**1]** $\quad S = \{ S_1, S_2 \} \quad A = \{ \text{stay}, \text{go} \}$



$r = -2$ (around $S_1$), $r = -3$ ($S_1 \to S_2$), $r = +5$ ($S_2 \to S_1$) terminate, $r = -2$ (around $S_2$)

**a]** $\displaystyle \sum_{t=0}^{\infty} \gamma^t r_t(S_1, a_t) = \sum_{t=0}^{\infty} \gamma^t r_t(S_1, \text{stay})$

$\displaystyle = \sum_{t=0}^{\infty} \gamma^t (-2), \quad$ since $0 < \gamma \leq 1$, this is a geometric series.

$$= \boxed{\dfrac{-2}{1-\gamma}}$$

**b]** Consider $\Pi_e$ that goes to terminate ASAP.

i.e. $\Pi_e(S_1) = \text{go}, \quad \Pi_e(S_2) = \text{go}.$

Now consider its sum of discounted rewards.

$$\sum_{t=0}^{\infty} \gamma^t r_t(s_t, a_t) = \gamma^0 r_0(s_1, go)$$
$$+ \gamma^1 r_1(s_2, go)$$

$$= \underline{-3 + \gamma(5)}$$

Let's compare this with result from 1a. $(\pi_{1a})$

We can use this $\pi_e$ if the following holds.

$$-3 + \gamma(5) - \left( \frac{-2}{1-\gamma} \right) \geq 0$$

$$\frac{(-3 + 5\gamma)(1-\gamma) + 2}{\boxed{1-\gamma}} \geq 0$$

$\longrightarrow$ denominator is always $> 0$ since $\gamma < 1$.

So inequality depends on numerator.

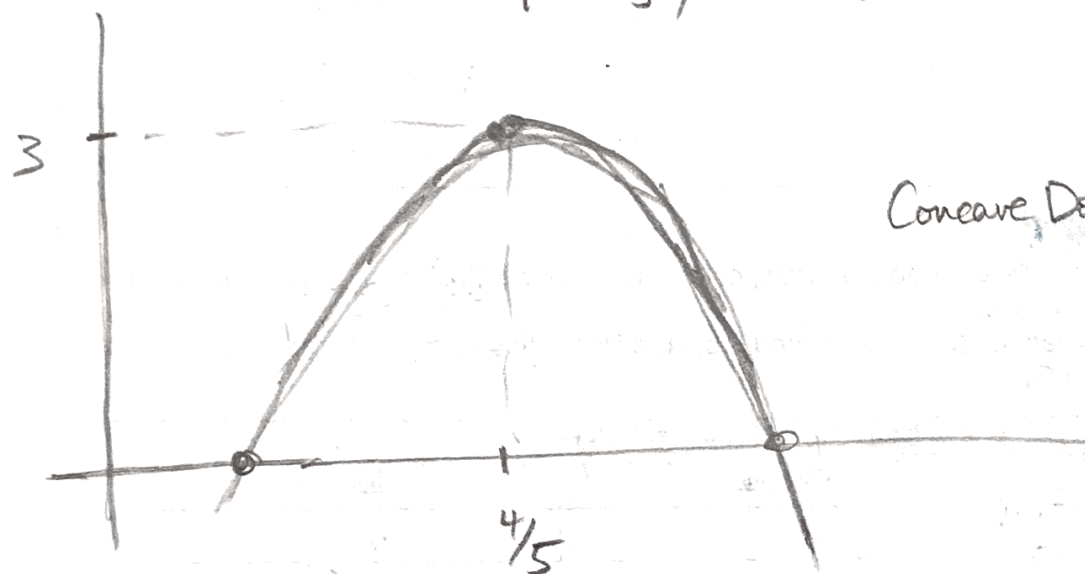numerator = $-3 + 5\gamma + 3\gamma - 5\gamma^2 + 2$

$$f(\gamma) = -5\gamma^2 + 8\gamma - 1$$

$$f(\gamma) = -5\left(\gamma^2 - \frac{8}{5}\gamma + \frac{4}{5}\right) + 3$$

$$= -5\left(\gamma - \frac{4}{5}\right)^2 + 3.$$

$$\therefore \quad \text{when } \gamma = \frac{4}{5}, \quad f(\gamma) = 3$$

3

Concave Down, since $-5$

4/5

need to find roots to determine

when $f(\gamma) \geq 0$; this is when we should

use $\pi e$, o.w. when $f(\gamma) < 0$, use $\pi_{1a}$

$$\text{roots} = \frac{-b \pm \sqrt{b^2 - 4ae}}{2a} = \frac{-8 \pm \sqrt{64 - 4(-5)(-1)}}{-10}$$

$$= \frac{-8 \pm 2\sqrt{11}}{-10} = \frac{4 \pm \sqrt{11}}{5} \approx (.137, 1.46)$$

$$\text{Optimal Policy} = \begin{cases} (go, go) & \text{if } .137 \leq \gamma \leq 1.46 \\ (stay, stay, \cdots, stay) & o.w. \end{cases}$$

__|c|__ $V^0 = [0, 0]$

$V^1(s_1) = \max \left\{ r(s_1, stay) + V^0(s_1), \right.$

$$\left. r(s_1, go) + V^0(s_2) \right\}$$

$$= \max \left\{ -2 + 0, -3 + 0 \right\} = -2$$

$V^1(s_2) = \max \left\{ r(s_2, stay) + V^0(s_2), r(s_2, go) \right\}$

$$= \max \left\{ -2 + 0, 5 \right\} = 5$$

$\therefore V^1 = [-2, 5]$

$V^2(s_1) = \max \left\{ r(s_1, stay) + V^1(s_1), r(s_1, go) + V^1(s_2) \right\}$

$$= \max \left\{ -2 - 2, -3 + 5 \right\} = 2$$

$V^2(s_2) = \max \left\{ r(s_2, stay) + V^1(s_2), r(s_2, go) \right\}$

$$= \max \left\{ -2 + 5, 5 \right\} = 5$$

$\therefore V^2 = [2, 5]$

[c]

$$V^3(s_1) = \max \left\{ r(s_1, stay) + V^2(s_1), \ r(s_1, go) + V^2(s_2) \right\}$$

$$= \max \left\{ -2+2, \ -3+5 \right\} = 2$$

$$V^3(s_2) = \max \left\{ r(s_2, stay) + V^2(s_2), \ r(s_2, go) \right\} = 5$$

$$V^3 = [2, 5]$$

$$\therefore V^* = [2, 5]$$

## 2)

**a)** $\|v^i - v^*\|_\infty = \max_{s \in \mathcal{F}} |v^i(s) - v^*(s)|$

$i=1; \quad \max\left\{|v^1(s_1) - v^*(s_1)|, |v'(s_2) - v^*(s_2)|\right\} = 4$

$\qquad\qquad\qquad\quad 4 \qquad\qquad\qquad\quad 0$

$i=2; \quad \max\left\{|v^2(s_1) - v^*(s_1)|, |v^2(s_2) - v^*(2)|\right\} = 0$

$i=3; \quad \max\left\{|v^3(s_1) - v^*(a)|, |v^3(s_2) - v^*(s_2)|\right\} = 0$

**b)** $T(v) = \max_a \sum_{s'} p(s'|s,a)\left[r(s,a) + r v(s')\right]$

$\qquad T(v') = \max_a \sum_{s'} p(s'|s,a)\left[r(s,a) + r v'(s')\right]$

$\|Tv - Tv'\|_\infty$

$= \left\| \max_a \sum_{s'} p(s'|s,a)\left[r(s,a) + r v(s')\right] - \max_a \sum_{s'} p(s'|s,a) \right.$
$\left[ r(s,a) + r v'(s')\right] \Big\|_\infty$

$\leq \max_a \left\| \sum_{s'} p(s'|s,a)\left[r(s,a) + r v(s')\right] - \sum_{s'} p(s'|s,a) \right.$
$\left[ r(s,a) + r v'(s')\right]\Big\|_\infty \quad \|\max_a f(a) - \max_a (g(a))\|_\infty$

$\leq \max_a \| f(a) - g(a) \|_\infty$

$$= \max_a \left\| \sum_{s'} p(s'|sa) \left[ rv(s') - rv'(s') \right] \right\|_\infty$$

$$= \max_a \left\| \sum_{s'} p(s'|sa) \, r(v(s') - v'(s')) \right\|_\infty$$

$$= r \max_a \left\| \sum_{s'} p(s'|sa)(v(s') - v'(s')) \right\|_\infty$$

$$\leq \gamma \left\| v(s') - v'(s') \right\|_\infty$$

---

2c] $\left\| V^{n+1} - V^* \right\|_\infty = \left\| V^n - V^{n+1} \right\|_\infty$

$$= \left\| \sum_{t=1}^{\infty} V^{n+1+t} - V^{n+t} \right\|_\infty$$

$$\leq \sum_{t=1}^{\infty} \left\| V^{n+1+t} - V^{n+t} \right\|_\infty$$

$$\leq \sum_{t=1}^{\infty} \gamma^t \left\| V^{n+1} - V^n \right\|_\infty$$

$$= \left\| V^{n+1} - V^n \right\|_\infty \cdot \frac{r}{1-r}$$

$\therefore \left\| V^{n+1} - V^* \right\|_\infty \leq$

$\dfrac{\gamma}{1-\gamma} \left\| V^{n+1} - V^n \right\|_\infty$

if $\left\| V^{n+1} - V^n \right\|_\infty < \varepsilon$

then

$\left\| V^{n+1} - V^* \right\|_\infty \leq \dfrac{\gamma}{1-\gamma} \varepsilon$

**2d]** Suppose $V^* = T^{(\infty)}(0)$

$$T(V^*) = T(T^\infty(0)) = T^\infty(0) = V^*$$

$$\therefore \quad T(V^*) = V^*$$

---

Suppose $V^* = T^\infty(0)$

$$\| T^n(v) - T^{(n)}(v') \|_\infty \leq \gamma^n \| v - v' \|_\infty \quad \text{by 2b}$$

$$\| T^\infty(v) - T^\infty(0) \|_\infty \leq \gamma \| v - 0 \|_\infty = 0$$

$$\| T^\infty(v) - V^* \|_\infty = 0$$

$$T^\infty(v) = V^*$$

**4a]**

$$\nabla_\theta J(\theta) = \nabla_\theta E_{\tau \sim \pi_\theta} [R(\tau)] \quad - (A)$$

if $R(\tau)$ is changed to $R(\tau) - b$,

then $\nabla_\theta J(\theta) = \nabla_\theta E_{\tau \sim \pi_\theta} [R(\tau) - b]$

s.t. $b$ is not a function of $\theta$

$$= E_{\tau \sim \pi_\theta} [\nabla_\theta (R(\tau) - b)]$$

$$= E_{\tau \sim \pi_\theta} [\nabla_\theta R(\tau)]$$

$$= \nabla_\theta E_{\tau \sim \pi_\theta} [R(\tau)] \quad ----- (B)$$

$$\therefore \quad (A) = (B)$$