

# **Temporal Difference Methods**

Rupam Mahmood

February 24, 2020





# for Prediction



## Prediction as estimating value functions

Forming a predictive question:

- Predictions are building blocks for many control methods
  - The usefulness of predictions goes beyond control
    - 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



E.g., Iterative policy evaluation

## From MC to TD(0)

#### Monte Carlo estimator for on-policy pre-

#### **Incremental Monte Carlo estimator:** $V(S_t)$

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

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

**diction:** 
$$V(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} G_t}{|\mathcal{T}(s)|}$$

$$(A) \leftarrow V(S_t) + \frac{1}{N(S_t)} \left[ G_t - V(S_t) \right]$$

 $S_t$ )

$$S_t$$
)

# Unlike Monte Carlo, TD(0) works online

#### 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  $s \in S^+$ , arbitrarily except that V(terminal) = 0Loop for each episode: Initialize SLoop for each step of episode:  $A \leftarrow action given by \pi \text{ for } S$ Take action A, observe R, S' $V(S) \leftarrow V(S) + \alpha \left[ R + \gamma V(S') - V(S) \right]$  $S \leftarrow S'$ until S is terminal

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?

### From Monte Carlo error to TD error

**Constant-** $\alpha$  **MC:**  $V(S_t)$ 

#### **TD(0):** $V(S_t) \leftarrow V(S_t) + o$

Based on this equivalence, can you think of an implementation of Monte Carlo method that computes partial updates in an online manner and completes the full update at the end of the episode?



$$(t) \leftarrow V(S_t) + \alpha \quad [G_t - V(S_t)]$$
  
**MC error:**  $\Delta_t$ 

$$\alpha \left[ R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right]$$

#### **TD error:** $\delta_t$

$$\sum_{k=t}^{t-1} \gamma^{k-t} \delta_k$$