

Model Free Prediction

a) Monte-Carlo Learning

MC methods learn directly from episodes of experience, we don't need knowledge of the MDP/transitions or rewards.

We look at complete episodes, e.g. we start a process like play a game then after a while it does finish.

We then go through our episodes, take sample returns for each episode, & then we average over them.

The only caveat is this works for episodic MDPs only, that is the episodes must terminate.

b) Monte-Carlo Policy Evaluation

Goal: Learn V_{π} from episodes of experience under policy π

$$S_1, A_1, R_1, \dots, S_T \sim \pi$$

So given a policy π we take a look at the episode of the state, action, reward that follows

We can calculate the return from a particular state as well as either the total discounted reward

$$g = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

and our value function is the expected return

$$V\pi(s) = E\pi [Y_t | S_t = s]$$

However MC methods use empirical mean return instead of expected return

(c) First Visit - Monte Carlo Policy Evaluation

Ideally a state appears at different times in different episodes, so every time it appears, we compute the return from that point onwards and then for all episodes we average out return values for each state

Now lets say in a particular episode we have multiple visits to the state, so we can have a value for this state at a single particular time step onwards. So the question is whether to treat the multiple visits as multiple or one sample

In first visit MC, in each episode we only use the first time a state appears

So here if our goal is to evaluate a particular state,

We maintain an $N(s)$ & $S(s)$ for each state
 $N(s)$ tells us how many samples has
 appeared in s & $G(s)$ tells us sum of
 returns from those samples

So the first time-step if that state s is visited in an episode we increment $N(s) \leftarrow N(s) + 1$,
 we calculate return for state s . given that
 time step in the episode
 Then we also update $G(s) \leftarrow G(s) + S(s)$

Then after multiple such episodes the
 value is estimated as $V(s) = \frac{S(s)}{N(s)}$

The law of large numbers says if we do this
 enough times or as $N(s) \rightarrow \infty$, $V(s) \rightarrow V_{\pi}(s)$

D) Every-Visit Monte Carlo Policy Evaluation

The only difference is that at every visit of s in an episode we update $N(s)$ & $S(s)$ at all t's we encountering s

E) Incremental Mean:

The mean $\mu_1, \mu_2, \mu_3, \dots$ of a sequence x_1, x_2, \dots
 can be computed incrementally,

$$\begin{aligned} \mu_K &= \frac{1}{K} \sum_{j=1}^K x_j \\ &= \frac{1}{K} \left(x_K + \sum_{j=1}^{K-1} x_j \right) \end{aligned}$$

$$= \frac{1}{k} (x_k + (k-1)u_{k-1})$$

$$= u_{k-1} + \frac{1}{k} (x_k - u_{k-1})$$

F) Incremental Monte-Carlo Updates

Now, we take over MC Algorithm 2 instead of maintaining the sum, we update episodically now

- Update $V(s)$ incrementally after episode $s_1, s_2, s_3, \dots, s_t$
- For each state s_t with return y_t

$$N(s_t) \leftarrow N(s_t) + 1$$

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (y_t - V(s_t))$$

In nonstationary problems, it can be useful to track a running mean i.e forget old episodes

$$V(s_t) \leftarrow V(s_t) + \alpha (y_t - V(s_t))$$

with alpha we track more recent avg's
 if α is large we highly prioritize to recent returns
 forget past episodes quickly, if its small
 it means there are slow stable updates

g) Temporal Difference Learning

- TD methods learn directly from episodes of experience → it also applies to non-episodic cases
- TD is model free when no knowledge of MDP transition or rewards are required
- We can learn from incomplete episodes, we can take a partial trajectory & then take an estimate for the actions ahead
- TD updates a guess towards a guess

H) ML & TD

Goal: Learn π_t online from experience under policy π

Simplist Temporal-Difference learning Algorithm: TD(0)

→ Update value $V(s_t)$ toward estimated return

$$R_{t+1} + \gamma V(s_{t+1})$$

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

$R_{t+1} + \gamma V(s_{t+1})$ is called the TD target

$s_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$ is called the TD error

Also here we simple wanna get or the estimated return to be the immediate reward R_{t+1} plus the discounted value for the reward of the next state

2) Advantages & Disadvantages of MC vs TD

- TD can learn before knowing the final outcome
TD learns online after every step & MC must wait until end of episode before return is known
- TD can learn without the final outcome
It learns from incomplete sequences
MC only learns from complete ones
TD works in continuing (non-terminating) environment
MC only works for episodic (terminating) environment

3) Bias / Variance trade off

- Between $g_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-1} R_T$ is unbiased estimate of $\sqrt{\pi}(s_t)$
- True TD target $R_{t+1} + \gamma V_{\pi}(s_{t+1})$ is unbiased estimate of $V_{\pi}(s_t)$

Moreover TD target $R_{t+1} + \gamma V(s_{t+1})$ is a biased estimate of $V_{\pi}(s_t)$

TD target is much lower variance than return which depends on many random actions, transitions towards, however target depends on just one random action, transition reward

K) MC has high variance, zero bias

- Good convergence properties
- Not very sensitive to initial value
- Very simple to understand & use

TD has low variance & high some bias

- Usually more efficient than MC
- TD(0) converges to $V_{\pi}(s)$
(but not always with function approximation)
- More sensitive to initial value

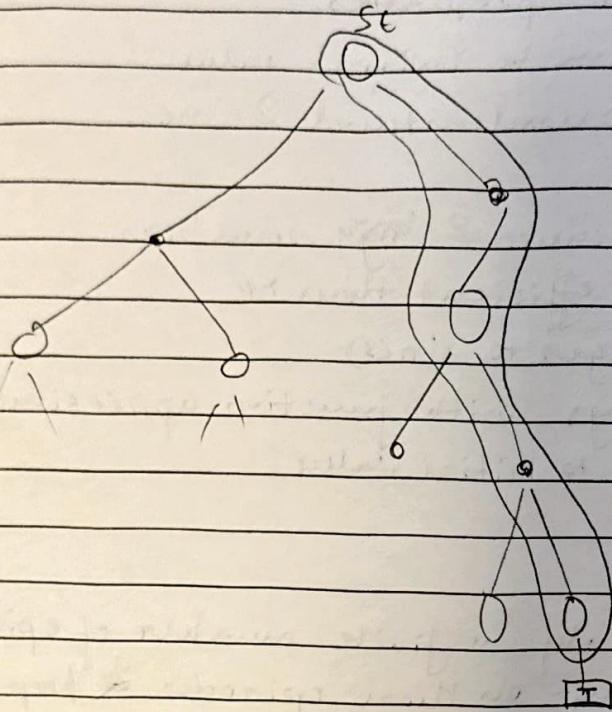
L) Batch TD & MC

If we have only a finite number of episodes & keep iterating on those episodes & keep learning based on the same data, would MC & TD find the same solution?

Monte-Carlo tries to converge to solution which finds the best fit to the entire universe whereas TD(0) converges to the solution that best explains the data

M) TD exploits Markov Property. So this is usually more efficient in markov environments
But MC does not exploit Markov property, but usually more efficient in non-markov environments

a) MC Backup

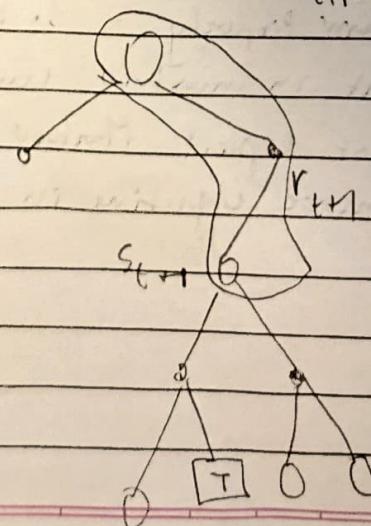


$$V(s_t) \leftarrow V(s_t) + \alpha (y_t - V(s_t))$$

here for every state we update $V(s_t)$
given the returns from that point to the
end

b) TD Backup

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$



how we update the value for s_t just given a single look ahead.

(b) n-Step Returns

Instead of a single step look ahead like we have for TD(0) we can have a multi-step n-step look ahead.

\therefore for n-step returns $n=1, 2, \infty$

$$n=1 \quad (\text{TD}) \quad g_t^{(1)} = R_{t+1} + \gamma V(s_{t+1})$$

$$n=2 \quad g_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(s_{t+2})$$

$$n=\infty \quad (\text{MC}) \quad g_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

So we can define the n-step return

$$g_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(s_{t+n})$$

n-step temporal difference learning

$$V(s_t) \leftarrow V(s_t) + \alpha (g_t^{(n)} - V(s_t))$$

(8) Averaging n Step Returns

We can avg n-step returns over different n
eg average the 2 step & 4 step returns

$$\frac{1}{2} y^{(2)} + \frac{1}{2} y^{(4)}$$

We combine information from two different timesteps

(9) TD - λ

The λ -return y_t^λ combines all n-step returns

$$y_t^{(n)}$$

using weight $(1-\lambda) \lambda^{n-1}$

$$y_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} y_t^{(n)}$$

Forward view

$$v(s_t) \leftarrow v(s_t) + \lambda (y_t^\lambda - v(s_t))$$

The issue with the forward view is we still need to wait for n steps to compute $y_t^{(n)}$ or y_t^λ

S) Batch update view

For every state we maintain an eligibility trace
Update value $v(s)$ for every state

$$E_t(s) = \gamma E_{t-1}(s) + \eta (s_t = s)$$

In proportion to TD-error of Eligibility trace
 $E_t(s)$

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

$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$

T) Understanding with an example

A \rightarrow B \rightarrow C \rightarrow D \rightarrow A \rightarrow terminal
+1 +1 +1 +1 +0

$$\gamma = 1$$

$$\lambda = 0.5$$

$$\alpha = 0.1$$

We initialize $V(A) = V(B) = V(C) = V(D) = 0$

Now we trace $E(s)$ & $V(s)$

With eligibility trace the goal is to update it for every state which has appeared till the given time step

initial

- Thus at step 1

$$E(A) = 1$$

$$E(B) = 0$$

$$E(C) = 0$$

$$E(D) = 0$$

$$\left(\text{uptake } E_t(s) = s \wedge E_{t-1}(s) + 1 (s = S_t) \right)$$

$$(S = S_{t+1} \neq X \vee (S_{t+1}) - V(S_t))$$

$$V(S) = V(S) + \alpha \times \delta_t \times E_t(s)$$

For $A \rightarrow B$ $r=1$

$$TD\text{-error} = \delta = 1 + V(B) - V(A) = 1$$

$$V(A) \leq 0.1 \times 1 \times 1 = 0.1$$

For Step 2

$$B-C, r=1$$

$$E(A) = 0.5$$

$$E(B) = 1$$

$$E(C) = 0$$

$$E(D) = 0$$

TD error

$$\delta = 1 + V(C) - V(B) = 1$$

$$V(A) \leq 0.1 \times 1 \times 0.5 = 0.15$$

$$V(B) \leq 0.1 \times 1 \times 1 = 0.1$$

For Step 3

$$C \rightarrow D, \alpha_1 = 1$$

$$E(A) = 0.25$$

$$E(B) = 0.5$$

$$E(C) = 1$$

$$E(D) = 0$$

TD error

$$\delta = 1 + V(D) - V(C) =$$

Value Update

$$V(A) \leftarrow 0.25 = 0.175$$

$$V(B) \leftarrow 0.5 = 0.15$$

$$V(C) \leftarrow 1 = 0.1$$

Step 4 (~~you made a mistake~~) (There is an error here
you caught it too late)

$$D \rightarrow A, \alpha_1 = 1$$

$$E(A) = 0.125 + 0.125 = 0.25$$

(The ~~+1~~ is wrong)

$$E(B) = 0.25$$

(approximated)

$$E(C) = 0.5$$

$$E(D) = 1$$

$$\delta = 1 + V(A) - V(D) = 1 + 0.175 - 0 = 1.175$$

Value update

$$V(A) \leftarrow 0.1 \times 1.175 \times 0.125 = +0.0125 \text{ (needs to be added)}$$

$$V(B) \leftarrow 0.1 \times 1.175 \times 0.25 = +0.029$$

$$V(C) \leftarrow 0.1 \times 1.175 \times 0.5 = +0.059$$

$$V(D) \leftarrow 0.1 \times 1.175 \times 1 = +0.118$$

$$\therefore V(A) = 0.1896$$

$$V(B) = 0.179$$

$$V(C) = 0.159$$

$$V(D) = 0.118$$

Step 5 $A \rightarrow \text{term}$

$$E(A) = 0.0625 + 1 = 1.0625$$

$$E(B) = 0.125$$

$$E(C) = 0.25$$

$$E(D) = 0.5$$

TD error

$$\delta = 0 + 0 - V(A) = -0.1896$$

Then update $V(A), V(B), V(C), V(D)$

The point is at every time step you update E for all states run till the given time step & calculate the error/TD error & propagate to back to every state seen in the episode.

$$\therefore V(A) = 0.1896$$

$$V(B) = 0.179$$

$$V(C) = 0.159$$

$$V(D) = 0.118$$

Step 5 $A \rightarrow \text{turn}$

$$E(A) = 0.0625 + 1 = 1.0625$$

$$E(B) = 0.125$$

$$E(C) = 0.25$$

$$E(D) = 0.5$$

τ_D error

$$\delta = 0 + 0 - V(A) = -0.1896$$

Then update $V(A), V(B), V(C), V(D)$

The point is at every time step you update E for all states run till the given time step & calculate the error (τ_D) & propagate to back to every state seen in the episode.