

# CMPUT 365: Temporal Difference Methods for Prediction

Rupam Mahmood

Feb 28, 2022

## Admin

Due dates C2M2:

- Practice quiz: Tue Mar 1
- Graded notebook: Sat Mar 5

Midterm:

- Marks and feedback are on eclass and assign2

## Review of MC

- MC methods come with a huge advantage: model-free estimation of value functions
- Disadvantage is that they make per-episode updates, not every time step
- They are computationally and memory-wise expensive and cumbersome

### First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy  $\pi$  to be evaluated

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$

$Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$ :

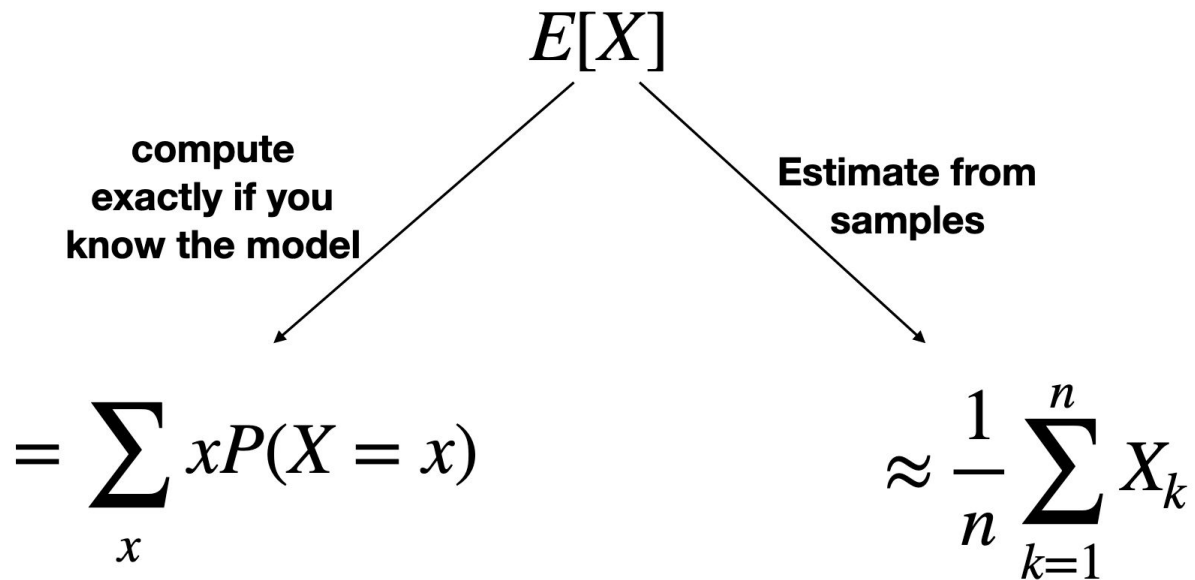
Append  $G$  to  $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

## Prediction as estimating value functions

- Predictions are building blocks for many control methods
- The usefulness of predictions goes beyond control
- Forming a predictive question: How many times will you get honked at today?
- (Pseudo-) reward: +1 for each honk
- Termination of episode: end of the day
- Behavior: the way you drive (think of your average speed, frequency of changing lanes, etc.)
- Can be answered by estimating  $v_{\pi}(s) \doteq E_{\pi}[G_t | S_t = s]$

Much of prediction is about estimating expected values



Much of prediction is about estimating expected values (cont'd)

$$v_{\pi}(s) \doteq E_{\pi}[G_t | S_t = s]$$

**Dynamic  
programming**

**E.g., Iterative  
policy evaluation**

**Sample-based**

**Monte Carlo  
(MC)**

# From MC to TD(0)

Monte Carlo estimator for on-policy prediction:  $V(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} G_t}{|\mathcal{T}(s)|}$  *→ sample average.*

Incremental Monte Carlo estimator:  $V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} [G_t - V(S_t)]$

Constant- $\alpha$  MC:  $V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$

*The target (the return is unbiased from the true value)*

*The more important MC update !!*

*→ the MC error (only update when episode ends).*

TD(0):  $V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)] \Rightarrow$  TD update.

*value estimate*

*this is slightly incorrect, b/c  $V(S_{t+1})$  is true value  $\therefore G$  can't be unbiased !!*

*(not the true value since we don't know it!!)*

Parham Golestaneh Talaei 对所有人说 下午1:48

PG

$V(S')$  is the value of the next state that was computed during the last trajectory right? Just like the in place policy iteration that we use the previous state values  $\Rightarrow$  Yes!!

# Unlike Monte Carlo, TD(0) works online

## First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy  $\pi$  to be evaluated

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$

$Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Unless  $S_t$  appears in  $S_0, S_1, \dots, S_{t-1}$ :

Append  $G$  to  $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

How would you  
characterize the differences  
between MC and TD(0)  
algorithmically?

## Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated

Algorithm parameter: step size  $\alpha \in (0, 1]$

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop for each episode:

Initialize  $S$

Loop for each step of episode:

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

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

$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

until  $S$  is terminal



## Unlike Monte Carlo, TD(0) works online (cont'd)

### Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated

Algorithm parameter: step size  $\alpha \in (0, 1]$

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop for each episode:

Initialize  $S$

Loop for each step of episode:

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

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

$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

until  $S$  is terminal

doing learning (more calculation than MC) during the episode

Aaron Skiba 对所有人说

下午1:50

AS

So 'offline' just means not learning as it's computing?  $\rightarrow$  Yes.

Say an oracle gives us return  $G$  from future at each step. Replace  $R + \gamma V(S')$  with  $G$ . This is an online but acausal Monte Carlo method. Will it be first-visit or every-visit?

## Just like MC, TD(0) estimates the value function but differently

$$\text{TD}(0): V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

the TD target (green box)  
 the estimate (orange circle)  
 the TD temporal difference error (green arrow)

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \quad (\text{from (3.9)}) \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]. \end{aligned}$$

the MC Target (pink underline)

Roughly speaking, Monte Carlo methods use an estimate of (6.3) as a target, whereas DP methods use an estimate of (6.4) as a target. The Monte Carlo target is an estimate because the expected value in (6.3) is not known; a sample return is used in place of the real expected return. The DP target is an estimate not because of the expected values, which are assumed to be completely provided by a model of the environment, but because  $v_{\pi}(S_{t+1})$  is not known and the current estimate,  $V(S_{t+1})$ , is used instead. The TD target is an estimate for both reasons: it samples the expected values in (6.4) and it uses the current estimate  $V$  instead of the true  $v_{\pi}$ . Thus, TD methods combine the sampling of Monte Carlo with the bootstrapping of DP. As we shall see, with care and imagination

# 公式 policy $\rightarrow$ Value

MC prediction, for estimating  $V \approx v_\pi$

Input: a policy  $\pi$  to be evaluated

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in \mathcal{S}$

$Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow \gamma G + R_{t+1}$

Append  $G$  to  $Returns(S_t)$

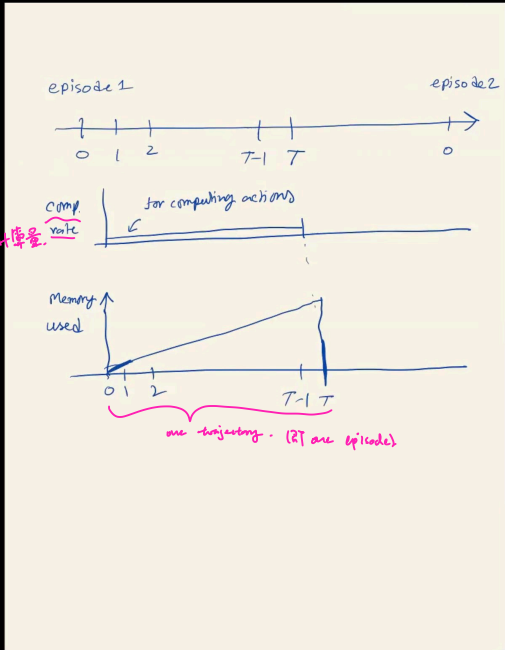
$V(S_t) \leftarrow \text{average}(Returns(S_t))$

after one episode is over and before the other episode is done.  
(Should the agent just wait until this computation is done?)

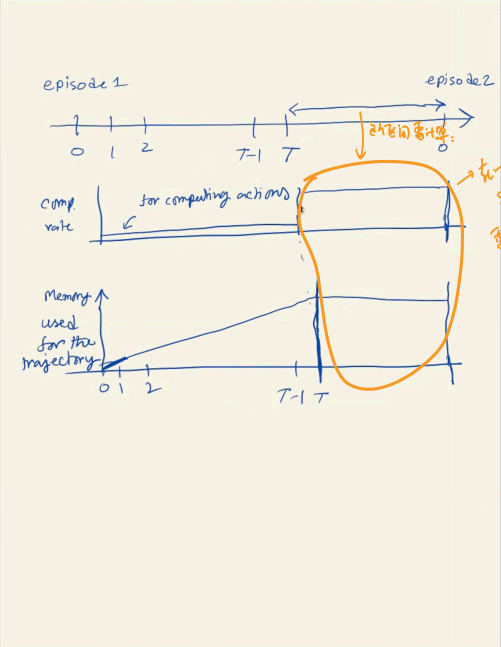
lots of computation!!  
expensive!!

the memory only stores one trajectory.  
at the end of the trajectory, you need to remember the whole trajectory.  
 $\Rightarrow$  still policy  $\pi$ .  
MC: 需要记住整个 episode 的所有信息。  
(when you generate this sequence, you are also remembering it.)

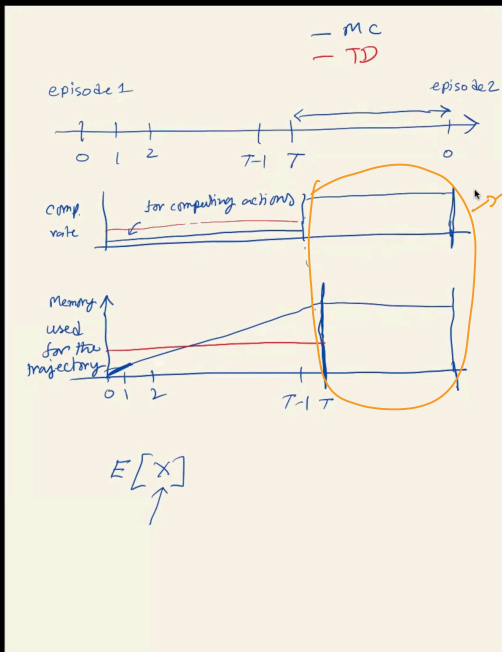
$\Rightarrow$  在已知  $P$  的情况下, 使用  $\pi$  来计算  $V$ -value.



正在发言: Rupam Mahmood



在每个 episode 结束时，  
第  $t$  个 episode 的总计算量  $G = G_t + R_{t+1}$ 。  
Append  $G_t$ 。  
 $V_{t+1} = \arg \max_{a \in A} (Q_t(s, a) + R_{t+1})$ 。  
 $\therefore$  所需计算量很大。



MC is better understanding than TD.