Question 1 (Exercise 5.4)                                          page - 1

$$Q_t(q,a) = R_{t+1} + \gamma R_{t+2}$$

$$Q_t(q,a) = (G_1 + G_2 - - + G_t)/t$$

$$Q_{t+1}(q,a) = (G_1 + - - G_t + G_{t+1})/t+1$$

$$= \frac{G_1 + - - + G_{t-1}}{t+1} + \frac{G_{t+1}}{t+1}$$

$$= \left(\frac{1}{t+1}\right)\left[t\cdot\left(\frac{G_1 + - - G_t}{t-1}\right) + \frac{G_{t+1}}{t+1}\right]$$

$$= \left(\frac{1}{t+1}\right)\left[t\cdot Q_t(q,a) + \frac{G_{t+1}}{t+1}\right]$$

$$Q_{t+1}(q,a) = \left(\frac{1}{t+1}\right)^2\left[t(t+1)Q(q,a) + G_{t+1}\right] \quad {\scriptstyle +Q_t(q,a) - Q_t(q,a)}$$

$$= Q_t(q,a) - \frac{(G_t - Q_t(q,a))}{t} \qquad \Big| \quad G_t \Rrightarrow R_t$$

$G_i \longrightarrow$ It represent return when state-action
pair is visit first in
and then discounted average return
from there.

$Q_T(q,a) \rightarrow$ Action - state value after T'th update
of state - action pair $(q,a)$

$t \rightarrow$ ~~It is t~~

$t \rightarrow$ It is number of times current
state - action pair has been updated.
Excluding current update

Initialize

$\pi(s) \in A(s)$

$Q(s,a) \in R$

Returns $(s,a) \leftarrow$ empty list, for all $s \in S$, $a \in A(s)$
$n(s,a) = 1$        $s, \in S, a \in A(s)$

~~Loop~~

Loop forever

     choose $S_0 \in S$ $A_0 \in A(S_0)$ randomly such that all pair have
     probability $> 0$

     Generate an episode from $S_0, A_0$ following $\pi : S_0, A_0, R_1, \dots$
                                      $S_{T-1}, A_{T-1}, R_T$

     $G \leftarrow 0$

     Loop for each step of episode, $t = T-1, T-2, \dots 0$

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

     Unless the pair $S_t, A_t$ appear in $S_0, A_0, S_1 \dots S_{t-1}, A_{t-1}$

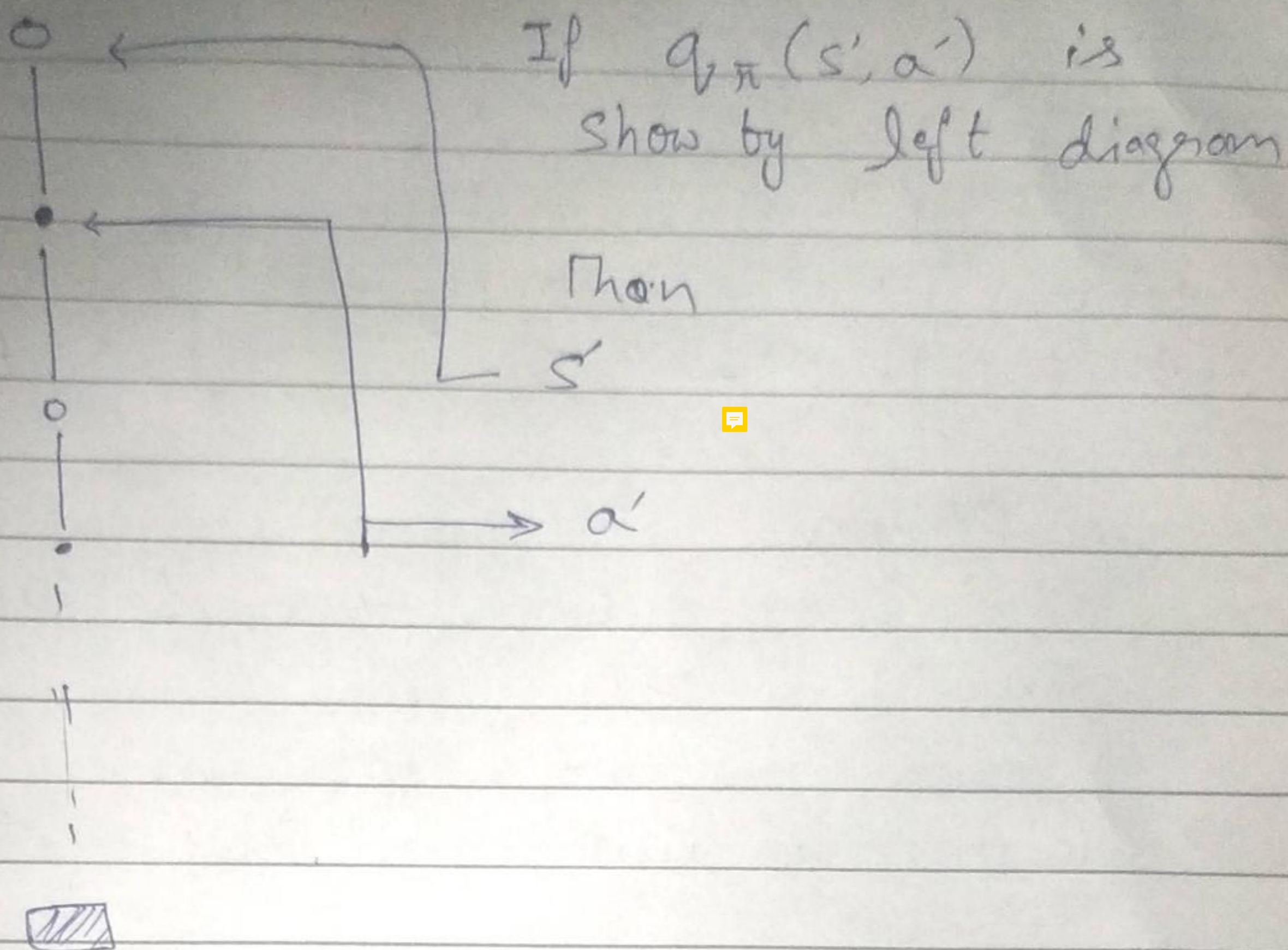     $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \dfrac{1}{n(s_t, a_t)} [G - Q(S_t, A_t)]$

     $n(s_t, a_t) = n(s_t, a_t) + 1$

It is single visit & Incremental approach

Question 2.
Exercise 5.3:-

If $q_\pi(s', a')$ is show by left diagram

Then

$\llcorner$ $s'$

$\rightarrow a'$

Question 5 (Exercise 6.2):-

Any task which is completely markov in nature.
  Like:-

Moves in chess. Given we are in certain state in chess and have experiance then we can predict our chances of wining without awaiting to game.
  End

Single visit MC :-

Assume at time $t$ we are in states $S_t$ and taking action $A_t$. Now, probability of subsequent trajectory

trajectory $\{(S_{t+1}, A_{t+1}) \quad\quad (S_{T-1}, A_{T-1}), S_T\}$

$$P_r\{ \downarrow | S_t, A_t; \pi\} = \prod p(S_{t+1} | S_t, A_t) \, p(A_{t+1} | S_{t+1})$$

$$= \prod p(S_k | S_{k-1}, A_k)$$

$$= \left[ \prod_{k=t}^{T-1} p(S_{k+1} | S_k, A_k) \, \pi(A_{k+1} | S_{k+1}) \right] p(S_T | S_{T-1}, A_{T-1})$$

So, Sampling Importance Ratio

$$= \prod_{k=t}^{T-1} \frac{\pi(A_{k+1} | S_{k+1})}{b(A_{k+1} | S_{k+1})}$$

Replace this Ratio in Eq of off-line state value.

$$v(s,a) = \frac{\sum_{t \in T(s,a)} P_{t:T(\theta-1)} \, G_t}{\sum_{t \in T(s,a)} P_{t:T(t)-1}} \quad\square\quad \tau(s,a)$$

mean in pair $(s,a)$

Consider $A_{k+1}$
Consider $b(A_t | S_t)$ only after taking action
$A$ in state $s$.

Question 8 (Exercise 6.12)

No,
   In Greedy SARSA behaviour Policy is dynamic.
But, In offline q-learning, policy is static.