

Temporary Difference

Saturday, May 7, 2022 10:36 AM

→ Value iteration: requires P_{MT} and r_F to be known to compute either $V(s)$ or $Q(s,a)$ by approximation (Bellman Equations) In memory we needed: P, R and V/Q
 ↳ (Iteration)

→ Steps:
 1. 100 random transitions to compute P and R always an approximation (learned)
 2. With P and R , we do 1 value iteration we can do this 10000 times in the beginning
 3. Again 100 random to adjust P and R and then do all the value iterations together (same)
 ... repeat

Another way: Temporary Difference

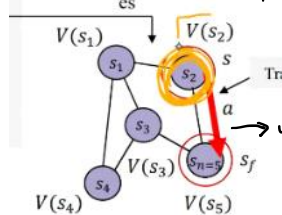
↳ Adjustment approximation per transition

→ Temporary Difference is another way of solving the Bellman Equations

↳ Where each $V(s)$ cell or $Q(s,a)$ row is calculated

means transition by transition:

↳ We modify only one $V(s)/Q(s,a)$ value during each transition (for only one state per transition/at a time).



advantage: we no longer need P_{MT} in advance (implicit in V/Q)

when we transition from s_2 to s_5 using action a : receives r and with r plus the transition $s_2 \rightarrow s_5$ we adjust $V(s_2)$

If we want the agent to learn the whole $V(s)$ the agent needs to experience all states (unlike in Value Iterations)
 { We adjust $V(s_4)$ when the agent is in $V(s_4)$ and transitions: to calculate the whole $V(s)$ if we got, say 10 states, we need 10 transitions, so that the ten $V(s)$ are updated.

each $V(s)$ is approximated step by step

a kind of value iteration

Since $V(s)$ doesn't tell us the action (unlike Q) and we need to choose the a with the reward function

needed

we don't need to know the reward function completely

↳ not needed

$V(s) \leftarrow V(s) + \alpha \Delta V(s)$
 ↳ error → increment in the direction of $V(s)$: in the direction that fixes the error
 ↳ decides how much of the error we consider
 ↳ α can be greater than 1
 ↳ $\alpha \in [0, \infty]$

where $\Delta V(s) = [f_R(s,a,s_f) + \gamma V(s_f)] - V(s)$
 ↳ also approximated
 ↳ what we think/expect V to have
 ↳ what we currently have

$$Q(s,a) \leftarrow Q(s,a) + \alpha \Delta Q(s,a)$$

$$\Delta Q(s,a) = [f_R(s,a,s_f) + \gamma \max_{a'} [Q(s_f, a')]] - Q(s,a)$$

↳ the immediate reward ↳ $Q(s,a)$ (approximated)

$\alpha \rightarrow$ defines the speed of learning

→ α small: slow learning, takes more time to reach desired value
 → α big: fast learning, takes less time to reach desired value, but if it reaches this value, it oscillates in a larger interval than it would with a small α .

Method

a) For $V(s)$ version $V(s) \leftarrow V(s) + \alpha \Delta V(s)$ where $\Delta V(s) = [f_R(s,a,s_f) + \gamma V(s_f)] - V(s)$
 ↳ The transition model function does not intervene, only reward func.

↳ needed

To Solve by Temporary Difference Iterations:

↳ This method makes the agent take random actions and we adjust $V(s)$ value per transition of the agent during those random actions OR we put random values in $V(s)$ and start to learn.

↳ when we compute $\max()$ in $V(s)$ to know the action, that \max is over all possible $V(s)$ values for all possible actions in the world.

↳ 2 actions $\max(V(a_1), V(a_2))$

