

Lecture 4

- No MDP given, although the environment is an MDP like system.
- We will discuss model free prediction - i.e. estimate the value function of an unknown MDP.

Methods

- Monte Carlo Learning → Learns directly from the episodes of experience.
 - ↳ we take sample return from the episodes and average over them
 - ↳ It only works for episodic experience i.e. we need to terminate to be able to work
 - ↳ Value = mean return over many episodes

Goal: Learn V_π from episodes of experience under policy π .

- ↳ we will look at the total reward in each time step ($G_{1:t}$) and then take an average of it to calculate $\tilde{V}(s)$ i.e. $E(G_{1:t})$
- ↳ use empirical mean return instead of expected return.

- First Visit Monte Carlo Policy Estimation →
- ↳ to evaluate state s
 - ↳ While over the series of episodes we look at the # number of time it we visited the state s the first time $N(s) \rightarrow N(s) + 1$
 - ↳ Increment the total return as we visit it $s(s) \leftarrow s(s) + G_t$
 - ↳ Value is estimated by mean return $V(s) = s(s) / N(s)$
 - ↳ By law of large numbers $V(s) \rightarrow V_\pi(s)$ as $N(s) \rightarrow \infty$
- Every Visit Monte Carlo Policy Estimator →
- ↳ Same as the above except that the visit to the episode need not only be the first to be included in the mean. If we visit a state multiple times in an episode, it will be included in the $N(s)$ and $s(s)$ ~~those many times~~ every time.
 - Incremental mean → Update the mean after every time step and not necessarily after the run is over.

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

$$= \frac{1}{K} (x_{k-1} - (k-1)u_{k-1})$$

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

$$u_k = u_{k-1} + \frac{1}{c} (\Delta u_k)$$

↓ old mean () ↓ constant → error term

→ Incremental Monte Carlo Updates \rightarrow update $V(s)$ incrementally after episode $s_1, A_1, R_2, \dots, s_T$

↳ For each state s_t with return G_t

$$\cancel{N(s)} \quad N(s_t) \leftarrow N(s_t) + 1$$

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

If it is a non stationary problem
i.e. the MDP changes with time steps, we
need a method which gives more weight
to latest observation in the mean.

$$V(s_t) \leftarrow V(s_t) + \underline{\alpha} (G_t - V(s_t))$$

$$0 < \alpha < 1$$

$$\alpha$$

→ Temporal Difference Learning



- ↳ learn directly from actual experience
- ↳ TD is model-free
- ↳ learns from incomplete episodes
- we can learn while the episode is still running and ~~does~~ doesn't require the episode to end to update the value functions. It utilises "bootstrapping"; we update the value function using the value functions of the other states.

Goal: learn V_π online from experience under policy π .

Online implies that we will update our value function while the episode is still running.

- ↳ we update our value function towards the ~~optimal~~ estimated return $R_{t+1} + \gamma V(s_{t+1})$

Note In TD(0) we take one step from state s to next state s' and get a return R . The estimated value of s is updated using the R and the estimated value of s' i.e. $V(s')$ using the bellman equation.



Then an ~~mean~~ incremental mean is performed on the $V(s)$ using the new bootstrapped value of $V(s)$ from R and $V(s')$

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

$$\Rightarrow \underbrace{V(s_t)}_{\text{New Value}} \leftarrow \underbrace{V(s_t)}_{\text{old Value}} + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

Error term (TD error)

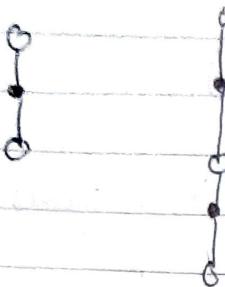
$R_{t+1} + \gamma V(s_{t+1})$ is the TD target

→ Bias / Variance trade off \Rightarrow

- ↳ TD return $R_{t+1} + \gamma V(s_{t+1})$ is biased estimate of $V_\pi(s_t)$, it would be unbiased if $V_\pi(s_{t+1})$ is the actual value of the value function. But the variance is low because the ~~randomness~~ randomness is only due to the reward after the state s. It compared to monte carlo which considers all the rewards till the episode terminates and therefore will have very different values everytime we sample it.

→ n steps Prediction →

TD(0) TD(1)



Monte Carlo



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

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

$$n=\infty \quad (\text{MC}) \quad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-1} R_T$$

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

n step temporal difference Learning

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

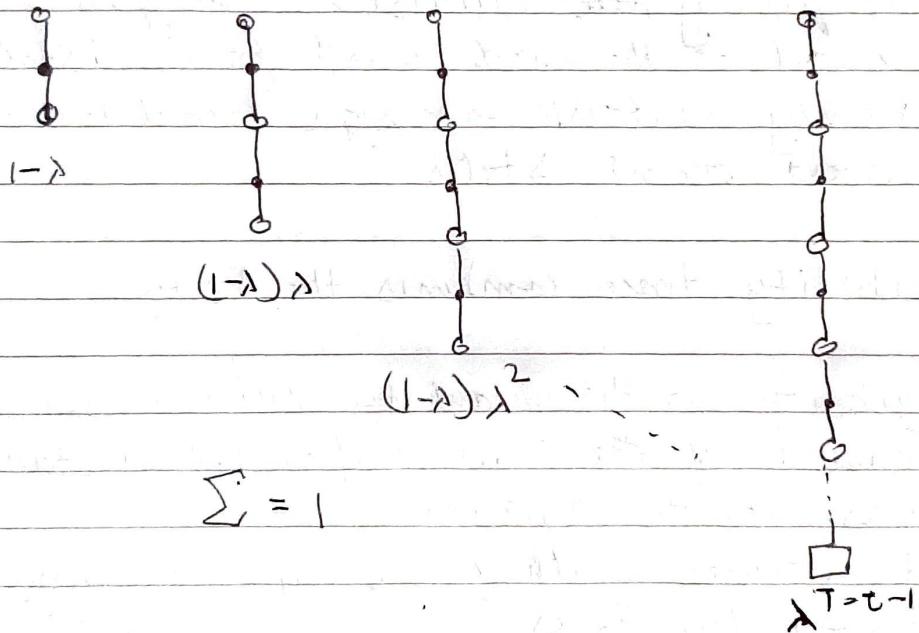
→ Average n step Returns

e.g. we can average 2 step and 4 step returns

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

(can we do it for all n)

→ TD (λ)



$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} g_{t+n} c_n$$

Forward view

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^\lambda - V(s_t))$$

→ Forward-view TD(λ)

Like MC it can only be computed from complete episode

→ Backward-view TD(λ)

What caused a reward?

- ~~but~~ or which state was most responsible
- Frequency ~~but~~ heuristic → assigns credit to the most recent frequent state
- Recency heuristic → assigns credit to the most recent states

Eligibility trace combines the both

Eligibility is the quantity which is a measure of ~~how~~ how responsible a state was for the reward.

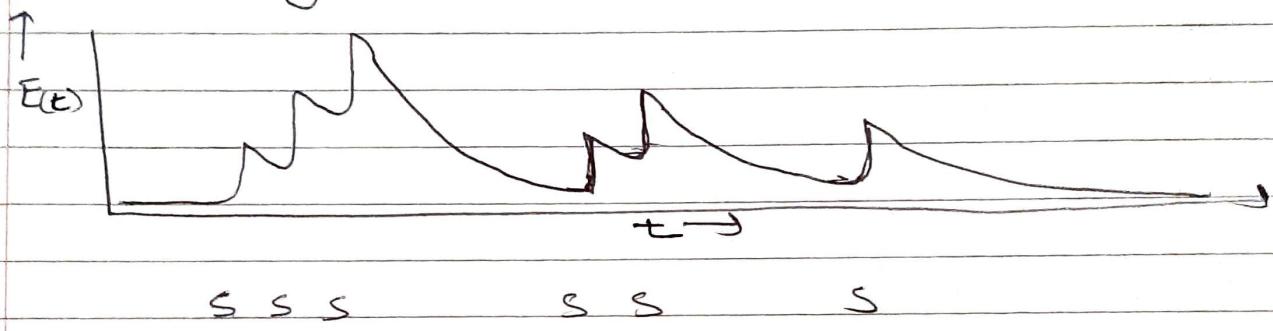
It increases with frequency and recency
→ denoted by $E_t(s)$

default value = 0 ; $E_0(s) = 0$

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

When ever we encounter the state s we increase its $\gamma \lambda E_{t-1}$ by 1 and make it the new $E_t(s)$

otherwise $E_t(s)$ is updated by multiplying with γ and δ both are less than 1. Therefore as the time step increases, the eligibility of the state decreases until it occurs again when its value is increased by $\delta \cdot 1$ after reducing it by γ and δ time.



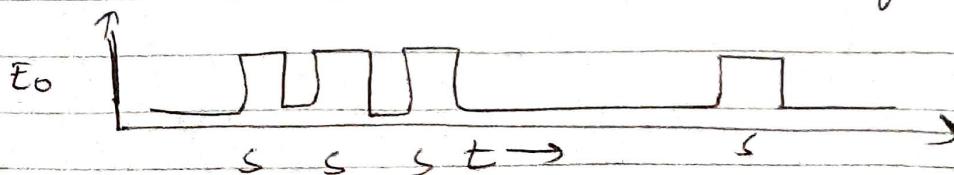
- keep an eligibility trace for every state s .
- update value $V(s)$ for every state s in proportion to the TD error and the eligibility trace $E_t(s)$

$$\begin{aligned}\delta_t &= R_{t+1} + \gamma V(s_{t+1}) - V(s_t) \\ V(s) &\leftarrow V(s) + \alpha \delta_t E_t(s)\end{aligned}$$

when $\alpha = 0$, it reduces to TD(0)

$$E_t(s) = 1(s_t = s)$$

~~and $E_t =$~~ and $E_t(s) = 0$ if $s_t \neq s$



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

as $\lambda = 0$,

$$\Rightarrow V(s) \leftarrow V(s_t) + \alpha \delta_t$$

If $\alpha = 1 \rightarrow$ we get same behaviour as MC