

② Temporal difference

Issue: What happens if trajectories are too long?

Typical update rule:

$$\text{NEW} = \text{OLD} + \alpha (\text{CURRENT} - \text{OLD})$$

Example:

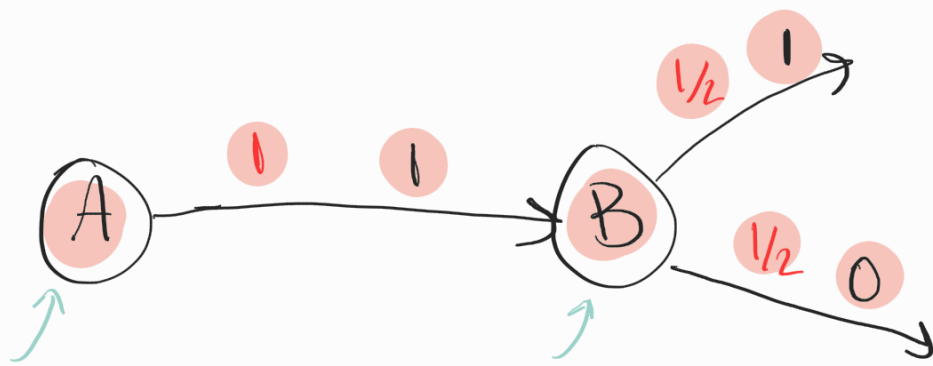
4 samples:

discount γ

A, B states

- $[A, 1, B, 1]$
- $B, 0$
- $B, 1$
- $B, 0$

MDP:



$$\text{val}(B) = 0.5$$

$$\begin{aligned}\text{val}(A) &= 1 \cdot (1 + \gamma \cdot \text{val}(B)) \\ &= 1 + \gamma \cdot 0.5\end{aligned}$$

temporal difference answer

Bootstrapping: estimating using estimates

Temporal difference update

Sample trajectories:

For each $s_0, a_0, \boxed{r_0}, s_1$:

$$\boxed{\widehat{\text{val}}(s_0)} = \widehat{\text{val}}(s_0) + \alpha \left(\underbrace{\boxed{r_0} + \gamma \widehat{\text{val}}(s_1)}_{\text{CURRENT}} - \widehat{\text{val}}(s_0) \right)$$

$$\boxed{\text{NEW}} = \text{OLD} + \alpha (\text{CURRENT} - \text{OLD})$$

Key advantage:

we update at each step (!)

Theorem: TD converge to $val(s)$

$$\hat{val}(s) \longrightarrow val(s)$$

Recap: two approaches for evaluating \hat{val}_π for π a policy:

Monte Carlo: requires full trajectories

Temporal Difference: step by step.