

## Question 4:

We want to show that

$$\begin{aligned} V^{\pi}(s) &\triangleq E^{\pi} \left( \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right) \\ &= E^{\pi} \left( \sum_{t=0}^{\infty} \gamma^{t-1} r(s_t, a_t) \mid s_1 = s \right) \end{aligned}$$

We can prove this lemma by assumption that  $\pi$  is stationary and changing the sum index

$$\begin{aligned} E^{\pi} \left( \sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t) \mid s_1 = s \right) &= E^{\pi} \left( \sum_{t=0}^{\infty} \gamma^t r(s_{t+1}, a_{t+1}) \mid s_1 = s \right) \\ &= E^{\pi} \left( \sum_{t=0}^{\infty} \gamma^t (r(s_{t+1}, \pi(s_{t+1}))) \mid s_1 = s \right) = E^{\pi} \left( \sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) \mid s_0 = s \right) \\ &= V^{\pi}(s) \end{aligned}$$

## Question 2

1. We can see that

$$r(s_0, a_1) = E[\text{bernoulli}(0.2)] = 0.2$$

$$r(s_1, a_1) = E[\text{normal}(0, 1)] = 1$$

$$r(s_2, a_1) = 0.5$$

$$r(s_0, a_2) = 0.4 \quad r(s_1, a_2) = 0$$

$$r(s_2, a_2) = 0.5$$

Jack's expected reward after 3 rounds is:

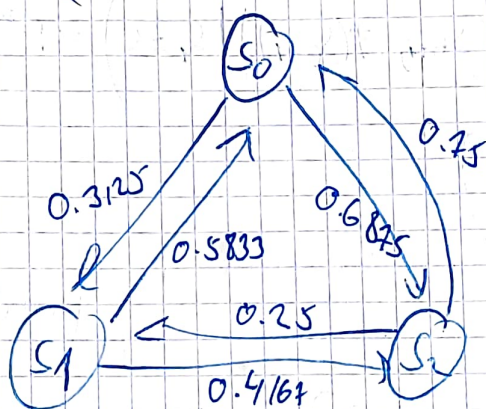
$$\begin{aligned} E^{\pi_{21}} \left[ \sum_{t=0}^3 R_t \right] &= r(s_0, a_2) + (0.125 r(s_1, a_1) + 0.875 r(s_2, a_1)) \\ &+ \left( (0.125 \cdot \frac{2}{3} + 0.875 \cdot 0.45) r(s_0, a_1) + 0.875 \cdot 0.25 r(s_1, a_2) \right. \\ &\quad \left. + 0.125 \cdot \frac{1}{3} r(s_2, a_2) \right) = 1.801 \end{aligned}$$

b. Now we get

$$\begin{aligned}
 P_{ij} &= P(s_{t+1} = s_j | s_t = s_i) = P(s_j | s_i, a_1) P(a_1) \\
 &\quad + P(s_j | s_i, a_2) P(a_2) \\
 &= 0.5 (P(s_j | s_i, a_1) + P(s_j | s_i, a_2))
 \end{aligned}$$

and with that we get

$$P = \begin{pmatrix} 0 & 0.3125 & 0.6875 \\ 0.5833 & 0 & 0.4167 \\ 0.45 & 0.25 & 0 \end{pmatrix}$$



Now the expected reward for each state is

$$r(s_i) = r(s_i, a_1) P(a_1) + P(a_2) r(s_i, a_2) = 0.5 (r(s_i, a_1) + r(s_i, a_2))$$

$$r = [0.45, 0.5, 0.5]$$

and the total expected reward is

$$E^{\pi} \left[ \sum_{t=0}^{\infty} R_t \right] = r(s_0) + (0.3125 r(s_1) + 0.6875 r(s_2))$$

$$\begin{aligned}
 &+ ((0.3125 \cdot 0.5833 + 0.6875 \cdot 0.75) r(s_0) + 0.6875 \cdot 0.25 r(s_1) \\
 &\quad + 0.3125 \cdot 0.4167 r(s_2)) = 1.415
 \end{aligned}$$



C. Let's write the Bellman equation for 3 rounds

$$V_3(s) = 0$$

$$V_k(s) = \max_{a \in \{a_1, a_2\}} \left\{ r(s, a) + \sum_{s' \in \{s_0, s_1, s_2\}} P(s'|s, a) \cdot V_{k+1}(s') \right\}$$

we get:

$$V_2(s_0) = \max\{0.2, 0.4\} = 0.4 \quad \{a_2\}$$

$$V_2(s_1) = \max\{1, 0\} = 1 \quad \{a_1\}$$

$$V_2(s_2) = \max\{0.5, 0.5\} = 0.5 \quad \{a_1/a_2\}$$

$$V_1(s_0) = \max\{0.2 + 0.5 \cdot 1 + 0.5 \cdot 0.5, 0.4 + 0.125 \cdot 1 + 0.875 \cdot 0.5\} \\ = 1.2625 \quad \{a_2\}$$

$$V_1(s_1) = \max\left\{1 + \frac{2}{3} \cdot 0.4 + \frac{1}{3} \cdot 0.5, 0 + 0.5 \cdot 0.4 + 0.5 \cdot 0.5\right\} \\ = 1.633 \quad \{a_1\}$$

$$V_1(s_2) = \max\{0.5 + 0.15 \cdot 0.4 + 0.25 \cdot 1, 0.5 + 0.15 \cdot 0.4 + 0.5 \cdot 1\} \\ = 2.2 \quad \{a_1/a_2\}$$

$$V_0(s_0) = \max\{0.2 + 0.5 \cdot 1.633 + 0.5 \cdot 2.2, 0.4 + 0.125 \cdot 1.633 + 0.875 \cdot 2.2\} \\ = 2.83 \quad \{a_2\}$$

and the optimal policy

$$\pi_0^*(s_0) = a_2 \quad \pi_1^*(s) = \pi_2^*(s) = \begin{cases} a_2 & s = s_0 \\ a_1 & s = s_1 \\ a_1 \text{ or } a_2 & s = s_2 \end{cases}$$

- d. The probability to stay in the casino after  $t$  rounds is  $(1-\beta)^t$ . hence the infinite horizon cumulative reward:

$$J_{\beta}^{\pi}(s) = E^{\pi, s_0} \left( \sum_{t=0}^{\infty} P(\text{stay after } t \text{ rounds}) \cdot R_t \right)$$

$$= E^{\pi, s_0} \left( \sum_{t=0}^{\infty} (1-\beta)^t R_t \right) = E^{\pi, s_0} \left( \sum_{t=0}^{\infty} (1-\beta)^t r(s_t, a_t) \right)$$

$\gamma \triangleq 1-\beta$  defines the connection between the discount factor and the death rate

- e. for the infinite horizon case we have the following bellman equations:

$$V(s) = \max_{a \in \{a_1, a_2\}} \left\{ r(s, a) + (1-\beta) \sum_{s' \in \{s_0, s_1, s_2\}} P(s'|s, a) \cdot V(s') \right\}$$

$$\pi^*(s) = \arg \max_{a \in \{a_1, a_2\}} \{ \dots \}$$