

① Policy Iteration

$$A = \{B, O, D\}, S = \{N, F, A\}$$

② $|A|^{151} = 3^3 = 27$ Policies

b)

	N	F	A
π_0	0	B	D
π_1	$\frac{-1}{4}$	0	$\frac{1}{100}$
π_2	B	B	B
π_3	B	B	B
	Done ✓		

$$V^{\pi_0}(\text{None}) = \frac{1}{4} - \frac{1}{2} + 0 = -\frac{1}{4}$$

$$V^{\pi_0}(F) = 0$$

$$V^{\pi_0}(A) =$$

$$\pi_1(\text{None}) = \max \{ B: -0.113, O: -0.284, D: 0.131 \}$$

$$\pi_1(F) = \max \{ B: -0.113, O: -0.284, D: 0.131 \}$$

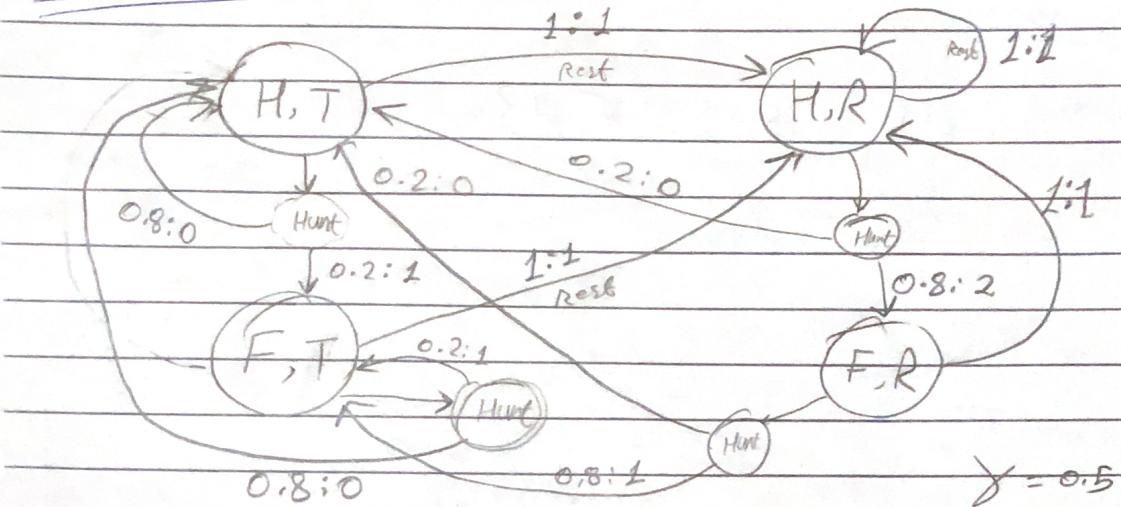
$$\pi_1(A) = \max \{ B: -0.113, O: -0.284, D: 0.131 \}$$

Best Policy is BBB

c) No, the Reward is a lot larger than the γ which

will cause it to converge no matter the γ & also
it's a cyclic graph all point to all

Q) Value Iteration:-



$$\text{a) } V^\pi(HUNgry) = \sum_s \pi(s,a,s) [R(s,a,s') + \gamma V^\pi(s')]$$

$$= 1 \times [1 + 0.5 \times V(H,R)] \\ = 1 + 0.5 V(H,R)$$

$$V(H,R) = 1 \times [0.8 \times 1 + 0.5 V(H,R)] \\ = 1 + 0.5 V(H,R)$$

$$V(H,T) = 1 + 0.5(1 + 0.5 V(H,R)) \\ = 1 + 0.5 + 0.25 V(H,R) \\ = 1 + \gamma + \gamma^2 V(H,R)$$

Geometric series sum

$$V(H,T) = \sum_{i=0}^{\infty} \gamma^i$$

$$V(H,T) = \frac{1}{1-\gamma}$$

#

$$\text{b) } V_1(H, T) = \max \begin{cases} \text{Rest: 1} \\ \text{Hunt: 0.2} \end{cases}$$

= 1

$$V_1(H, R) = \max \begin{cases} \text{Rest: 1} \\ \text{Hunt: 1.6} \end{cases}$$

= 1.6

$$V_1(F, T) = \max \{ \text{Rest: 1}, \text{Hunt: 0.2} \}$$

= 1

$$V_1(F, R) = \max \{ \text{Rest: 1}, \text{Hunt: 0.8} \}$$

= 1

	H, T	H, R	F, T	F, R
No	0	0	0	0
V ₁	1	1.6	1	1
V ₂	2.6	2.6	2.6	2.6
V ₃	3.6	4.2	3.6	3.6
V ₄	5.2	5.2	5.2	5.2

extra # ✓