

## Temporal Difference Learning

Temporal difference learning is central to reinforcement learning, allowing for incremental updates to value estimates without waiting for the end of an episode.

The TD error, denoted as  $\delta_t$ , represents the difference between the estimated value of a state and the value of the next state.

$$\begin{aligned} V(S_t) &\leftarrow V(S_t) + \alpha[G_t - V(S_t)] \\ V(S_t) &\leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)] \leftarrow \delta_t \end{aligned}$$

Reminder:  $G_t \approx R_{t+1} + \gamma V(S_{t+1})$

### 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 the episode:

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

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

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

$S \leftarrow S'$

    Until  $S$  is terminal

### The advantages of Temporal Difference

Temporal Difference combines key ideas from Monte Carlo methods and Dynamic programming together, making it able to update from the previous parameters and learn directly from experience at once.

Feature	Monte Carlo (MC)	Dynamic Programming (DP)	Temporal Difference (TD)
Model Requirement	Does not need a model	Requires full model	Does not need a model
When It Updates	Only after episode ends	Iteratively over full state space	Step-by-step, during episode

<b>Use in Continuing Tasks</b>	Not suitable	Suitable	Suitable
<b>Works in Real-Time</b>	No (waits for episode to finish)	No	Yes
<b>Learning Efficiency</b>	Slower (uses full returns)	Fast if model is accurate	Faster (bootstraps estimates)
<b>Sample Efficiency</b>	Less efficient	High (requires full sweeps)	More efficient
<b>Bias/Variance</b>	Low bias, high variance	Low variance, higher bias	Moderate bias and variance
<b>Scalability</b>	Moderate (with enough samples)	Poor in large state spaces	High (learns from samples)
<b>Flexibility</b>	Works with unknown environments	Only in known environments	Works with unknown environments
<b>Responsiveness to Change</b>	Slower to adapt	Depends on full environment update	Quickly adapts with each step