

## LECTURE 4 - MODEL FREE PREDICTION.

$$\text{PLANNING} = \underbrace{\text{PREDICTION}}_{\text{EVALUATE } V_{\pi} \text{ FOR } \pi} + \underbrace{\text{CONTROL}}_{\text{FIND } \pi^* \text{ and } V_{\pi^*}}$$

GOAL: Given a Policy  $\pi$ , we want to EVALUATE the policy (Find  $V_\pi$ ) without knowledge of MDP.

### MONTE-CARLO METHODS.

- DOES NOT require knowledge of the dynamics of the MDP. instead, requires EXPERIENCE. (BEING ABLE to obtain SAMPLES).

- ONLY WORKS FOR EPISODIC TASKS that terminate no matter what actions we take.

$$\cdot \text{We wish to estimate } V_\pi(s) = \mathbb{E}_{\substack{T \sim \pi \\ \text{TRAJECTORY SAMPLED FROM } \pi}} [R_{t+1} + \gamma R_{t+2} + \dots | S_t = s].$$

WE ESTIMATE  $V_\pi(s)$  BY EMPIRICAL AVG.

### FIRST-VISIT MONTE CARLO ON POLICY EVALUATION

$N(s)$ : #times we visit state  $s = 0$ .

$$G(s) = 0$$

loop:

- SAMPLE EPISODE  $i$ :  $(S_{i0}, A_{i0}, R_{i0}, S_{i1}, A_{i1}, \dots, R_{iT}, S_{iT})$ .

- Define  $G_{i,T} = r_{i,T} + \gamma r_{i,T+1} + \dots$  : Return from time  $T$  onwards

- For Each state visited in Episode  $i$ :

- For first time step  $t$  where  $s$  is visited:

- $N(s) \leftarrow N(s) + 1$
- $G(s) \leftarrow G(s) + R_{it}$
- $V^\pi(s) \leftarrow \underline{G(s)}$

- N(s)
- this Estimator is
    - { UNBIASED.  $E[V^*(s)] = E[G_t | s_t = s]$
    - { Consistent. By Law of Large Numbers.  $V^*(s) \rightarrow E[G_t | s_t = s]$   
as  $N(s) \rightarrow \infty$ .

### EVERY-VISIT MC ON POLICY EVALUATION.

- We USE the SAME IDEA but we consider Every sample Return we see in an EPISODE in our AVERAGE.
- this Estimator is Biased (Returns are not IID)
  - { is consistent.

### INCREMENTAL MC ON POLICY EVALUATION.

- WE DO NOT NEED to wait UNTIL the END OF ALL EPISODES to obtain the Estimate of  $V^*(s)$ . We can UPDATE our ESTIMATE AFTER EVERY EPISODE.

Suppose we want to compute the mean of  $x_1, \dots, x_m$  incrementally.

$$\begin{aligned} \bar{x}_m &= \frac{1}{m} \sum_{i=1}^m x_i = \frac{1}{m} \left[ x_m + \sum_{i=1}^{n-1} x_i \right] = \\ &= \frac{1}{m} \left[ x_m + (n-1) \bar{x}_{n-1} \right] = \frac{1}{m} \left[ x_m + m \bar{x}_{n-1} - \bar{x}_{n-1} \right] \\ &= \frac{1}{m} x_m + \bar{x}_{n-1} - \frac{\bar{x}_{n-1}}{m} = \frac{1}{m} \left[ x_m - \bar{x}_{n-1} \right] + \bar{x}_{n-1} \end{aligned}$$

$\therefore \bar{x}_m = \underbrace{\bar{x}_{n-1}}_{\substack{\text{new} \\ \text{Estimate}}} + \underbrace{\frac{1}{m} \left[ x_m - \bar{x}_{n-1} \right]}_{\substack{\text{OLS} \\ \text{Estimate} \\ \text{S.E.} \\ \text{size}}}$	between what you thought the mean was and what it actually saw.
---	---

- therefore, we can UPDATE our Estimate of  $V^*(s)$  AFTER EACH EPISODE!

loop:

- SAMPLE EPISODE  $i$ :  $(S_{i0}, A_{i0}, R_{i0}, S_{i1}, A_{i1}, \dots, R_{iT}, S_{iT})$ .
- Define  $G_{i,\tau} = r_{i,\tau} + \gamma r_{i,\tau+1} + \dots$  : Return from time  $\tau$  onwards
- For Each state visited in Episode  $i$ :
  - For first time step  $\tau$  where  $s$  is visited:
    - $N(s) \leftarrow N(s) + 1$

#### UPDATE ESTIMATE

$$\boxed{\cdot V^{\pi}(s) \leftarrow V^{\pi}(s) + \frac{1}{N(s)} [G_{i,\tau} - V^{\pi}(s)]}$$

Running mean: Non-stationary domains. (Policy changes over time)

. In some cases, we want to 'forget' older data.  
Therefore, we would update as follows:

$$V^{\pi}(s) \leftarrow V^{\pi}(s) + \lambda [G_{i,\tau} - V^{\pi}(s)]$$

if  $\lambda > \frac{1}{N(s)}$ , we are giving more importance to  $G_{i,\tau}$   
and 'forgetting' old values.

#### CONS OF MC METHODS:

- usually high variance  $\rightarrow$  requires lots of data.
- only useful for episodic MDPs.

#### TEMPORAL DIFFERENCE LEARNING (TD)

- BOOTSTRAPS (USES ESTIMATES to compute new Estimate)
- SAMPLES
- UPDATES  $V^{\pi}(s)$  AFTER every step.  $(s, a, r, s')$
- USEFUL in EPISODIC OR IMPERFECT-HORIZON MDPs.

Aim: Estimate  $V^{\pi}(s) = \mathbb{E}_{\pi} [R_{T+1} + \gamma V^{\pi}(s_{t+1}) | s_t = s]$

Aim: Estimate  $V_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s]$   
 By Sampling Returns from  $\pi$ . TNT.

$$V(s) \leftarrow V(s) + \underbrace{\lambda [R_{t+1} + \gamma V(s_{t+1}) - V(s)]}_{\text{TD TARGET}} \quad \text{TO ERROR}$$

- We are not sampling one FULL Return  $G_t$ , instead we are ONLY SAMPLING 1 STEP, AND USING our current estimate  $V(s_{t+1})$ .

- $V_{\pi}(s) = \mathbb{E}[G_t | s_t = s] \Leftrightarrow = \mathbb{E}[R_{t+1} + \gamma V(s_{t+1}) | s_t = s] \Leftrightarrow$
- In Monte Carlo we use  $G_t$  as target, and in TD-Learning we use  $R_{t+1} + \gamma V(s_{t+1})$  as a target.
- It can be shown, under certain assumptions, that TD will converge to  $V_{\pi}(s)$  for a fixed  $\lambda$ .

### Adv/Dissav Monte-Carlo / TD-Learning.

- MC is unbiased, because we are taking an average over returns sampled from the true distribution.

- MC has high variance, because in each trajectory we have noise, and it stacks up over time steps

- TD has high bias because we bootstrap.

- TD has low variance. Only have 1 step of noise

- Monte-Carlo is more DATA-EFFICIENT than TD-Learning !!

Why? → Let's say we are in state  $s$  and we follow a trajectory of length  $L$  till the end of the EPISODE. A reward is only given at the end of the EPISODE. WI/Monte-Carlo, we immediately propagate that info back to our estimate of  $V_{\pi}(s)$ , because we consider FULL returns  $G_t$ .

However, WI/TD it might take  $L$  samples over that trajectory to propagate back the reward to state  $s$ .

- MC is EPISODIC, while TD isn't

- MC UPDATES ONLY AFTER EVERY EPISODE. TD REKS EACH STEP.
- MC DOES NOT USE MARKOV-ASSUMPTION. TD DOES.

## BATCH MC AND TD LEARNING - / DATA EFFICIENCY

### BATCH UPDATING.

- SUPPOSE ONLY FINITE AMOUNT OF EXPERIENCE  
↳ EPISODES / DOO TIME STEPS
- PROCESS BATCH  $\Rightarrow$  UPDATE VALUE  $\Rightarrow \dots \Rightarrow$  CONVERGE!
- WE WILL GO THROUGH THE BATCH, APPLY MC OR TD,

But will only update the value function at the END OF THE BATCH. REPEAT UNTIL CONVERGENCE.

But, What do these algorithms converge to in  
the batch settings?

- BATCH MC CONVERGES TO  $V_T$  THAT FITS THE ACTUAL DATA THE BEST. I.e.  $V_T(S) = \frac{1}{m} \sum_{i=1}^m G_i$ .

$$V_T(S) = \text{argmin}_\gamma \frac{1}{n} \sum_{i=1}^m (G_i - \gamma)^2$$

Returns seem for state  $S$   
MSE. ( $V_T$  minimizes MSE on DATA)

- BATCH TD CONVERGES TO THE VALUE-FUNCTION OF THE MAXIMUM-LIKELIHOOD MODEL FOR THE MDP.