1.) Input : a Policy $\pi$ to be evaluated

Initialize :

$V(s) = 0$ for all $s \in S$

$N(s) = 0$ for all $s \in S$

Loop forever (for each episode):

Generate an episode following $\pi$: $S_0, A_0, R_1, S_1, A_1, R_2 \cdots$

$S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

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

$N(S_t) \leftarrow N(S_t) + 1$

$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} \left( G - V(S_t) \right)$

2.a)

$\Rightarrow$ A State in **Blackjack** is the players total sum, the deleares visible card and whetha there is a usable ace.

Considering there are no splitting or doubling down, according to special rules.

It can be considered that there wont be much of a difference because the episodes are ended in the game of blackjack when the player wins or losses or got a draw.

lets consider a episode where the usa has

$(14, 2, \text{usable Ace})$ and if he hits

and got a 6 the total would be

$\textcircled{21}$

another episode $(K, 2, \text{usable Ace})$ and he strikes

the total would be

$\textcircled{21}$
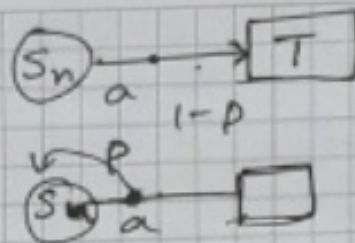
which is not usable as the episode is already happend and the game is won which will not impact the policy improvement as required.

1

2b)



The first visit only updates the returns for the non-terminal state only once

$$Q_{FVMC}(s, a) = average\ (10) = 10$$

But in case of every visit - MC it updates the non-terminal state every time it visits with an increased reward value.

$$\therefore Q_{EVMC}(s, a) = average\ (1+2+3 \cdots +10) = 5\cdot 5$$