**1.**

(a) <u>State Space</u>: $X = \{0,1,2,3,4,11,12,22\}$ , which means the money we have, e.g. if $x=1$ then we have $\$10.000$. Then terminal states are $T = \{x \leq 0 \text{ or } x \geq 3 \mid x \in X\}$.

<u>Control Space</u>: $U = \{1(r), 0(b)\} \times \{m\} \subseteq \mathbb{R}^2$, where $1$ means betting on red and $0$ means betting on black. $m$ is integer and $m \geq 1$, this is the money we bet, e.g. if $m = 2$ then we bet $\$20.000$ this round.

<u>Motion model</u>: If $x_t \in X \setminus T$

$$x_{t+1} \mid x_t, u_t = (1, m) = \begin{cases} x_t + m, & \text{with prob. } 0.7 \\ x_t - m, & \text{with prob. } 0.3 \end{cases}$$

$$x_{t+1} \mid x_t, u_t = (0, m) = \begin{cases} x_t + 10m, & \text{with prob. } 0.2 \\ x_t - m, & \text{with prob. } 0.8 \end{cases}$$

If $x_t \in T$.

$$x_{t+1} \mid x_t, u_t = x_t, \quad \text{i.e.} \begin{cases} P_f(x_t \mid x_t, u_t) = 1 \\ P_f(x_{t+1} \neq x_t \mid x_t, u_t) = 0 \end{cases}$$

<u>Stage cost</u>: $\ell(x, u, x') = -(x' - x)$

i.e. Once we get into terminal costs, we get $-x'$ reward.

Then $\tilde{\ell}(x, u) = -\sum_{x'} P_f(x' \mid x, u)(x' - x)$

<u>Terminal cost</u>: $q(x_T) = -x_T$

(b) $V^*(x) = \min\limits_{u \in U(x)} \tilde{\ell}(x,u) + \mathbb{E}_{x' \sim \tilde{P}_f(\cdot|u,x)}[V^*(x')]$

In this case, $\pi = u_t = (1 \cdot 1)$. then
$$X = \{0, 1, 2, 3\}. \text{ terminal states } T = \{0, 3\}$$
Suppose at iteration $k=0$, $V_0^\pi(x) = 0$ for all $x \in X \backslash T$
$$V_0^\pi(x_T) = -x_T \text{ for all } x_T \in T.$$
When $k=1$.
$$V^\pi(x) = \tilde{\ell}(x,u,x') + \sum_{x' \in X} \tilde{P}_f(x'|x,u) V^\pi(x')$$
By program.
$$\overline{V}^\pi = [0, -0.4, -2.5, -3, -4, -11, -12, -22]^T$$
for $x = 0, \quad 1, \quad 2, \quad 3, \quad 4, \quad 11, \quad 12, \quad 22$
Also. we get precise estimate
$$\hat{V}^\pi = [0, -2.7215, -3.3165, -3, -4, -11, -12, -22]^T$$
for $x = 0, \quad 1, \quad 2, \quad 3, \quad 4, \quad 11, \quad 12, \quad 22$


(c) State space $X = \{0, 1, 2, 3, 4, 11, 12, 22\}, T = \{0, 4, 11, 12, 22\}$
Given $\overline{V}^\pi$. using computer we get
$$\overline{\pi}'(1) = (0, 1), \quad \overline{\pi}'(2) = (0,2), \quad \overline{\pi}'(x) = \{\text{go home}\}, x \in T.$$
which both means to bet on black and bet all money.


Given $\hat{V}^\pi$, using computer we get
$$\hat{\pi}'(1) = (0, 1), \quad \hat{\pi}'(2) = (0, 2), \quad \hat{\pi}'(x) = \{\text{go home}\} \; x \in T.$$
so if we have $\$10,000$. we bet on black and bet $\$10,000$;
if we have $\$20,000$. we bet on black and bet $\$20,000$.

$\bar{\pi}'$ and $\hat{\pi}'$ are the same. Compared to $\pi = (1, 1)$. the policy $\bar{\pi}'$, $\hat{\pi}'$ become bolder, trying to earn more money.

(d) Using Policy Iteration (PI). we'll have

$$V^*(1) = -6.6364 \quad , \quad \pi^*(1) = (1, 1) \longrightarrow \text{bet red. } \$10,000$$

$$V^*(2) = -8.9091 \quad . \quad \pi^*(2) = (0, 1) \longrightarrow \text{bet black. } \$10.000$$

$$V^*(x) = x, \quad \pi^*(x) = \{stay\} \quad \text{for } x \in 7.$$

Therefore. we bet red, $10.000 when we have $10,000;
and we bet black, $10,000 when we have $20,000.

Under this policy. I simulate the gambling for 50.000 times. and the average money I take home is $38.614.

Celebration!