# Question 6

## 6.3

Lets start from C

$$V(B) = 0.5 \qquad V(D) = 0.5$$

Reward is 0 for either side.

$$So \ V(C) = 0.5 + 0.1(0 + 0.5 - 0.5)$$
$$= 0.5 \qquad\qquad estimate \ remains \ same$$

Suppose we went to D
$$V(D) = 0.5 + 0.1(0 + 0.5 - 0.5)$$
$$= 0.5 \quad, \quad whether \ we \ went \ to \ C \ or \ E$$
$$estimate \ remains \ same$$

Suppose we went to E
reward is 1 for going to the terminal state and 0 otherwise.
If we went to the terminal state,

$$V(E) = 0.5 + 0.1(1 + \overset{0}{\cancel{0.5}} - 0.5)$$
$$V[E] = 0.55$$

But according to the figure, $V[E] = 0.5$ which means that this episode terminated at the bottom.

$$V[A] = 0.5 + 0.1(0 + 0 - 0.5)$$
$$= 0.5 - 0.05$$
$$= 0.45$$

Change is -0.05

64

Yes. For example, using $\alpha = 1$ would make Monte

MC at $\alpha = .01$ performs better than TD at $\alpha = 0.01$ ~~for these~~ For higher values of $\alpha$, TD seems to perform better than MC. This could be due to the fact that individual rewards don't affect $V(s)$ as much as episode returns.

6-5

This may be happening due to the whole stochasticity of the process, and due to the updates happening before the final return is generated

$$V(s) \leftarrow V(s) + \alpha (R + V(s') - V(s))$$

$V(s)$ is affected more when $\alpha$ is larger and R is fluctuating as well. This may be why we observe the increase in error. With smaller $\alpha$, the learning takes place more steadily, but slowly