

16)

$$S = \{ \text{Fresh}, \text{Stale} \}$$

$$A = \{ \text{Query}, \text{Silent} \}$$

s	s'	a	$P[s' s, a]$	r
Stale	Fresh	Query	0.8	-8
Stale	Stale	Query	0.2	-8
Fresh	Stale	Query	0.1	-4
Fresh	Fresh	Query	0.9	-4
Stale	Fresh	Silent	0	-
Stale	Stale	Silent	1	-4
Fresh	Stale	Silent	0.5	4
Fresh	Fresh	Silent	0.5	4

(b)) $\text{General } V_{\pi}(n) = \sum_a \pi(a|n) \sum_{s, s'} P[n', s | n, a] (r + \gamma V_{\pi}(n'))$
 say knowledge practice

State $\rightarrow s$	Delivery $\rightarrow a_1$	Notation
Fresh - f	Silent $\rightarrow a_2$	

Let's say we start in State state.

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

n' can be 'f' or 's' depending on our action.

Starting In State

$$\pi(a|s) = \frac{1}{2} \text{ always}$$

Assuming this policy

$$V_{\pi}(s) = \left(P[a_1 | s] \sum_{s', n'} P[n', s | s, a_1] (r + \gamma V_{\pi}(n')) \right)$$

$$+ \left(P[a_2 | s] \sum_{s', n'} P[n', s | s, a_2] (r + \gamma V_{\pi}(n')) \right)$$

if I choose delivery

$$V_{\pi}(s) = \frac{1}{2} \left(P[s, -8 | s, a_1] (-8 + \gamma V_{\pi}(s)) \right)$$

if I choose to stay silent

$$+ P[f, -8 | s, a_1] (-8 + \gamma V_{\pi}(f))$$

stay silent

$$+ \frac{1}{2} \left(P[s, 4 | s, a_2] (4 + \gamma V_{\pi}(s)) \right)$$

$$+ P[f, -1 | s, a_2] (-1 + \gamma V_{\pi}(f))$$

$$U_{\pi}(s) = \frac{1}{2} \left[\left(0 \cdot 2 (-8 + \gamma U_{\pi}(s)) \right. \right. \\ \left. \left. + (0 \cdot 8 (-8 + \gamma U_{\pi}(d))) \right) \right]$$

$$+ \frac{1}{2} \left(1 (4 + \gamma U_{\pi}(s)) \right)$$

$$U_{\pi}(s) = \frac{1}{2} \left(-1 \cdot 6 + \frac{\gamma}{5} U_{\pi}(s) - 6 \cdot 4 + \frac{4}{5} \gamma U_{\pi}(d) \right. \\ \left. + \frac{1}{2} (4 + \gamma U_{\pi}(s)) \right)$$

$$U_{\pi}(s) = \frac{1}{2} \left(-8 + \frac{\gamma}{5} U_{\pi}(s) + \frac{4}{5} \gamma U_{\pi}(d) \right. \\ \left. + 4 + \gamma U_{\pi}(s) \right)$$

$$U_{\pi}(s) = \frac{1}{2} \left(-4 + \frac{6\gamma}{5} U_{\pi}(s) + \frac{4}{5} \gamma U_{\pi}(d) \right)$$

$$U_{\pi}(s) = -2 + \frac{2}{5} \gamma U_{\pi}(d) + \frac{3\gamma}{5} U_{\pi}(s)$$

$$U_{\pi}(s) - \frac{3}{5} \gamma U_{\pi}(s) = \frac{2}{5} \gamma U_{\pi}(d) - 2$$

Now we start w/t $n = f$ &

find $u_{\pi}(f)$

if I choose
any

$$u_{\pi}(f) = \frac{1}{2} \left[P[f, -4 | f, a_1] (-4 + r u_{\pi}(f)) \right]$$

$$+ P[s, -4 | f, a_1] (-4 + r u_{\pi}(s))$$

$$+ \frac{1}{2} \left[P[f, 4 | f, a_2] (4 + r u_{\pi}(f)) \right.$$

$$\left. + P[s, 4 | f, a_2] (4 + r u_{\pi}(s)) \right]$$

$$u_{\pi}(f) = \frac{1}{2} \left(0.9 (-4 + r u_{\pi}(f)) + 0.1 (-4 + r u_{\pi}(s)) \right)$$

$$+ \frac{1}{2} \left(0.5 (4 + r u_{\pi}(f)) + 0.5 (4 + r u_{\pi}(s)) \right)$$

$$u_{\pi}(f) = \frac{1}{2} \left(-3.6 + \frac{9r}{10} u_{\pi}(f) - 0.4 + \frac{r u_{\pi}(s)}{10} \right)$$

$$+ \frac{1}{2} (4 + \frac{r}{2} u_{\pi}(f) + \frac{r}{2} u_{\pi}(s))$$

$$U_{\bar{J}}(f) = \frac{1}{2} \left(-4 + \frac{9r}{10} U_{\bar{J}}(f) + \frac{r}{10} U_{\bar{J}}(s) \right)$$

$$+ \frac{1}{2} \left(4 + \frac{r}{2} U_{\bar{J}}(f) + \frac{r}{2} U_{\bar{J}}(s) \right)$$

$$U_{\bar{J}}(f) = \cancel{-2} + \frac{9r}{20} U_{\bar{J}}(f) + \frac{r}{20} U_{\bar{J}}(s)$$

$$+ \frac{r}{4} U_{\bar{J}}(f) + \frac{r}{4} U_{\bar{J}}(s)$$

$$U_{\bar{J}}(f) = \frac{14r}{20} U_{\bar{J}}(f) + \frac{6r}{20} U_{\bar{J}}(s)$$

$$U_{\bar{J}}(f) - \frac{14r}{20} U_{\bar{J}}(f) = \frac{6r}{20} U_{\bar{J}}(s)$$

$$U_{\pi}(s) - \frac{3}{5} r U_{\pi}(s) = \frac{2}{5} r U_{\pi}(f) - 2$$

$$U_{\pi}(s) \rightarrow i$$

$$U_{\pi}(f) \rightarrow j$$

$$r = 1/2$$

$$U_{\pi}(f) - \frac{14r}{20} U_{\pi}(f) = \frac{6r}{20} U_{\pi}(s)$$

$$i - \frac{3}{10} i = \frac{2}{10} f - 2$$

$$\frac{7}{10} i = \frac{2}{10} f - 2$$

$$i = \frac{2}{7} f - \frac{20}{7}$$

$$j - \frac{14}{40} j = \frac{6}{40} i$$

$$\frac{6i}{40} = \frac{26}{40} j$$

$$i = \frac{13}{3} j$$

$$\frac{13}{3} j = \frac{2}{7} f - \frac{20}{7}$$

$$\frac{91}{21} j - 6j = -\frac{20}{7} \Rightarrow 8j = -60$$

$$j = -\frac{60}{85} = -\frac{12}{17}$$

$$i = \frac{13}{3} \left(-\frac{12}{17} \right) = -\frac{52}{17}$$

$$U_{\pi}(s) = -\frac{52}{17} \quad \text{An}$$

$$U_{\pi}(f) = -12 / 17$$

16)

$$S = \{ \text{Fresh, Stale} \}$$

Point 3)

$$A = \{ \text{Query, Silent} \}$$

s	s'	a	$P[s' s, a]$	r
Stale	Fresh	Query	0.8	-8
Stale	Stale	Query	0.2	-8
Fresh	Stale	Query	0.1	-4
Fresh	Fresh	Query	0.9	-4
Stale	Fresh	Silent	0	-
Stale	Stale	Silent	1	-4
Fresh	Stale	Silent	0.5	4
Fresh	Fresh	Silent	0.5	4

We found state values for policy π_0 .

$$\pi_0(a|x) = \frac{1}{2}$$

Stale $\rightarrow s$	query $\rightarrow a_1$
Fresh $\rightarrow f$	silent $\rightarrow a_2$

$$U_{\pi_0}(s) = -\frac{52}{17}$$

$$U_{\pi_0}(f) = -\frac{12}{17}$$

Now, we improve the policy :-

FOR STALE START

$$\pi_1(s) = \underset{a \in A(s)}{\operatorname{argmax}} \sum_{n, n'} P[s, n' | s, a] (\pi + \gamma U_{\pi_0}(n'))$$

for a_1 ,

$$\sum_{n, n'} P[s, n' | s, a_1] (\pi + \gamma U_{\pi_0}(n'))$$

$$= (0 \cdot 2) \left(-8 + \frac{1}{2} \left(-\frac{52}{17} \right) \right) + (6 \cdot 8) \left(-8 + \frac{1}{2} \left(-\frac{12}{17} \right) \right)$$

$$= -1 \cdot 9 - 6 \cdot 6 \cdot 8$$

$$= -8 \cdot 585$$

for a_2

$$\sum_{n, n'} P[s, n' | s, a_2] (\pi + \gamma U_{\pi_0}(n'))$$

$$= 1 \left(4 + \frac{1}{2} \left(-\frac{52}{17} \right) \right)$$

$$= 2 \cdot 47$$

For Stale Start

∴ Choosing silent is the greedy policy

FOR FRESH START

$$\sum_{n,n'} P[n, n' \mid f, \alpha_1] (n + r u_{\text{ss}}(n'))$$

for α_1

$$\sum_{n,n'} P[n, n' \mid f, \alpha_1] (n + r u_{\text{ss}}(n'))$$

$$= (0.1) \left(-4 + \frac{1}{2} \left(-\frac{52}{17} \right) \right) + (0.9) \left(-4 + \left(-\frac{12}{34} \right) \right)$$

$$= -0.553 - 3.91$$

$$= -4.47$$

for α_2

$$\sum_{n,n'} P[n, n' \mid f, \alpha_2] (n + r u_{\text{ss}}(n'))$$

$$= \frac{1}{2} \left(4 + \frac{1}{2} \left(-\frac{12}{17} \right) \right) + \frac{1}{2} \left(4 + \frac{1}{2} \left(-\frac{52}{17} \right) \right)$$

$$= 1.8235 + 1.2352$$

$$= \underline{\underline{3.0587}}$$

Pick α_2 (Silent)

∴ our policy is unchanged

$$\pi_2(s) = a_2$$

$$\pi_2(f) = a_2$$

$$U_{\pi_2}(s) = 1(4 + \frac{1}{2} U_{\pi_2}(s))$$

$$U_{\pi_2}(s) = 4 + \frac{U_{\pi_2}(s)}{2}$$

$$U_{\pi_2}(s) = 8$$

$$U_{\pi_2}(f) = \frac{1}{2} \left(4 + \frac{1}{2} U_{\pi_2}(s) \right) + \frac{1}{2} \left(4 + \frac{1}{2} U_2(f) \right)$$

$$U_{\pi_2}(f) = \frac{1}{2} (4+4) + 2 + \frac{U_{\pi_2}(f)}{4}$$

$$\frac{3}{4} U_{\pi_2}(f) = 6$$

$$U_{\pi_2}(f) = 8$$

Always choose
Silence!
Ans

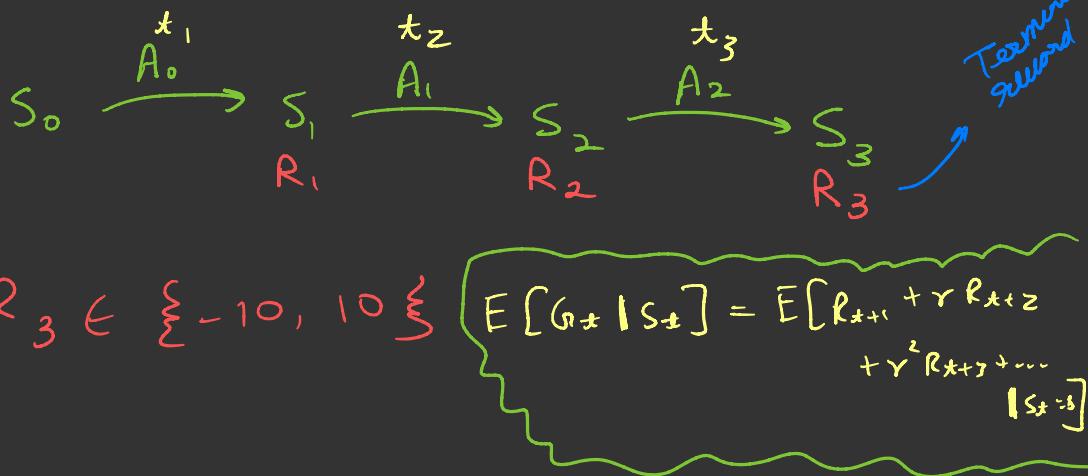
Since $U_{\pi_1}(s) = U_{\pi_2}(s)$ & $U_{\pi_1}(f) = U_{\pi_2}(f)$

∴ We have (concluded) : $\begin{cases} \pi(s) = a_2 \\ \pi(f) = a_2 \end{cases}$

Post 2/1

Question 16. A server requires information from a sensor. The server would like the information to be fresh. However, there is a cost to querying information from the sensor. Specifically, the state at the server can be either fresh or stale. The former indicates that the information at the server about the sensor is fresh and the latter indicates that it is stale. At any time, the server may choose to query or remain silent. A query makes stale information fresh with probability 0.8 and fresh information stale with probability 0.1. Staying silent keeps stale information stale with probability 1 and makes fresh information stale with probability 0.5. Also, a query when current information at the server is stale costs 8 dollars. Otherwise, it costs 4 dollars. Staying quiet has a reward of 4 dollars.

1. Draw and/or tabulate the MDP.
2. Consider a finite horizon problem in which the server gets to optimize costs over three time steps. Assume that if at the end of the three time steps the server has stale information it pays a cost of 10, else it receives a reward of 10. Also, assume that future costs are discounted with a factor $1/2$. Calculate the optimal costs for every starting state. Calculate the optimal policy. Show all your steps. Guesses and intuition will bring you unbounded costs.
3. Suppose the server is to optimize for ever and is looking for a stationary policy. Write down the first four iterations of value iteration that will help the server get closer to a stationary policy. Do the same for policy iteration (where one iteration includes evaluation and improvement). Clearly identify each iteration.



Bellman optimality status if we can find

π_2^* then for π_1^* & π_0^* , we can use state values given by π_2^* as it will be the optimum.



If route 1 + route 4 is the optimum from place 1 to 3. Then route 4 is opt for 2-3.

At the 3rd time step, we take an action & we either get reward 10 or -10.

$$\pi_2^* = \underset{a \in A(n)}{\operatorname{argmax}} \mathbb{E}[R_3 \mid x_2 = n, A_2 = a]$$

For Fresh ($x_2 = f$)

$$\pi_2^*(f) = \operatorname{argmax} \left\{ \begin{array}{l} 0.9(10) + 0.1(-10), \\ 0.5(10) + 0.5(-10) \end{array} \right\} \xrightarrow{\text{anxious}}$$

$$\boxed{\pi_2^*(f) = \text{anxious}} \quad \xrightarrow{\text{silent}}$$

For State ($x_2 = s$)

$$\pi_2^*(s) = \operatorname{argmax} \left\{ \begin{array}{l} 0.8(10) + 0.2(-10), \\ 0.1(10) + 0.9(-10) \end{array} \right\} \xrightarrow{\text{anxious}}$$

$$\therefore \boxed{\pi_2^*(s) = \text{anxious}} \quad \xrightarrow{\text{silent}}$$

Now, we have the optimum policy at third time step. we use this to find expectations at $t_2 \& t_1$.

at 2nd time step.

$$\gamma = \frac{1}{2}$$

$$E[R_2 + \gamma R_3 | X_1 = x]$$

we need to find the optimum policy.

$$\pi_1^*(x) = \underset{a \in A}{\operatorname{argmax}} \left(E[R_2 + \gamma R_3 | X_1 = x, A_1 = a] \right)$$

$$= \underset{a \in A}{\operatorname{argmax}} \left(E[R_2 | X_1 = x, A_1 = a] + \gamma E[R_3 | X_1 = x, A_1 = a] \right)$$

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

17)

$$U_{\pi_0}(s) = \sum_a \pi_0(a|s) \sum_{s', r} (r + \gamma U_{\pi_0}(s')) P[s', s' | s, a]$$

$$\pi_1(s) = \underset{a \in A(s)}{\operatorname{argmax}} \left[\sum_{s', r} (r + \gamma U_{\pi_0}(s')) P[s', s' | s, a] \right]$$

Since, this policy will get the action, which minimizes $U_{\pi_0}(s)$, This means that either $U_{\pi_1}(s)$ will be greater than $U_{\pi_0}(s)$ or $U_{\pi_0}(s)$ won't minimized w.r.t π_0 or they'll be equal if in fact $\pi_0(s)$ picked the minimizing action.

(8)

$G, 1, 2, 3, \dots, n, F$

(τ)

(T)

$i \in [1, n]$

$j \in [1, n)$

Actions: $\{D, U\}$

s	s'	a	r	$P[r, s' s, a]$
i	$i-1$	D	2	0.5
i	$i-1$	D	0	0.5
1	G	D	1	1
j	$j+1$	U	-1	1
n	F	U	$2n$	1

$$n=2$$

$$\pi(a|s) = 0.5 \quad \forall$$

G, I, 2, F

$$a \in A(s)$$

s	s'	a	r _{s,s'}	P[s', s' s, a]
1	2	U	-1	1
2	F	U	4	1
2	1	D	2	0.5
2	1	D	0	0.5
1	G	D	1	1

We have 4 states out of which two are terminal

Let's define the initial state value for each state as zero.

Since G, F are terminal states, their value will remain zero.

$$V_{\pi}(G) = V_{\pi}(F) = 0$$

$$\gamma = 1$$

first we find the state value for the given policy. $\gamma = 1$, given

$$V_{\pi}(s) = E_{\pi} \left[R_{t+1} + \gamma V_{\pi}(s') \mid s_t = s, A_t = a \right]$$

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

$$V_{\pi}(1) = \pi(u|1) \left[(\gamma V_{\pi}(2) + (-1)) 1 \right]$$

$$+ \pi(d|1) \left[(\gamma V_{\pi}(6) + 1) 1 \right]$$

$$V_{\pi}(1) = \frac{1}{2} (V_{\pi}(2) - 1) + \frac{1}{2} (0 + 1)$$

$$\boxed{V_{\pi}(1) = V_{\pi}(2)/2} \quad ①$$

If we start in 2

$$V_{\pi}(2) = \frac{1}{2} (V_{\pi}(F) + 4) + \frac{1}{2} \left((V_{\pi}(1) + 2) \frac{1}{2} + (V_{\pi}(1) + 0) \frac{1}{2} \right)$$

$$U_{\pi}(2) = 2 + \frac{1}{2} (U_{\pi}(1) + 1)$$

$$U_{\pi}(2) = 2 \cdot 5 + \frac{U_{\pi}(1)}{2}$$

(2)

from ① & ②

$$U_{\pi_0}(1) = \frac{5}{3}$$

$$U_{\pi_0}(2) = \frac{10}{3}$$

$$2U_{\pi}(1) = 2 \cdot 5 + \frac{U_{\pi}(1)}{2}$$

$$4U_{\pi}(1) = 5 + U_{\pi}(1)$$

$$U_{\pi}(1) = \frac{5}{3}$$

$$U_{\pi}(2) = \frac{10}{3}$$

Next, we update our policy, for the first time!

$$\pi_1(s) = \underset{a \in A(s)}{\operatorname{argmax}} U_{\pi_0}(s)$$

$$\pi_1(s) = \underset{a \in A(s)}{\operatorname{argmax}} \sum_{s', r} P[r, s' | s, a] (r + \gamma U_{\pi_0}(s'))$$

$$\pi_1(1) = \underset{a \in A(1)}{\operatorname{argmax}} \left[(-1 + U_{\pi_0}(2)) \frac{1}{2}, (1 + 0) \frac{1}{2} \right]$$

$$\pi_1(1) = \cup p$$

$$\overline{\pi}_1(2) = \underset{a \in A(2)}{\operatorname{argmax}} \left[(\underline{\pi} + u_{\pi_0}(1)) \frac{1}{2} + (\underline{\alpha} + u_{\pi_0}(0)) \frac{1}{2} \right]$$

$$\overline{\pi}_1(2) = \underset{a \in A(2)}{\operatorname{argmax}} \left[\frac{D}{3}, \frac{U}{4} \right]$$

$$\overline{\pi}_1(2) = U_p$$

$$\overline{\pi}_1(1) = U_p$$

$$\overline{\pi}_1(2) = U_p$$

$$u_{\pi_1}(1) = (-1 + u_{\pi_1}(2)) 1$$

$$u_{\pi_1}(1) = u_{\pi_1}(2) - 1$$

$$u_{\pi_1}(2) = (4+0) / 4 = 1$$

$$u_{\pi_1}(2) = 1$$

$$u_{\pi_1}(1) = 3$$

$$\pi_2(s) = \operatorname{argmax}_{a \in A(s)} \left[\sum_{s', r} (r + u_{\pi}(s')) P[s', r | s, a] \right]$$

$$\pi_2(1) = \operatorname{argmax}_{a \in A(1)} \left[(-1 + u_{\pi}(2)) \text{U}, (1 + u_{\pi}(0)) \text{D} \right]$$

$$\pi_2(1) = \operatorname{argmax}_{a \in A(1)} \left[\text{U} \overline{3}, \text{D} \overline{1} \right]$$

$$\boxed{\pi_2(1) = \text{Up}}$$

$$\pi_2(2) = \operatorname{argmax}_{a \in A(2)} \left[(4 + u_{\pi}(F)) \text{U}, \frac{(2 + u_{\pi}(1))}{2} + \frac{(0 + u_{\pi}(1))}{2} \right]$$

$$\pi_2(2) = \operatorname{argmax}_{a \in A(2)} \left[\text{U} \overline{4}, \text{D} \overline{1+3} \right]$$

$$\boxed{\pi_2(2) = \text{Up or Down}}$$

$$\begin{aligned} \pi(\text{U}|2) &= 1/2 \\ \pi(\text{D}|2) &= 1/2 \end{aligned}$$

$$U_{\pi_2}(1) = \left(-1 + U_{\pi_2}(2)\right)1$$

$$U_{\pi_2}(1) = U_{\pi_2}(2) - 1$$

$$U_{\pi_2}(2) = \frac{1}{2}(4+0) + \frac{1}{2} \left[(2 + U_{\pi_2}(1)) \frac{1}{2} + (0 + U_{\pi_2}(1)) \frac{1}{2} \right]$$

$$U_{\pi_2}(2) = 2 + \frac{1}{2} \left[1 + U_{\pi_2}(1) \right]$$

$$U_{\pi_2}(2) = \frac{5}{2} + \frac{U_{\pi_2}(1)}{2}$$

$$\therefore U_{\pi_2}(2) = \frac{5}{2} + \frac{U_{\pi_2}(2)}{2} - \frac{1}{2}$$

$$U_{\pi_2}(2) = 4$$

$$U_{\pi_2}(1) = 3$$

$$\rightarrow V^*(s)$$

$$U_{\pi_{k+1}}(s) = U_{\pi_k}(s) + s \quad \therefore \text{ we converge}$$

$$\pi^*(1) = \text{Up} \quad \pi^*(2) = \{\text{Up}, \text{Down}\}$$

If I have more states, I will have to iterate more, as there will be more paths to terminal states.

$$n=2$$

$$\pi(a|s) = 0.5 \quad \forall$$

G, I, 2, F

$$a \in A(s)$$

s	s'	a	r _{s,s'}	P[s', s' s, a]
1	2	U	-1	1
2	F	U	4	1
2	1	D	2	0.5
2	1	D	0	0.5
1	G	D	1	1

We have 4 states out of which two are terminal

Let's define the initial state value for each state as zero.

Since G, F are terminal states, their value will remain zero.

$$V_{\pi}(G) = V_{\pi}(F) = 0$$

$$\gamma = 1$$

