policy  $\pi$ .

When  $i_N = 0$ ,  $i_k = 0 \quad \forall k > N$

Set  $r(i_k, i_{k+1}) = 0 \quad \forall k \geq N$

expected reward Value current state =  $i_k$  & next state  $i_{k+1}$

$V(i_k)$  = Value estimate when starting from  $i_k$

$$V(i_k) := V(i_k) + \gamma (r(i_k, i_{k+1}) + \gamma V(i_{k+1}) - V(i_k)) \quad (2)$$

$$\begin{aligned} &= V(i_k) + \gamma \left( r(i_k, i_{k+1}) + V(i_{k+1}) - V(i_k) \right) \\ &\quad + \gamma \left( r(i_{k+1}, i_{k+2}) + V(i_{k+2}) - V(i_{k+1}) \right) \\ &\quad + \dots + \\ &\quad + \gamma \left( r(i_{N+1}, i_N) + V(i_N) - V(i_{N+1}) \right) \end{aligned}$$

$$\text{If } \delta_k = r(i_N, i_{N+1}) + V(i_{k+1}) - V(i_k)$$

$$\text{then, } V(i_k) := V(i_k) + \gamma (\delta_k + \delta_{k+1} + \dots + \delta_{N-1})$$

Incremented Updates.

$$V(i_k) := V(i_k) + \alpha \delta_k \quad \left| \gamma = 1 \text{ discount factor} \right.$$

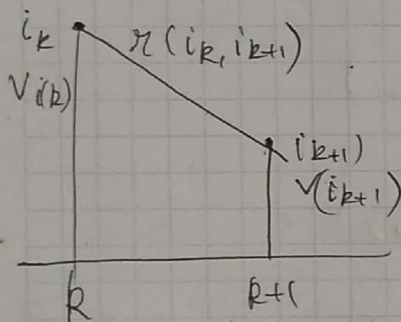
TD( $\lambda$ )

$$V(i_k) = E \left[ \sum_{m=0}^{\infty} \gamma^m r(i_k, i_{k+m+1}) \mid i_k \right], i_k \in S$$

\* Under the given policy, the process terminates in a finite (possibly random) time w.p 1 [proper policy]

Bellman Eqn:

$$\begin{aligned} V(i_k) &= E [r(i_k, i_{k+1}) + V(i_{k+1})] \\ &= E [r(i_k, i_{k+1}) + \gamma V(i_{k+1}, i_{k+2}) + V(i_{k+2})] \end{aligned}$$



$$= E \left[ \sum_{m=0}^l r(i_{k+m}, i_{k+m+1}) + V(i_{k+l+1}) \right]$$

Since the value of  $l$  is arbitrary, we form a weighted average of the bellman Equations. With ~~weight~~ weight  $\lambda$

Note that:  $\sum_{l=0}^{\infty} (1-\lambda) \lambda^l = 1$  we then form a weighted Bellman Equation.

$$V(i_k) = (1-\lambda) E \left[ \sum_{l=0}^{\infty} \lambda^l (r(i_{k+l}, i_{k+l+1}) + V(i_{k+l+1})) \right] \quad \text{--- (1)}$$

$$(1-\lambda) \sum_{l=0}^{\infty} \lambda^l V(i_k) = V(i_k)$$

$$\therefore V(i_k) = E \left[ \sum_{m=0}^{\infty} \gamma^m r(i_{k+m}, i_{k+m+1}) + V(i_{k+l+1}) \right], \lambda \geq 0$$



continuing with ①

$$= (1-\lambda) E \left[ \sum_{l=0}^{\infty} \lambda^l \sum_{m=0}^l r(i_{k+m}, i_{k+m+1}) \right] + (1-\lambda) E \left[ \sum_{l=0}^{\infty} \lambda^l v(i_{k+l+1}) \right]$$

I II

$$\text{I} = (1-\lambda) E \left[ \sum_{l=0}^{\infty} \lambda^l \sum_{m=0}^l r(i_{k+m}, i_{k+m+1}) \right]$$

$$= (1-\lambda) E \left[ \sum_{m=0}^{\infty} \left( \sum_{l=m}^{\infty} \lambda^l \right) r(i_{k+m}, i_{k+m+1}) \right] \quad \sum_{l=m}^{\infty} \lambda^l = \frac{\lambda^m}{1-\lambda}$$

$$= E \left[ \sum_{m=0}^{\infty} \lambda^m r(i_{k+m}, i_{k+m+1}) \right]$$

$$\text{II} = (1-\lambda) E \left[ \sum_{l=0}^{\infty} \lambda^l (v(i_{k+l+1})) \right]$$

$$= E \left[ \sum_{l=0}^{\infty} (\lambda^l - \lambda^{l+1}) v(i_{k+l+1}) \right]$$

$$= E \left[ \sum_{l=0}^{\infty} (\lambda^l - \lambda^{l+1}) v(i_{k+l+1}) \right]$$

$$= E \left[ v(i_{k+1}) - \lambda v(i_k) + \lambda (v(i_{k+2}) - v(i_{k+1})) + \dots + v(i_k) \right]$$

$$= E \left[ \sum_{m=0}^{\infty} \lambda^m (v(i_{k+m+1}) - v(i_{k+m})) \right] + v(i_k)$$

$$\text{Thus, } E \left[ \sum_{m=0}^{\infty} \lambda^m r(i_{k+m}, i_{k+m+1}) \right] +$$

$$E \left[ \sum_{m=0}^{\infty} \lambda^m (v(i_{k+m+1}) - v(i_{k+m})) \right] + v(i_k)$$

$$v(i_k) = E \left[ \sum_{m=0}^{\infty} \lambda^m (r(i_{k+m}, i_{k+m+1}) + v(i_{k+m+1}) - v(i_{k+m})) \right] + v(i_k)$$

$$\Rightarrow E \left[ \sum_{m=0}^{\infty} \lambda^m (r(i_{k+m}, i_{k+m+1}) + V(i_{k+m+1}) - V(i_{k+m})) \right] = 0$$

letting  $\delta_k = r(i_{k+m}, i_{k+m+1}) + V(i_{k+m+1}) - V(i_{k+m})$

Then,  $E \left[ \sum_{m=0}^{\infty} \lambda^m \delta_{k+m} \right] = 0$

or  $E \left[ \sum_{m=t}^{\infty} \lambda^{m-k} \delta_m \right] = 0$

not surprising since  $E[\delta_m] = 0$  anyways.

Algorithm

$$V(i_k) := V(i_k) + \gamma \sum_{m=k}^{\infty} \lambda^{m-k} \delta_m$$

$\lambda = 1$  TD(1) monte carlo based policy estimation.

$$V(i_k) := V(i_k) + \alpha (r(i_k, i_{k+1}) + r(i_{k+1}, i_{k+2}) + \dots + r(i_{N-1}, i_N) - V(i_k))$$

$\lambda = 0$  TD(0)

$$V(i_k) := V(i_k) + \gamma \delta_k$$

(check)

Example 6.1

Driving Home:

How does one select  $\lambda$

Bootstrapping happens when  $\lambda < 1$

Better efficiency of algorithm

don't wait for termination, updation happens on the fly