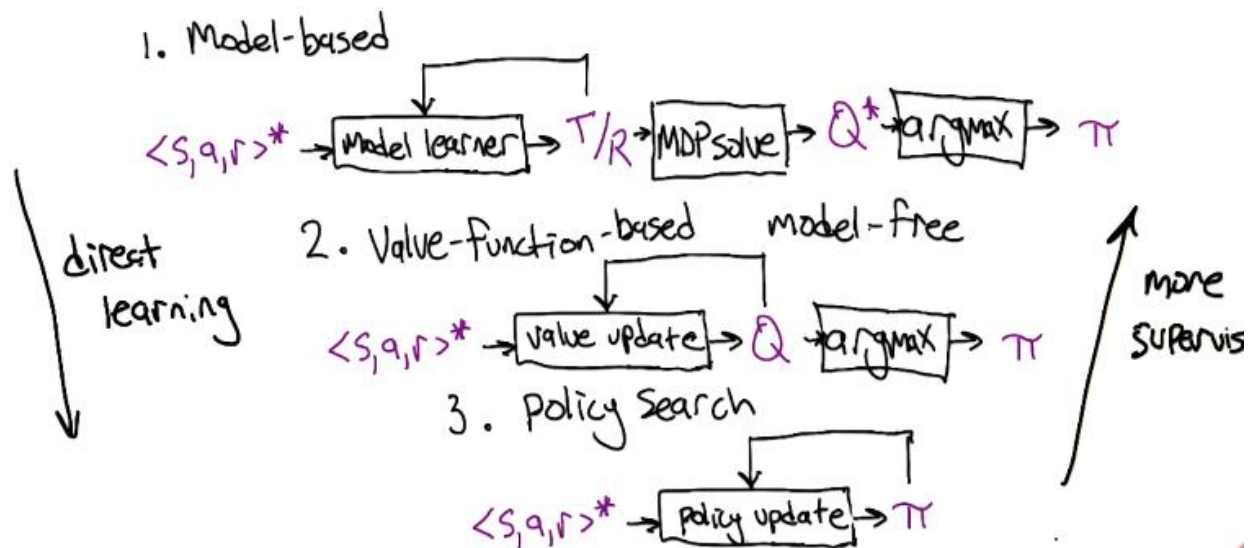


03. Temporal Difference Learning

Introduction:

- In Reinforcement Learning, the inputs are transitions (Initial state s , action a , reward r) and the intended output is to “learn” the policy π .
- Types of Reinforcement Learning algorithms:
 - Model-based: Sending the sequence of (state/action/reward) $\langle s, a, r \rangle^*$ to a Model Learner that learns the transitions and rewards. Then, an MDP solver can be used to infer an optimal value function Q^* from these transitions and rewards. Then, using argmax we can get the policy π . The model learner also takes into consideration the resulting transitions and rewards when producing new transitions and rewards.
 - Value Function based (model-free): Mapping states to values. Learning values from states is quite direct but turning this into a policy might be done in some sense using an argmax .
 - Policy Search: Taking the policy itself and feed it back to a policy update that modifies the policy according to the input $\langle s, a, r \rangle^*$. The problem with this approach is that the data doesn't reveal what actions need to be taken.



Temporal Difference Learning – TD(λ):

- TD Learning is learning to predict over time. We try to predict the expected sum of discounted rewards from a sequence of states where each transition has an associated reward.
- Computing estimates incrementally:

$$\begin{aligned}
 V_T(s_1) &= \frac{(T-1)V_{T-1}(s_1)R_T(s_1)}{T} \\
 &= \frac{(T-1)}{T}V_{T-1}(s_1) + \frac{1}{T}R_T(s_1) \\
 &= V_{T-1}(s_1) + \alpha_T(R_T(s_1) - V_{T-1}(s_1)) \quad \text{where } \alpha_T = \frac{1}{T}
 \end{aligned}$$

- This equation means that the update to the value is equal to the difference between the rewards we got at this step and the value we had at the previous step, multiplied by the learning rate α_T . α_T is getting smaller with higher number of episodes.

Properties of Learning Rates:

$$V_T(s_1) = V_{T-1}(s_1) + \alpha_T(R_T(s_1) - V_{T-1}(s_1)) \quad \text{where } \alpha_T = \frac{1}{T}$$

- If the number of episodes T is big enough ($\sim \infty$), the estimated value $V_T(s)$ will go to the true value $V(s)$:

$$\lim_{T \rightarrow \infty} V_T(s) = V(s)$$

- So, we have two conditions over the learning rate sequence:

$$\begin{aligned}
 \sum_T \alpha_T &= \infty \\
 \sum_T \alpha_T^2 &< \infty
 \end{aligned}$$

- If the learning rate doesn't get smaller over time, we'll never converge. This is why α_T should always be a value associated with T .

TD(1) Update Rule:

- Pseudocode for TD(1):

Episode T

For all s , $e(s) = 0$ at the beginning of the episode, $V_T(s_1) = V_{T-1}(s_1)$

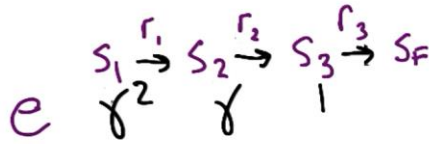
After $s_{t-1} \xrightarrow{r_t} s_t$: (step t)

$$e(s_{t-1}) = e(s_{t-1}) + 1$$

For all s ,

$$\begin{aligned}
 V_T(s) &= V_T(s) + \alpha_T(r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))e(s) \\
 e(s) &= \gamma e(s)
 \end{aligned}$$

TD(1) Example



$$\begin{aligned} \Delta V(s_1) &\propto (r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 V_{T-1}(s_F) - V_{T-1}(s_1)) \\ \Delta V(s_2) &\propto (r_2 + \gamma r_3 + \gamma^2 V_{T-1}(s_F) - V_{T-1}(s_2)) \\ \Delta V(s_3) &\propto (r_3 + \gamma V_{T-1}(s_F) - V_{T-1}(s_3)) \end{aligned}$$

Episode T

For all s , $e(s) = 0$ at start of episode, $V_T(s) = V_{T-1}(s)$

After $s_{t-1} \xrightarrow{r_t} s_t$: (step t)

$$e(s_{t-1}) = e(s_{t-1}) + 1$$

For all s ,

$$V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))e(s)$$

$$e(s) = \gamma e(s)$$

CLAIM:

TD(1) is the same as outcome-based updates (if no repeated states)

- The problem with TD(1) is that it's using individual episodes as opposed to using all the available data, and this can result in an estimate that is really far of the true value.
- To solve this issue, we investigate the TD(0) rule.

TD(0) Update Rule:

- Pseudocode for TD(0):

Episode T

For all s , at the beginning of the episode, $V_T(s_1) = V_{T-1}(s_1)$

After $s_{t-1} \xrightarrow{r_t} s_t$: (step t)

$$V_T(s_{t-1}) = V_{T-1}(s_{t-1}) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))$$

$$V_T(s_{t-1}) = \mathbb{E}_{s_t} [r + \gamma V_T(s_t)]$$

TD(λ) Update Rule:

- This is a general rule of TD(1) and TD(0) where $0 \leq \lambda \leq 1$
- Pseudocode for TD(λ):

Episode T

For all s , $e(s) = 0$ at the beginning of the episode, $V_T(s_1) = V_{T-1}(s_1)$

After $s_{t-1} \xrightarrow{r_t} s_t$: (step t)

$$e(s_{t-1}) = e(s_{t-1}) + 1$$

For all s ,

$$V_T(s) = V_{T-1}(s) + \alpha_T (r_t + \gamma V_{T-1}(s_t) - V_{T-1}(s_{t-1}))e(s)$$

$$e(s) = \lambda \gamma e(s)$$

k -step Estimators:

- It is useful to learn $TD(\lambda)$ by studying things that are not $TD(\lambda)$, like k -step estimators.
- We estimate the value of a state s_t as a function of what happens next.
- E_1 : Move it a little in the direction of the immediate reward plus the discounted estimated value of the state that we landed in.
 $TD(0) \rightarrow$ One-step estimator.
- E_2 : Move it a little in the direction of the immediate reward we received, plus the discounted reward that we received next, plus the double discounted value of the state we landed in two steps from now. Two-step estimator.
- E_∞ looks like $TD(1)$
- $TD(\lambda)$ is a weighted combination of all these estimators.

$$\begin{aligned}
 E_1 \quad V(s_t) &= V(s_t) + \alpha_T (r_{t+1} + \gamma V(s_{t+1}) - V(s_t)) \quad TD(0) \\
 E_2 \quad V(s_t) &= V(s_t) + \alpha_T (r_{t+1} + \gamma r_{t+2} + \gamma^2 V(s_{t+2}) - V(s_t)) \\
 E_3 \quad V(s_t) &= V(s_t) + \alpha_T (r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V(s_{t+3}) - V(s_t)) \\
 E_K \quad V(s_t) &= V(s_t) + \alpha_T (r_{t+1} + \dots + \gamma^{K-1} r_{t+K} + \gamma^K V(s_{t+K}) - V(s_t)) \\
 E_\infty \quad V(s_t) &= V(s_t) + \alpha_T (r_{t+1} + \dots + \gamma^{K-1} r_{t+K} + \dots - V(s_t)) \quad TD(1)
 \end{aligned}$$

k -step Estimators and $TD(\lambda)$:

- We will use weighted combination of all the estimators for any one of the $TD(\lambda)$ algorithms.
- When we change the value of a state, we are going to change it by all the different estimators. We will use a weighted combination of all of them each time that we do an update for a value of the state.
- We need to make sure that all the values coming from the estimators add up to one:

$$\begin{aligned}
 &\sum_{k=1}^{\infty} \lambda^{(k-1)} (1 - \lambda) \\
 &= (1 - \lambda) \sum_{k=1}^{\infty} \lambda^{(k-1)} = (1 - \lambda) \frac{1}{(1 - \lambda)} = 1
 \end{aligned}$$

$$\begin{aligned}
 E_1 &1 - \lambda \\
 E_2 &\lambda (1 - \lambda) \\
 E_3 &\lambda^2 (1 - \lambda) \\
 E_K &\lambda^{K-1} (1 - \lambda) \\
 &\dots \\
 E_\infty &\lambda^\infty
 \end{aligned}$$

- If $\lambda = 1$, we put all our weights on the last estimator, so we get TD(1).
- If $\lambda = 0$, we put all our weights on the first estimator, so we get TD(0).
- If $0 < \lambda < 1$, we get weights for all the estimators, and we combine all of them together.

TD(λ) Empirically:

- What happens if TD(λ) is run as an update algorithm on some finite amount of data?
- We run with different possible values of λ .
- Usually the error goes down as λ goes to some intermediate value and then starts to shoot up again and eventually hit error for TD(1) as λ approaches one.
- This is the justification on why we want to consider TD(λ), it gives you benefit of both TD(0) and TD(1). We are rapidly propagating information and not getting a biased estimate because of the way we are using 1 step information.