

RL - Week - 6

N-Step Prediction & TD-lambda

TD(0)

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

only next step reward,

$$G_t^{(1)} = r_{t+1} + \gamma V_t(s_{t+1})$$

1 step truncated corrected returns

Obv. n step

$$V(s_t) \leftarrow V(s_t) + \alpha [r_{t+1} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n V(s_{t+n}) - V(s_t)]$$

similarly $G_t^{(n)}$ → consider rewards till the end of the episode (MC)

beginning the value funcⁿ & returns could be wrong, it gets fine later by discounting over eventually.

TD(1 step) ——— n step MC



- Instead of n -step, let's consider the average of n -step returns.

$$G_t^{\text{avg.}} = \frac{1}{2} G_t^{(2)} + \frac{1}{2} G_t^{(4)}$$

$$\text{or } 0.8 \quad " \quad + 0.2 \quad "$$

This is $\text{TD}(\lambda)$. The average contains all the n -step backups each weighted α to λ^{n-1} , $0 \leq \lambda \leq 1$

$$\lambda\text{-return} \Rightarrow G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

it combines

1 -step to MC.

length of episode

$$\downarrow \quad \quad \quad \downarrow \quad T-t-1$$

$$1-\lambda \quad \quad \quad \lambda$$

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

$$G_t^{(\lambda)} = G_t + \lambda^{T-t-1}$$

so far we saw the forward view of $\text{TD}(\lambda)$

- Eligibility Traces (Backward view)

↳ helps to implement $\text{TD}(\lambda)$

These are traces related to each state given by $e_t(s)$. It is indicated the degree to which state is eligible for

learning changes.

$$e_t(s) = \begin{cases} r + e_{t-1}(s) & s \neq s_t \\ \gamma e_{t-1}(s) + 1 & s = s_t \end{cases}$$

$\nexists s \in S$, γ is discount if eligible rate.

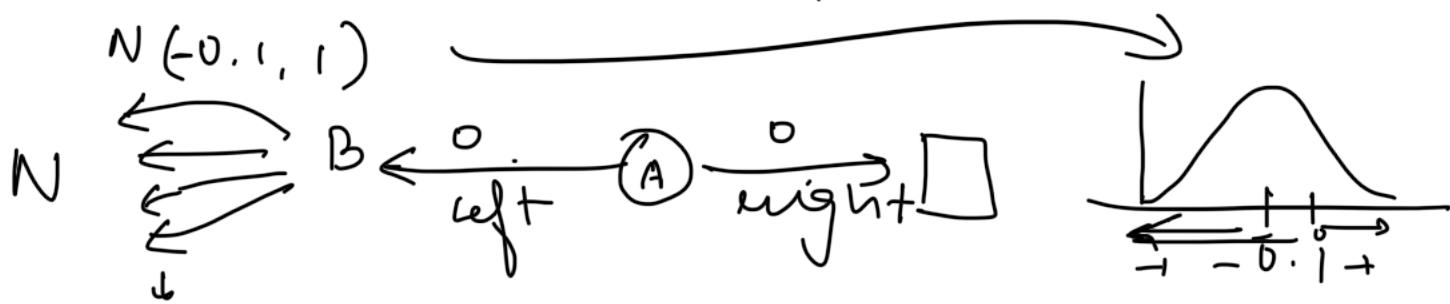
we accumulate the eligible traces where $s = s_t$.

Initially, $e(s) = 0 \rightarrow$ eligibility is initialized.

- both viewers accumulate the same updates. They are equivalent, bit a hand show.
- some eq.

L2 Double Q-learning

- we are actually estimating q func". The target uses the current max. So, therefore, it will propagate errors.
↳ maximian bias.



10,000 possibility

where to move from A

$$\checkmark E[G_t | s_0 = A, a_0 = \text{right}] = 0$$

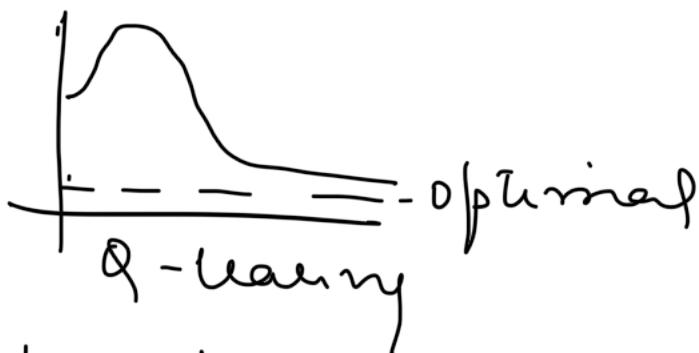
$$E[G_t | s_0 = A, a_0 = \text{left}] = 0 + (-0.1) = -0.1$$

when we use Q-learning.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a)]$$

few acts have + value, then it cause to go for left because of max, so it always overestimates.

leads to the bias, called maximum bias. All approaches with max provides this.



but we want



- the problem is using the same samples both to determine the max. act & to estimate the value.

→ Soln → use diff. estimates for maximising the act & estimate the value.

with $1/2$ probab.

$$Q_1(s, a) \leftarrow Q_1(s, a) + \alpha (R + \gamma Q_2(s', \operatorname{argmax}_a Q_1(s', a)) - Q_1(s, a))$$

similarly $Q_2(s, a) \propto \gamma Q_1$,

terminates $\Rightarrow Q_1 \approx Q_2 \approx q^*$

it reduces maximization bias
and not removes it.

