

1. (30, 30%) Consider an undiscounted Markov Reward Process with two states  $A$  and  $B$ . The transition matrix and reward function are unknown, but you have observed two sample episodes:

$A + 3 \rightarrow A + 2 \rightarrow B - 3 \rightarrow A + 4 \rightarrow B - 4 \rightarrow \text{terminate}$

$B - 2 \rightarrow A + 3 \rightarrow B - 4 \rightarrow \text{terminate}$

In the above episodes, sample state transitions and sample rewards are shown at each step, e.g.,  $A + 3 \rightarrow A$  indicates a transition from state  $A$  to state  $A$ , with a reward of  $+3$ .

- (a) (5, 10%) Using first-visit Monte-Carlo evaluation, estimate the state-value function  $V(A), V(B)$ .

$$V(A) = \frac{(3 + 2 - 3 + 4 - 4)}{2} + \frac{(3 - 4)}{2}$$

$$= \frac{2 - 1}{2} = 0.5$$

$$V(B) = \frac{(-3 + 4 - 4)}{2} + \frac{(-2 + 3 - 4)}{2}$$

$$= \frac{-3 - 3}{2} = -3$$

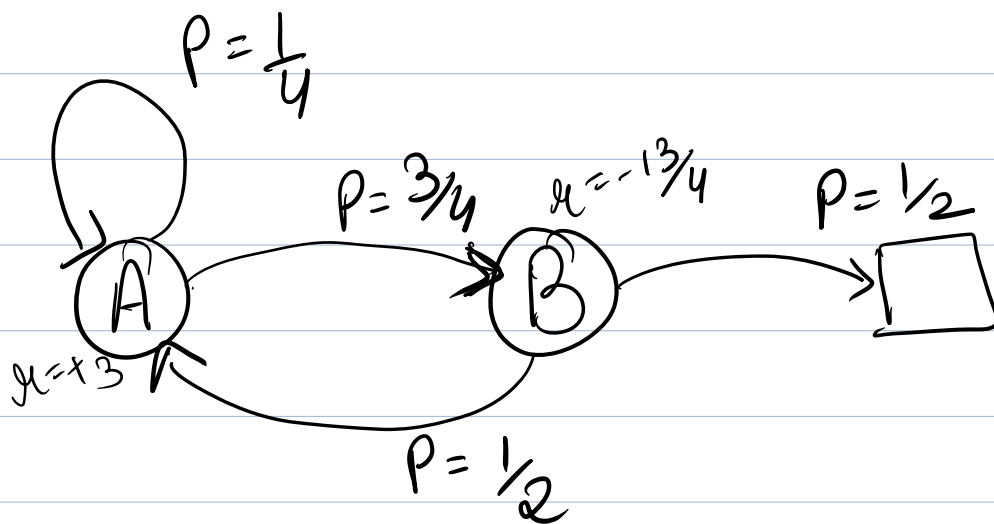
- (c) (5, 10%) Draw a diagram of the Markov Reward Process that best explains these two episodes (i.e. the model that maximises the likelihood of the data - although it is not necessary to prove this fact). Show rewards and transition probabilities on your diagram.

$A \xrightarrow{+3} A, A \xrightarrow{+2} B, B \xrightarrow{-3} A, A \xrightarrow{+4} B, B \xrightarrow{-4} \square \text{ (term)}$

$B \xrightarrow{-2} A, A \xrightarrow{+3} B, B \xrightarrow{-4} \square \text{ (term)}$

reward = avg??

In state  $A = 4$  times  $\frac{1}{4}$  back to  $A$ ,  $\frac{3}{4}$  to  $B$   
 In state  $B = 4$  times  $\frac{2}{4}$  terminates,  $\frac{2}{4}$  to  $A$



22  
23  
24

- (d) (5, 10%) Define the Bellman equation for your above Markov reward process. Solve the Bellman equation **directly**, rather than iteratively, to find the true state-value function  $V(A), V(B)$

$$V(s) = E \{ G_t | S_t = s \}$$

//

$$= E \{ r_{t+1} + V(s_{t+1}) | S_t = s \} \quad \text{undiscounted } (\gamma = 1)$$

$$= E \{ r_{t+1} | S_t = s \} + E \{ V(s_{t+1}) | S_t = s \}$$

$$= r_s + \sum_{s'} P(s' | s) V(s')$$

$$\therefore V(A) = 3 + \frac{1}{4} V(A) + \frac{3}{4} V(B) \quad - \textcircled{1}$$

$$\therefore V(B) = -\frac{13}{4} + \frac{1}{2} V(A) \quad - \textcircled{2}$$

①

$$V(A) - \frac{1}{4} V(A) = 3 + \frac{3}{4} V(B)$$

$$3V(A) = 12 + 3V(B)$$

$$V(A) = 4$$

②

$$V(B) = -\frac{13}{4} + \frac{1}{2} V(A)$$

$$\Rightarrow V(B) = -\frac{13}{4} + 2 + \frac{1}{2} V(B)$$

$$\Rightarrow \frac{1}{2} V(B) = -\frac{5}{4}$$

$$\Rightarrow V(B) = -\frac{5}{2} = \underline{\underline{-2.5}}$$

$$\therefore V(A) = 4 - 2.5 = \underline{\underline{1.5}}$$