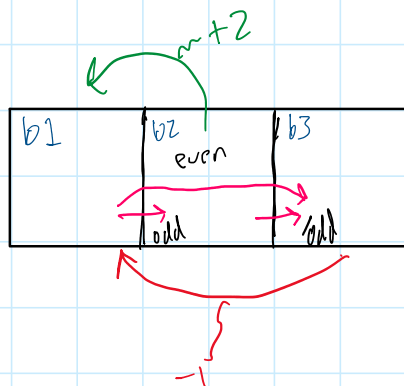1 a) The reward function represents the expected reward received from executing an action from a given state.
The state-action function Q is similar, but represents the expected reward from starting at a given state, and taking an action, and continuing with a greedy policy with respect to Q. So, Q represent expected total reward throughout a process, and the reward function only represents the reward for one step.

b)

$$Q(s,a) = R(s,a) + \gamma \sum_{s'} T(s'|s,a) V(s')$$

2)



a) $S = \{b_1, b_2, b_3\}$

$$A = \{ Roll, Reset \}$$

assuming fair die

$$T_{roll} = \begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix} \begin{bmatrix} 0 & 0.5 & 0.5 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{matrix} b_1 & b_2 & b_3 \end{matrix}$$

$$T_{reset} = \begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix} \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{matrix} b_1 & b_2 & b_3 \end{matrix}$$

$$T_{roll, unfair} = \begin{bmatrix} 0 & P & 1-P \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R(S, a, S') = \begin{cases} +2 & \text{if } a = reset \text{ and } S = b_2 \\ -1 & \text{if } a = reset \text{ and } S = b_3 \end{cases}$$

$$R = \begin{bmatrix} 0 \\ +2 \\ -1 \end{bmatrix} \begin{matrix} b_1 \\ reset\ b_2 \\ reset\ b_3 \end{matrix}$$

b)

$$U = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix}$$

☆ assuming $P = 0.5$ for this problem

Guessing $\pi = \begin{bmatrix} roll \\ reset \\ reset \end{bmatrix}$

$$U^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') U^{\pi}(s)$$

$$\left[ U^{\pi}(b_1) = 0 + \left( P \cdot U^{\pi}(b_2) + (1-P) U^{\pi}(b_3) \right) \right.$$

$$\left. U^{\pi}(b_2) = 2 + \gamma (U^{\pi}(b_1)) \right. \rightarrow$$

$$U^\pi(b_2) = 2 + \gamma(U^\pi(b_1)) \quad \rightarrow$$

$$U^\pi(b_3) = -1 + \gamma(U^\pi(b_1)) \quad \rightarrow$$

$$\rightarrow b_1 = \left( P(2 + \gamma b_1) + (1-P)(-1 + \gamma b_1) \right) \cdot \gamma$$

$$b_1 = (2P + \gamma P b_1 + \gamma b_1 - 1 + P - \gamma P b_1) \cdot \gamma$$

$$b_1 = 2P\gamma + \gamma^2 P b_1 + \gamma^2 b_1 - \gamma + P\gamma - \gamma^2 P b_1$$

$$b_1 - \gamma^2 b_1 = 2P\gamma - \gamma + P\gamma$$

$$b_1 = \frac{3P\gamma - \gamma}{1 - \gamma^2} \xrightarrow[\gamma=0.95]{P=0.5} b_1 = 4.872$$

$$U^\pi(b_1) = 4.872$$

play in $\rightarrow$ $U^\pi(b_2) = 6.63$

$$U^\pi(b_3) = 3.63$$

Check that these values satisfy $U^*(s) = \max\limits_{a}\left( h(s,a) + \gamma \sum\limits_{s'} T(s'|s,a) U^*(s') \right)$

$b_1 \rightarrow$ 
roll: $0 + 0.95(0.5(6.63) + 0.5(3.63)) = 4.874$ } max = 4.874
reset: $0 + 0.95(4.872) = $ less than $4.874$

$b_2 \rightarrow$ roll: $0 + 0.95(3.63) = $ less than $3.63$ } max = 6.63
reset: $2 + 0.95(4.872) = 6.63$

$b_3 \rightarrow$ roll: $0 + 0.95(3.63) = $ less than 3.63 } max $= 3.63$
reset: $-1 + 0.95(4.872) = 3.63$

- All of these values are the same, except for $s_1$. This optimal value is very close though, wholly resulting from a rounding error. So, the policy $\pi = \begin{Bmatrix} \text{roll} \\ \text{reset} \\ \text{reset} \end{Bmatrix}$ is $\underline{\text{optimal}}$

c) I will use the same policy as (b) because in the case of all evens, reward will be minimized by repeatedly reaching $b_3$ and resetting. Similarly, with all odd rolls reward will be maximized by resetting only on $b_2$, with $+2$ reward. I am able to compute bounds by using the expression found earlier for $U^\pi(b_1)$ and varying the probability $(\rho)$ term.
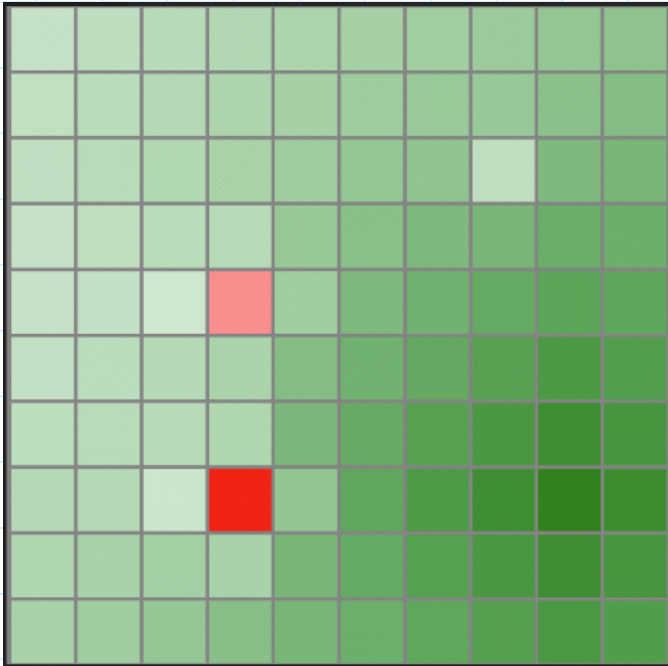
$$U^\pi(b_1) = \frac{3\rho\gamma - \gamma}{1 - \gamma^2}$$, where $\rho = 0$ simulates all-even rolls, and $\rho = 1$ simulates all-odd rolls.

$U^\pi(b_1)\big|_{\substack{\rho = 0 \\ \gamma = 0.95}} = \underline{-9.744}$

$U^\pi(b_1)\big|_{\substack{\rho = 1 \\ \gamma = 0.95}} = \underline{19.487}$

So, discounted score $\quad -9.744 \leq U(b_1) \leq 19.487$

3)



4) Passed autograder with n = 7