

se no.	s	s'	a	r	$P[s' s, a]$	$P[r s, a, s']$
1	H	H	SH	10	0.95	0.95
2	H	H	SH	0	0.95	0.05
3	H	̄H	SH	-10	0.05	0.90
4	H	̄H	SH	0	0.05	0.10
5	H	H	G	20	0.70	0.90
6	H	H	G	0	0.70	0.10
7	H	̄H	G	-10	0.30	0.90
8	H	̄H	G	0	0.30	0.10
9	̄H	̄H	M	-2	0.1	1
10	̄H	H	M	-1	0.9	1
11	̄H	̄H	̄M	-1	0.4	1
12	̄H	H	̄M	0	0.6	1

$$P[s'|s, a] = \sum_{r \in R(s)} P[r, s' | s, a] = \sum_{r \in R(s)} P[r | s, s', a] P[s' | s, a]$$

$$P[s', a | s, a] = P[a_1 | s', s, a] P[s' | s, a]$$

s	s'	a	a_1	$P[s', a_1 s, a]$
H	H	SH	10	0.9025
H	H	SH	0	0.0475
H	H	SH	-10	0.045
H	H	SH	0	0.005
<hr/>				
H	H	G	20	0.63
H	H	G	0	0.07
H	H	G	-10	0.27
H	H	G	0	0.03
<hr/>				
H	H	M	-2	0.1
H	H	M	-1	0.9
H	H	M	-1	0.4
H	H	M	0	0.6
<hr/>				

$$A = \{ S, H, M, \bar{M} \}$$

\downarrow \downarrow \downarrow \downarrow
 a_1 a_2 a_3 a_4

$$S = \{ H, \bar{H} \}$$

\downarrow \downarrow
 h \bar{h}

$$U_{\pi}(s) = \sum_{a} \pi(a|s) \sum_{s', r} P[s', r | s, a] (r + \gamma U_{\pi}(s'))$$

Given Policy: $\forall a \in A(s)$
 $\pi(a|s)$ is equal

Let's start w/t current state healthy.

$$U_{\pi}(h) = \sum_{a} \pi(a|h) \sum_{s'} P(s', r'|h, a) (r + \gamma U_{\pi}(s'))$$
$$A(h) = \{ \text{SM, G} \} = \{ a_1, a_2 \}$$

$$U_{\pi}(h) = \frac{1}{2} \left(0.9025 \left(10 + \frac{9}{10} U_{\pi}(h) \right) \right. \\ + 0.0475 \left(0 + \frac{9}{10} U_{\pi}(h) \right) \\ + 0.045 \left(-10 + \frac{9}{10} U_{\pi}(h) \right) \\ \left. + 0.005 \left(0 + \frac{9}{10} U_{\pi}(h) \right) \right)$$
$$+ \frac{1}{2} \left(0.63 \left(20 + \frac{9}{10} U_{\pi}(h) \right) \right. \\ + 0.07 \left(0 + \frac{9}{10} U_{\pi}(h) \right) \\ + 0.27 \left(-10 + \frac{9}{10} U_{\pi}(h) \right) \\ \left. + 0.03 \left(0 + \frac{9}{10} U_{\pi}(h) \right) \right)$$

$$0.2575 u_{\pi}(h) = 9.2375 + 0.1575 u_{\pi}(h)$$

① $u_{\pi}(h) = 35.874 + 0.611 u_{\pi}(\bar{h})$

For starting with ' \bar{h} '

$$u_{\pi}(\bar{h}) = \sum_{\alpha} \pi(\alpha | s) \left[p(s, s' | \bar{h}, \alpha) (s + \gamma u_{\pi}(s')) \right]$$

$$A(h) = \{S, H, M, \bar{M}\} = \{\alpha_1, \alpha_2\}$$

$$u_{\pi}(\bar{h}) = \frac{1}{2} \left(0.1(-2 + \gamma u_{\pi}(\bar{h})) + 0.9(-1 + \gamma u_{\pi}(h)) \right) \quad \begin{matrix} \\ \end{matrix} \rightarrow M$$

$$+ \frac{1}{2} \left(0.4(-1 + \gamma u_{\pi}(\bar{h})) + 0.6(0 + \gamma u_{\pi}(h)) \right) \quad \begin{matrix} \\ \end{matrix} \rightarrow \bar{M}$$

$$u_{\pi}(\bar{h}) = \frac{1}{2} \left[-0.2 + \frac{9}{100} u_{\pi}(\bar{h}) - 0.9 + \frac{81}{100} u_{\pi}(h) - 0.4 + \frac{36}{100} u_{\pi}(\bar{h}) + \frac{54}{100} u_{\pi}(h) \right]$$

$$U_{\pi}(\bar{h}) = \frac{1}{2} \left[-1.5 + \frac{45}{100} U_{\pi}(\bar{h}) + \frac{135}{100} U_{\sigma}(\bar{h}) \right]$$

$$U_{\sigma}(\bar{h}) = -0.75 + \frac{45}{200} U_{\pi}(\bar{h}) + \frac{135}{200} U_{\sigma}(\bar{h})$$

$$\frac{U_{\pi}(\bar{h})}{200}(135) = 0.75 + \frac{155}{200} U_{\sigma}(\bar{h})$$

$$U_{\sigma}(\bar{h}) = \frac{150}{135} + \frac{155}{135} U_{\sigma}(\bar{h})$$

$$U_{\pi}(\bar{h}) = 1.11 + 1.148 U_{\sigma}(\bar{h}) \quad \textcircled{2}$$

From ① + ②

$$\overline{U_{\pi}(\bar{h})} = 75.42$$

$$U_{\sigma}(\bar{h}) = 64.737$$

$$21) \quad t \rightarrow H \quad f = E(G_t | S_t = h, S_{t+1} = \bar{h})$$

$$t+1 \longrightarrow \bar{H}$$

$$f = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = h, S_{t+1} = \bar{h}]$$

$$f = E[R_{t+1} | S_t = h, S_{t+1} = \bar{h}] + \gamma E[G_{t+1} | S_t = h, S_{t+1} = \bar{h}]$$

$$f = \sum_a \pi_a P[a | S_t = h, S_{t+1} = \bar{h}] + \gamma E[G_{t+1} | S_t = h, S_{t+1} = \bar{h}]$$

$$f = (-10(0.9)(\frac{1}{2}) + 0(0.1)(\frac{1}{2}) + \frac{1}{2}(0.9)(-10))$$

↓ ↓ ↓ ↓
 get a reward probability of getting -10 probability of taking action SH. + $\frac{1}{2}(0.1)(0)$
 -10, when $s = h$ & $s' = \bar{h}$ reward when $s = h, s' = \bar{h}$
 $s = h$ & $a = SH$ $s = h, s' = \bar{h}$

$$\gamma E[G_{t+1} | S_t = h, S_{t+1} = \bar{h}]$$

$$f = -4 \cdot 5 - 4 \cdot 5 + \gamma E[G_{t+1} | S_t = h, S_{t+1} = \bar{h}]$$

$$f = -9 + \gamma E[G_{t+1} | S_t = h, S_{t+1} = \bar{h}]$$

$$\hookrightarrow \gamma u_\pi(\bar{h})$$

$$f = -9 + \frac{9}{10}(75 \cdot 42) \Rightarrow f = 58.878$$

$$P[s', a' | s, a] = P[a_1 | s', s, a] P[s' | s, a]$$

s	s'	a	a_1	$P[s', a_1 s, a]$
H	H	SH	10	0.9025
H	H	SH	0	0.0475
H	H	SH	-10	0.045
H	H	SH	0	0.005
<hr/>				
H	H	G	20	0.63
H	H	G	0	0.07
H	H	G	-10	0.27
H	H	G	0	0.03
<hr/>				
H	H	M	-2	0.1
H	H	M	-1	0.9
H	H	M	-1	0.4
H	H	M	0	0.6
<hr/>				

$$\gamma = 0.9$$

22) Let's start with π_0 , where each action is equally probable.

V_{π_0} FOR HEALTHY STATE

$$U_{\pi_0}(H) = \sum_{a \in A(H)} \pi_0(a|H) \sum_{s, s'} (r_a + \gamma V_{\pi_0}(s')) P[s, s' | H, a]$$

$$V_{\pi_0}(H) = \frac{1}{2} \left((10 + 0.9 V_{\pi_0}(H)) 0.9025 \right.$$

\nearrow

$$\qquad \qquad \qquad + (0 + 0.9 V_{\pi_0}(H)) 0.0475$$

$$\qquad \qquad \qquad + (-10 + 0.9 V_{\pi_0}(\bar{H})) 0.045$$

$$\qquad \qquad \qquad \left. + (0 + 0.9 V_{\pi_0}(\bar{H})) 0.005 \right)$$

\nearrow

$$+ \frac{1}{2} \left((20 + 0.9 V_{\pi_0}(H)) 0.63 \right.$$

$$\qquad \qquad \qquad + (0 + 0.9 V_{\pi_0}(H)) 0.07$$

$$\qquad \qquad \qquad + (-10 + 0.9 V_{\pi_0}(\bar{H})) 0.27$$

$$\qquad \qquad \qquad \left. + (0 + 0.9 V_{\pi_0}(\bar{H})) 0.03 \right)$$

$$U_{\pi_0}(H) = \frac{1}{2} \left(9.025 - 0.45 + (U_{\pi_0}(H)) \left(\frac{9}{10} \right) (0.95) + (U_{\pi_0}(\bar{H})) \left(\frac{9}{10} \right) (0.05) \right)$$

$$+ \frac{1}{2} \left(12.6 - 2.7 + (U_{\pi_0}(H)) \left(\frac{9}{10} \right) (0.7) + (U_{\pi_0}(\bar{H})) \left(\frac{9}{10} \right) (0.3) \right)$$

$$U_{\pi_0}(H) = \frac{1}{2} \left(18.475 + (U_{\pi_0}(H)) (1.485) + (U_{\pi_0}(\bar{H})) (0.315) \right)$$

$$U_{\pi_0}(H) = 9.2375 + 0.7425 U_{\pi_0}(H) + 0.1575 U_{\pi_0}(\bar{H})$$

$$0.2575 U_{\pi_0}(H) = 9.2375 + 0.1575 U_{\pi_0}(\bar{H})$$

$$U_{\pi_0}(H) = 35.8738 + 0.6117 U_{\pi_0}(\bar{H})$$

①

$$U_{\pi_0}(H) = 0.6117 U_{\pi_0}(\bar{H}) + 35.8738$$

\mathcal{U}_{π_0} FOR NOT HEALTHY STATE

$$\mathcal{U}_{\pi_0}(\bar{H}) = \frac{1}{2} \left[(-2 + 0.9 \mathcal{U}_{\pi_0}(\bar{H})) 0.1 \right. \\ \left. + (-1 + 0.9 \mathcal{U}_{\pi_0}(H)) 0.9 \right] \\ + \frac{1}{2} \left[(-1 + 0.9 \mathcal{U}_{\pi_0}(\bar{H})) 0.4 \right. \\ \left. + (0 + 0.9 \mathcal{U}_{\pi_0}(H)) 0.6 \right]$$

$$\mathcal{U}_{\pi_0}(\bar{H}) = \frac{1}{2} \left(-1.5 + \mathcal{U}_{\pi_0}(H) (1.35) \right. \\ \left. + \mathcal{U}_{\pi_0}(\bar{H}) (0.45) \right)$$

$$\mathcal{U}_{\pi_0}(\bar{H}) = -0.75 + 0.675 \mathcal{U}_{\pi_0}(H) \\ + 0.225 \mathcal{U}_{\pi_0}(\bar{H})$$

$$\boxed{\mathcal{U}_{\pi_0}(\bar{H}) = -0.9677 + 0.871 \mathcal{U}_{\pi_0}(H)}$$

$$U_{\pi_0}(H) = 0.6117 U_{\pi_0}(\bar{H}) + 35.8738$$

$$U_{\pi_0}(\bar{H}) = -0.9677 + 0.871 U_{\pi_0}(H)$$

$$U_{\pi_0}(H) = 75.516$$

$$U_{\pi_0}(\bar{H}) = 64.807$$

Now, we improve the policy

$$\pi_t(s) = \underset{a \in A(s)}{\operatorname{argmax}} \sum_{s', r} P[s', a | s, a] (r + \gamma U_{\pi_0}(s'))$$

For $s = H$

$$\begin{aligned} \pi_t(H) = \underset{a \in A(s)}{\operatorname{argmax}} & \left\{ \begin{aligned} & (10 + 0.9 U_{\pi_0}(H)) 0.9025 \\ & + (0 + 0.9 U_{\pi_0}(H)) 0.0475 \\ & + (-10 + 0.9 U_{\pi_0}(\bar{H})) 0.045 \\ & + (0 + 0.9 U_{\pi_0}(\bar{H})) 0.005 \end{aligned} \right. \end{aligned}$$

$$a = SH$$

$$\begin{aligned} \pi_t(H) \rightsquigarrow & (20 + 0.9 U_{\pi_0}(H)) 0.63 \\ & + (0 + 0.9 U_{\pi_0}(H)) 0.07 \\ & + (-10 + 0.9 U_{\pi_0}(\bar{H})) 0.27 \\ & + (0 + 0.9 U_{\pi_0}(\bar{H})) 0.03 \end{aligned}$$

$$\pi_1(h) = \underset{a \in A(h)}{\operatorname{argmax}} \left\{ \begin{array}{l} 76.0575, 74.973 \\ a = S \\ a = G \end{array} \right\}$$

$$\therefore \boxed{\pi_1(h) = S}$$

$$\text{For } s = \bar{H} \quad a = M \rightarrow$$

$$\pi_1(\bar{H}) = \underset{a \in A(\bar{H})}{\operatorname{argmax}} \left\{ \begin{array}{l} (-2 + 0.9 u_{\pi_0}(\bar{H})) 0.1 \\ + (-1 + 0.9 u_{\pi_0}(H)) 0.9, \\ (-1 + 0.9 u_{\pi_0}(\bar{H})) 0.4 \\ + (0 + 0.9 u_{\pi_0}(H)) 0.6 \end{array} \right\}$$

$$a = \bar{M} \rightarrow$$

$$\pi_1(\bar{H}) = \underset{a \in A(\bar{H})}{\operatorname{argmax}} \left\{ 65.9, 63.709 \right\}$$

$$\therefore \boxed{\pi_1(\bar{H}) = M}$$

Now we find $\mathbb{U}_{\pi_1}(H)$ & $\mathbb{U}_{\pi_1}(\bar{H})$

$$\pi(SH|H) = 1, \text{ but is zero}$$

$$\mathbb{U}_{\pi_1}(H) = \sum_{s', r} p[r, s' | H, a] \binom{r+s'}{\mathbb{U}_{\pi_1}(s')}$$

$$\begin{aligned}\mathbb{U}_{\pi_1}(H) &= (10 + 0.9(\mathbb{U}_{\pi_1}(H))) 0.9025 \\ &\quad + (0 + 0.9(\mathbb{U}_{\pi_1}(H))) 0.0475 \\ &\quad + (-10 + 0.9(\mathbb{U}_{\pi_1}(\bar{H}))) 0.045 \\ &\quad + (0 + 0.9(\mathbb{U}_{\pi_1}(\bar{H}))) 0.005\end{aligned}$$

$$\begin{aligned}\mathbb{U}_{\pi_1}(H) &= (9.025 - 0.45) + (\mathbb{U}_{\pi_1}(H))(0.855) \\ &\quad + (\mathbb{U}_{\pi_1}(\bar{H}))(0.045)\end{aligned}$$

$$\begin{aligned}\mathbb{U}_{\pi_1}(H) &= 8.575 + 0.855(\mathbb{U}_{\pi_1}(H)) \\ &\quad + 0.045(\mathbb{U}_{\pi_1}(\bar{H}))\end{aligned}$$

$$\boxed{\mathbb{U}_{\pi_1}(H) = 59.138 + 0.31 \mathbb{U}_{\pi_1}(\bar{H})}$$

①

$$V_{\pi_1}(\bar{H}) = -1.22 + 0.89 V_{\pi_1}(H)$$

$$V_{\pi_1}(H) = 59.138 + 0.31 V_{\pi_1}(\bar{H})$$

$$V_{\pi_1}(H) = 84.1487$$

$$V_{\pi_1}(\bar{H}) = 71.00$$

Our State value increased
after Policy improvement.

$$P[s', a' | s, a] = P[a_1 | s', s, a] P[s' | s, a]$$

s	s'	a	a_1	$P[s', a_1 s, a]$
H	H	SH	10	0.9025
H	H	SH	0	0.0475
H	H	SH	-10	0.045
H	H	SH	0	0.005
<hr/>				
H	H	G	20	0.63
H	H	G	0	0.07
H	H	G	-10	0.27
H	H	G	0	0.03
<hr/>				
H	H	M	-2	0.1
H	H	M	-1	0.9
H	H	M	-1	0.4
H	H	M	0	0.6
<hr/>				

$$\gamma = 0.9$$