# Theory Answers

## Solution!

Initialise:

$\pi(s) \in \beta(s) \quad \forall s \in S$

$Q(s, a) \in R, \quad \forall s \in S, \quad a \in \beta(s)$

Returns $(s, a) \leftarrow$ empty list $\quad \forall s \in S, \quad a \in \beta(s)$.

reward-t $\leftarrow$ be zero array $\quad \forall s \in S, \quad a \in \beta(s)$

Loop for each episode:

Random choose $S_0, \beta_0$

Generate an episode $S_0, A_0,$ following $\pi$

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair $S_t, A_t$ appears

$Q(s, a) = Q(s, a) + 0.1 (G - Q(s, a))$

$\text{Returns}(s, a)$

$idx = np. \underset{a}{argmax} (Q(s, a))$

reward-t $[S] = $ reward-t$[S] + \frac{1}{episode} (G - \text{reward-t}[S])$
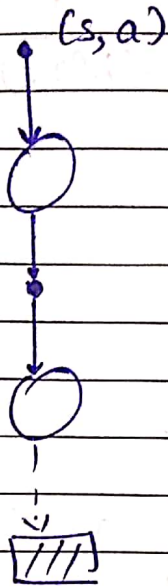
Calculating mean incrementally.

And appending to the list

~~reward-store $[S], $ episode~~

Returns $(s, a) \leftarrow$ ~~G~~ reward-t $[S]$

Thus $Q(s, a)$ is incrementally updated with returns for each state and returns take the mean value of $G$ and stores it, calculated incrementally.

**Solution 2**    Backup diagram for $q_\pi$

(s, a)

Problem is to find an estimate for $q_\pi(s, a)$

Start in state $s \rightarrow$ taking action $a$ and then follow a policy $\pi$.

**Solution 3**

We know $v_\pi(s) = \sum\limits_{a \in A(s)} q_\pi(s, a) \, \pi(a/s)$.

$$\frac{v_\pi(s)}{\sum\limits_{a \in A(s)} \pi(a/s)} = q_\pi(s, a)$$

$$\frac{1}{\sum b(a/s)} \frac{\sum b(a/s) \, v_\pi(s)}{\sum \pi(a/s)} = q_\pi(s, a)$$

$$\frac{v_\pi(s)}{\sum b(a/s) \, S_{t:T-1}} \quad \cancel{v_\pi(s)} = q_\pi(s, a)$$

$$v_\pi(s) = \sum q_\pi(s, a) \, S_{t:T-1} \, b(a/s).$$

$$q_\pi(s, a) = \frac{v_\pi(s)}{\sum b(a/s) \, S_{t:T-1}}$$

Thus $Q(s,a) = \dfrac{\sum_{t=T-1} G(t)}{|I(s)|} \cdot \dfrac{1}{\sum b(a/s)} \sum_{t=T}$

$$= \dfrac{G(t)}{|I(s)| \cdot \sum b(a/s)}$$

Thus $Q(s,a) = \dfrac{\sum G(t)}{b(a/s) \, |I(s)|}$

where $|I(s)|$ can be calculated as the number of times a state is visited for first visit or every visit.

## Solution 5 -

Starting with the problem described in question, once the user enter into a new building, some of the data encountered now will be same, . Thus starting with an initial guess, convergence can be achieved faster.

Similar way can be seen over packet delivery in a channel. In case of any delays in delivery, the updates can't be done as and when packets are delivered on the client side. Thus with an initial guess for fw values of these packet delay, the delay scene are the n/w can be seen faster.

# Solution 6

**Solution 6.3** — If the state value of A only changed in initial episode, it signifies that the game ended in the left terminal state with reward 0.

Since episode = 1, there's no room for exploration and considering the aim to reach terminal state, it ends in left terminal state with reward 0, having a short action benefit.

True value for state $A = \frac{1}{6} = 0.17$

Estimated value for state $A = 0.425$

% change = $\frac{0.425 - 0.17}{0.17} \times 100 = 25\%$

**Solution 6.4** — Had the step size parameter be wider, TD would perform better because on an average it will have more steps to explore before updating at next state and will push the agent towards the terminal state. Thus the algorithms would have definetly worked better had there been a wider range of $\alpha$. The improvement will be better in the initial phase but eventually there will be more fluctuations for RMSE.

| Notes | Appointment |
|---|---|
| | |

## Solution 6.5

The approximate value functions are initialized as 0.5 for all states and 0 for the terminal states.

$$RMSE = \sqrt{(true\ value - estimated\ value)^2}$$

with higher $\alpha$, the agent will be pushed more towards the terminal state in starting but eventually the fluctuations will start due to same value for all states and RMSE becomes more due to higher $\alpha$.

Large value of $\alpha$ implies more of $V(s)$ is updated at each timestep. Thus TD(0) becomes heavily dependent on the returns at each step of specific random walk sequence. At smaller $\alpha$, although learning takes more time to do but is affected by any perturbations at any time step.

## Solution 8

Suppose action selection is greedy, then post each time step,

$$action \leftarrow arg\ max\ Q(s', a)$$

which will be chosen for next episode. The action selection will be same always.