mm 는게산으로 가시반복 한고기금은 성생활수 있다.

역신 t=0 이 734 모든 상태는 0으로 크기와 해결다.

ø	0	0	0
Ø	0	0	0
0	0	0	0
0	0	0	0
+=0			

S F F F F H F F F G

그러는 는 이 대 유민 흑필지점에 인접한 15번 칸에서

(각 해공은 앞는 각 표저 L,D,R,U 으로 표현된다, stochastic P=1)

 $\mathcal{V}^{(1)}(|F\rangle = \max \left\{ P(|I|, o||F, U) \left(o + V^{(o)}(|I|) \right), P(|B|, ||F, R) \left(|+V^{(o)}(|B|) \right), P(|F, o||F, D) \left(o + V^{(o)}(|F|) \right) \right\}$

= max { 0,1,0,0} =[

: A(15) = { Right }

H에서는 이름이 인지보는 없으므로 제의하는 각 칸에서 만간 가능을 개름해보면 다음과 같다.

$$\mathcal{V}^{(1)}_{(1)=0}, \mathcal{V}^{(1)}(2)=0$$
 , $\mathcal{V}^{(1)}(3)=0$, $\mathcal{V}^{(1)}(4)=0$, $\mathcal{V}^{(1)}(5)=0$, $\mathcal{V}^{(1)}(9)=0$

$$\mathcal{V}^{(1)}(q) = 0$$
, $\mathcal{V}^{(1)}(l^0) = 0$, $\mathcal{V}^{(1)}(l^1) = 0$, $\mathcal{V}^{(1)}(l^2) = 0$, $\mathcal{V}^{(1)}(l^2) = 0$

S. 0 0 0 0 H. 0 H. 0 0 0 H. 0 0 1 G.

til

Ho

50

0

Ho

Ho

이어서 고반과에 반복을 수해한다. 이번 전라에 이어서 때부분의 간의 전라는 이에고 15간라 근처찮은만 다시구해보다.

$$V^{(2)}(11) = \max_{x \in \mathbb{R}} \{ |x| (0+V^{(15)}), 0, 0, 0 \} = | , A(11) = \{ Down \}$$

 $V^{(2)}(14) = \max_{x \in \mathbb{R}} \{ |x| (0+V^{(15)}), 0, 0, 0 \} = | , A(14) = \{ Right \}$

3 비교N 사보에에도 각 관에서의 가시하는 구5시는다.

$$V^{(3)}(0) = 1$$

$$A(0) = \frac{1}{2} Down, Right$$

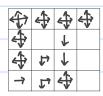
$$V^{(3)}(13) = 1$$
 $A(13) = 2 Right3$

$$\mathcal{V}^{(3)}(\eta) = \{ A(\eta) = \{ Down \} \}$$

$$V^{(3)}(11) = 1$$
, A(11) = $\frac{1}{2}$ Down $\frac{1}{3}$ / $V^{(3)}(14) = 1$, A(14) = $\frac{1}{2}$ Down, Right $\frac{1}{3}$ $V^{(3)}(15) = 1$, A(15) = $\frac{1}{2}$ Down, Right, Left, Up $\frac{1}{3}$

→ 인근 관의 보자이 귀개에 ILL 성상률이 축가로난다.

t=2



Ş	0	0	0
0	Ho	ł	F
0	1	(, H
Į	ı	ſ	J

4 발교씨 반복제제조 각 칸에서의 가시하는 구 구에는다.

$$v^{(t)}(\eta) = \{ A(\eta) = \{ Down \} \}$$

$$V^{(4)}(14) = 1$$
, $A(14) = 2 D_{own}$, Right, Left, V_{P} $V^{(4)}(15) = 1$, $A(15) = 2 D_{own}$, Right, Left, V_{P} $V^{(4)}(15) = 1$, V_{P}

4	4	1	4
4		1	
t,	Ţ	4	
4	4	4	

S	0	l	0
0	H	1	H
l	ı	(_o H
ſ	_	ſ	J

t=4

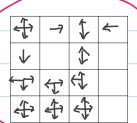
$$V^{(4)}(3) = 1$$
, $A(3) = 2 Down } / V^{(4)}(9) = 1$, $A(9) = 2 Down, Right }$

아니아 03 번복은 전체하던 가음과 받다.

$$\mathcal{V}^{(5)}(2)=1$$
, $A(2)=\{Right\}$ / $\mathcal{V}^{(5)}(4)=1$, $A(4)=\{Left\}$
 $\mathcal{V}^{(5)}(5)=1$, $A(5)=\{Down\}$

t=42+ HUS9 성공의 공유가 작가시는 7유는 또한 다음과 같다.

 $V^{(5)}(10)=1$, $A(10)=\frac{1}{2}$ Down, Right, Left $\frac{1}{2}$ $\int V^{(5)}(3)=1$, $A(3)=\frac{1}{2}$ Down, Up $\frac{1}{2}$ $V^{(5)}(13)=1$, $A(13)=\frac{1}{2}$ Down, Right, Left, Up $\frac{1}{2}$ $\int V^{(5)}(\eta)=1$, $A(\eta)=\frac{1}{2}$ Down, Right, Left $\frac{1}{2}$



이에 때에 코르를 심해하는 꺼라는 다음과 같은 꺼라를 풀었는 그림은 그들은 그는 놀라 같다.



1	7	4	\uparrow
4	4	←	←
1	1	←	(
_	7	1	←

print(optimal_policy)

In [15]: print(optimal_policy)

 $[1.\ 2.\ 0.\ 3.\ 0.\ 0.\ 0.\ 0.\ 3.\ 1.\ 0.\ 0.\ 0.\ 2.\ 1.\ 0.]$

