

## Policies & Value functions

Given a policy  $\pi$ ,

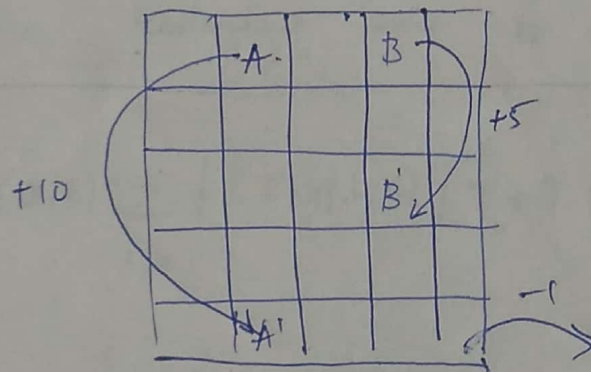
$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

$$\text{(State-Value } f^{\pi}) = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right], \forall s \in S$$

$$q_{\pi}(s, a) = E_{\pi} [G_t \mid S_t = s, A_t = a]$$

$$\text{(Action-value } f^{\pi}) = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

Example:- Grid World:



- The cells of the grid correspond to the states of the MDP
- At each cell, four actions are possible N, S, E & W
- An action that takes the state out of the grid results in "state quo" for the state but gives reward of -1
- In states A & B any action takes one to states A' & B' respectively, with a reward of +10 & +5
- All actions in other states result in a reward of "0"

Suppose the agent selects all four actions with equal distribution prob. Suppose discount factor = 0.9

Let  $V_{\pi}(s)$  be the value  $f^u$  under policy  $\pi$

	A	B		
	3.3	8.8	4.4	5.3
	1.5	3.0	2.3	1.9
	0.1	0.7	0.7	0.4
	-1.0	-0.4	-0.4	-0.6
	-1.9	-1.3	-1.2	-1.4

$A \rightarrow A' \text{ \& } B \rightarrow B'$   
regardless of position of  
 $A \text{ \& } B$ .

Value  $f^u$  table under policy  $\pi$

$$\frac{1}{4} \left[ (0, 1) + \gamma (V_{\pi}(s'_1)) \right]$$

$$0 + \gamma (V_{\pi}(s'_2))$$

$$0 + \gamma$$

$$V_{\pi}(s) = E_{\pi} [r_t + \gamma V_{\pi}(s') | s_t = s]$$

$$= E_{\pi} [R_{t+1} + \gamma V_{\pi}(s') | s_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [\gamma + \gamma V_{\pi}(s')] \quad \forall s \in S$$

Note that,  $V_{\pi}(A) = 8.8 < \text{immediate reward of } 10$

Also,  $V_{\pi}(B) = 5.3 > \text{immediate reward of } 5$

(i) Suppose we know that  $V_{\pi}(A') = -1.3$ , what is  $V_{\pi}(A)$

$$V_{\pi}(A) = \frac{1}{4} \left[ (0, 1) + \gamma V_{\pi}(A') \right] = 10 + 0.9(-1.3) = \approx 8.8$$

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

(ii) Suppose wkt,  $V_{\pi}(B') = 0.4$

$$V_{\pi}(B) = \frac{1}{4} \left[ (0, 1) + \gamma V_{\pi}(B') \right]$$

$$= 5 + 0.9 \times 0.4 \approx 5.4$$

## Center State

	2.3	
0.7	0.7	0.4
	0.4	

there are equal chance from each state to C

$$\frac{1}{4}(0) + \frac{1}{4}(0) + \frac{1}{4}(0) + \frac{1}{4}(0)$$

$$V_{\pi}(C) = \underbrace{0}_{\text{Immediate}} + 0.9 \left[ \frac{1}{4} \times 2.3 + \frac{1}{4} \times 0.4 + \frac{1}{4} \times (-0.4) + \frac{1}{4} \times 0.7 \right]$$

$$= 0.9 \times 0.75$$

$$\approx 0.7$$

corner:

$V_{\pi}(E)$

	1.2	
-1.4	-2.0	-1
	-1	

outside

$$V_{\pi}(E) = \frac{1}{4} \times (-1) + \frac{1}{4} \times (-1) + \frac{1}{4} (0) + \frac{1}{4} (0)$$

$$+ 0.9 \left[ \frac{1}{4} \times (-1.4) + \frac{1}{4} \times (-1.2) + \frac{1}{2} (V_{\pi}(E)) \right]$$

staying in same state if action is outside grid.

$$V_{\pi}(E) = -0.5 - 0.585 + 0.45 \times (V_{\pi}(E))$$

$$0.55 V_{\pi}(E) = -1.085$$

$$V_{\pi}(E) \approx -2.02$$

Ex:

	1.5	
1.9	0.5	
	-0.4	

D

$$V_{\pi}(D) = \frac{1}{4} (-1) + 0.9 \left( \frac{1}{4} (1.5) + \frac{1}{4} (1.9) + \frac{1}{4} (0.4) + \frac{1}{4} V_{\pi}(D) \right)$$

$$= -0.25 + 0.9 \left( \frac{1}{4} (3) + \frac{1}{4} V_{\pi}(D) \right)$$

$$0.75 V_{\pi}(D) = -0.25 + 0.675 (-1 + 2.7) \Rightarrow \frac{3.1}{4} (V_{\pi}(D)) \Rightarrow 0.566$$



Bellman Eqn for  $q_{\pi}(s, a)$

$$q_{\pi}(s, a) = E_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s, A_t = a]$$

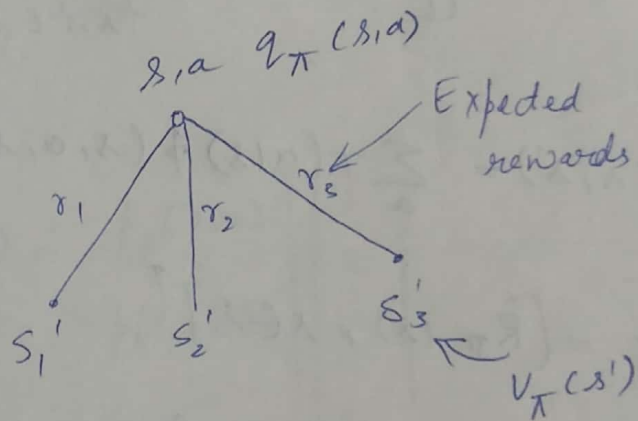
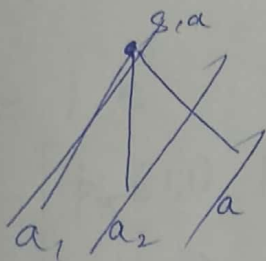
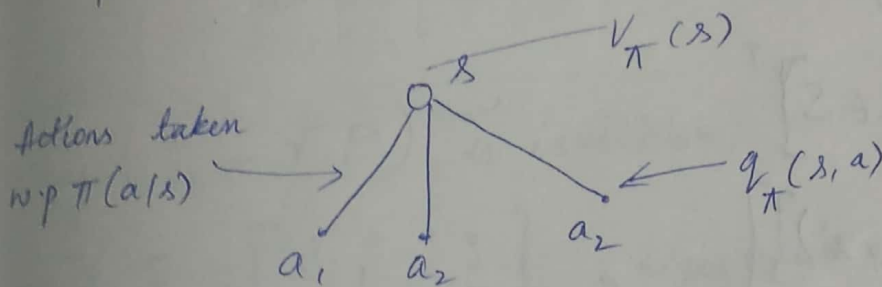
$$= \sum_{s', r} p(s', r | s, a) (\gamma + r + \gamma V_{\pi}(s'))$$

State Value

Action-value

How is  $V_{\pi}(s)$  related to  $q_{\pi}(s, a)$ ?

$$V_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a)$$



	$A_2$	$A_3$	$A_4$
$A_1$			
$A_2$			
$A_3$			
$A_4$			

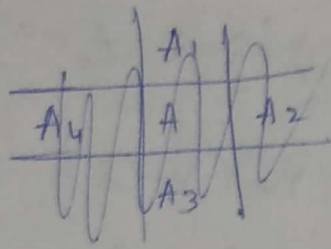
$$V_{\pi}(A) = \frac{1}{4} \times r \times V_{\pi}(A_1) +$$

$$\frac{1}{4} \times r \times V_{\pi}(A_2) +$$

$$\frac{1}{4} \times r \times V_{\pi}(A_3) +$$

$$\frac{1}{4} \times r \times V_{\pi}(A_4)$$

## Optimal Policies & optimal Value functions



Value functions define a partial ordering over policies. A policy  $\pi$  is defined to be better or equal

to policy  $\pi'$  if the expectation return under policy  $\pi$  is greater than or equal to the expected return under policy  $\pi'$

$$\Rightarrow \pi \succeq \pi' \quad \text{i/f} \quad V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s \in S$$

$$V_{\pi} = [V_{\pi}(s), s \in S]^T$$

$$P_{\pi} = \left[ \left[ p_{\pi}(s, s') \right] \right]_{s, s' \in S}$$

$$p_{\pi}(s, s') = \sum_a \pi(a|s) p(s', a, s')$$

$$R_{\pi} = [R_{\pi}(s), s \in S]^T$$

$$V_{\pi} = R_{\pi} + \gamma P_{\pi} V_{\pi}$$

$$\text{or } (I - \gamma P_{\pi}) V_{\pi} = R_{\pi}$$

$$\text{or } V_{\pi} = (I - \gamma P_{\pi})^{-1} R_{\pi}$$

$$R_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) r$$

## Value Iteration

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$V_{\pi}(s) = R_{\pi}(s) + \gamma \sum_{s'} P_{\pi}(s, s') V_{\pi}(s')$$

Initialize  $V_0(s) = 0 \quad \forall s$

Iterate:

$$V_{n+1}(s) = R_{\pi}(s) + \gamma \sum_{s'} P_{\pi}(s, s') V_n(s'), \quad n \geq 0$$

as  $n \rightarrow \infty$ ,  $V_n(s) \rightarrow V_{\pi}(s) \quad \forall s \in S$

Why  $(I - \gamma P)$  is invertible,

$$I = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \quad \text{eigen values of } I = 1$$

$$P_{\pi} \mathbf{e} = 1 \cdot \mathbf{e}$$

$$\begin{bmatrix} P_{\pi}(1,1) & P_{\pi}(1,2) & \dots & P_{\pi}(1,n) \\ P_{\pi}(2,1) & \dots & \dots & P_{\pi}(2,n) \\ \vdots & \vdots & \ddots & \vdots \\ P_{\pi}(n,1) & \dots & \dots & P_{\pi}(n,n) \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$I - \gamma P$  is always  $\neq 0$  since  $P=1$  &  $\gamma P < 1$   
 $\therefore 1 - \gamma P > 0 \therefore I - \gamma P$  has no zero  
e.v hence invertible.

It turns out that there is always at least one policy that is better than or equal to all other policies. That policy will be optimal policy.

We denote all optimal policies by  $\pi_*$ .

The state value  $f^*$  for each optimal policy  $\pi_*$  is the same - we ~~also~~ denote it by  $V_*$

$$V_*(s) = \max_{\pi} J_{\pi}(s), \forall s \in S.$$

Optimal policies also share same optimal action value  $q^*$

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

$$q^*(s, a) = E[R_{t+1} + \gamma V_*(s_{t+1}) | S_t = s, A_t = a]$$

~~for~~

$V^* \equiv$  optimal value function

Under a given policy  $\pi$ , value function  $V_{\pi}$  satisfies the Bellman equation.

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}[r_t | S_t = s] \\ &= \sum_{a \in A(s)} \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')] \end{aligned}$$



we expect the optimal value function  $V^*$  to also satisfy some form of the Bellman equation.

$$\begin{aligned} V^*(s) &= \max_{a \in A(s)} q_{\pi^*}^*(s, a) \\ &= \max_{a \in A(s)} E_{\pi^*} [q_t | s_t = s, A_t = a] = \max_{a \in A(s)} E [R_{t+1} + \gamma V_{t+1}^* | s_t = s, A_t = a] \\ &= \max_{a \in A(s)} E_{\pi^*} [R_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, A_t = a] \end{aligned}$$

$$V^*(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')] \quad \text{B.E for state } s$$

The Bellman optimality equation for  $q^*$  can similarly be written as

$$q^*(s, a) = E [R_{t+1} + \gamma \max_{a'} q^*(s_{t+1}, a') | s_t = s, A_t = a]$$

This follows from,

$$q^*(s, a) = E [R_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, A_t = a]$$

$$\text{Now } V^*(s_{t+1}) = \max_{a'} q^*(s_{t+1}, a')$$

$$q^*(s, a) = \sum_{s', r} p(s', r | s, a) \left[ r + \gamma \max_{a'} q^*(s', a') \right] \quad \text{B.E for action value } f^a$$

$$\text{Note that, } V^*(s) = \max_{a \in A(s)} q^*(s, a)$$

$$= \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) \left( r + \gamma \max_{a'} q^*(s', a') \right)$$

$$V^*(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) \left( r + \gamma V^*(s') \right)$$



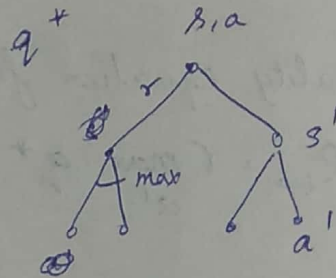
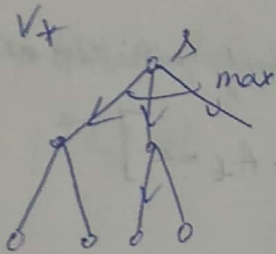
Suppose we let,

$$q^*(s,a) = \sum_{s',r} p(s',r|s,a) (r + \gamma v^*(s'))$$

then,  $v^*(s) = \max_a q^*(s,a)$

$$\Rightarrow q^*(s,a) = \sum_{s',r} p(s',r|s,a) (r + \gamma \max_{a'} q^*(s',a'))$$

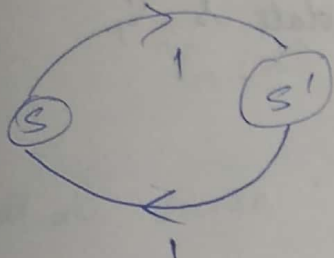
BE for action value function



In RL setup  $q^*(s,a)$  is more favourable since the max is within the expectation. Since the system is unknown (Stochastic Approx Rule),  $E$  of system is not known ~~the~~ hence  $E \text{ sys max}(E)$  is possible to calculate. Whereas  $E(\max())$  is possible via  $q^*$

There always exists a policy that is better than or equal to other policies.

Eg: Consider the following MDP with transitions



Under two different policies  $\pi$  &  $\pi'$   
assume  $\pi$  &  $\pi'$  are deterministic policies for simplicity.

The difference in the policies is in terms of rewards obtained.

Note that,

$$V_{\pi}(s) = r_{\pi}(s) + \gamma V_{\pi}(s') \quad \text{--- (1)}$$

$$V_{\pi'}(s) = r_{\pi'}(s) + \gamma V_{\pi'}(s') \quad \text{--- (2)}$$

Suppose that,

$$V_{\pi}(s) > V_{\pi'}(s) \quad \text{--- (3)}$$

$$\text{but } V_{\pi}(s') < V_{\pi'}(s') \quad \text{--- (4)}$$

From (1), (2) & (3), we have

$$\begin{aligned} r_{\pi}(s) + \gamma V_{\pi}(s') &> r_{\pi'}(s) + \gamma V_{\pi'}(s') \quad \text{from (3)} \\ \Rightarrow r_{\pi}(s) &> r_{\pi'}(s) - (5) \because V_{\pi}(s') < V_{\pi'}(s') \quad \text{from (4)} \end{aligned}$$

Now similar to (1)-(2) we can write,

$$V_{\pi}(s') = r_{\pi}(s') + \gamma V_{\pi}(s) \quad \text{--- (6)}$$

$$V_{\pi'}(s') = r_{\pi'}(s') + \gamma V_{\pi'}(s) \quad \text{--- (7)}$$

$$R_{\pi}(s') + \gamma V_{\pi}(s) < R_{\pi'}(s') + \gamma V_{\pi'}(s) \quad \text{--- (8)}$$

$$V_{\pi}(s) > V_{\pi'}(s)$$

$$\Rightarrow R_{\pi}(s') < R_{\pi'}(s') \quad \text{--- (9)}$$

It is optimal to pick action  $\pi(s)$  in state  $s$  &  
action  $\pi'(s')$  in state  $s'$

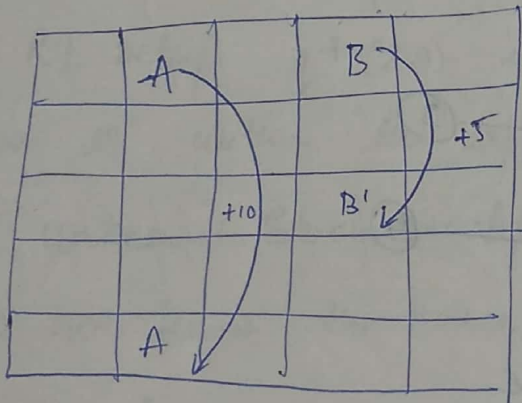
Let  $\pi''$  be another policy st

$$\pi''(s) = \pi(s)$$

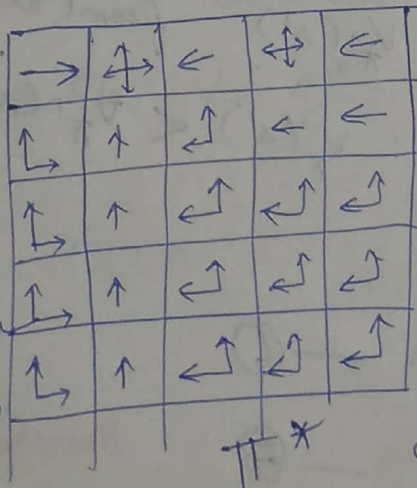
$$\& \pi''(s') = \pi'(s')$$

Thus one can construct a policy that is better than or  
equal to all other policies. In this case,  $\pi''$

Grid World Example:



22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7



which  
max  $V^*$   $\rightarrow V^*$

$$V(s) = R + \gamma V(s')$$

$$0 + 0.9$$

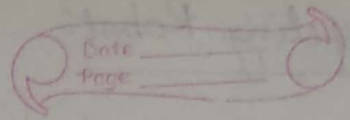
both action  
& if  
are optimal

$\pi^*$

so action need not be  
unique.



Init  $V_0(s) \quad \forall s \in S$



Iterate

$$V_{n+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma V_n(s')) \quad n \geq 0$$

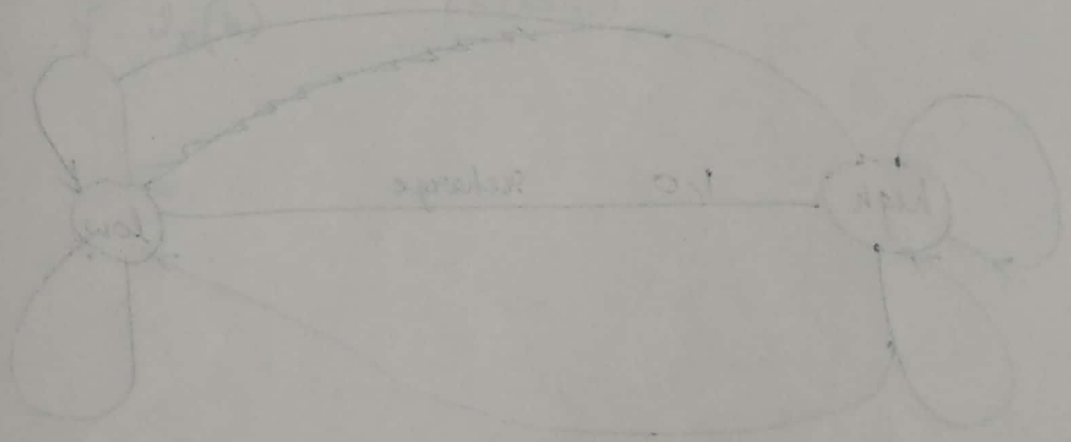
$$V_n(s) \longrightarrow V^*(s) \text{ as } n \longrightarrow \infty$$

Exercise: Implement value iterate for the grid example.

stop when,

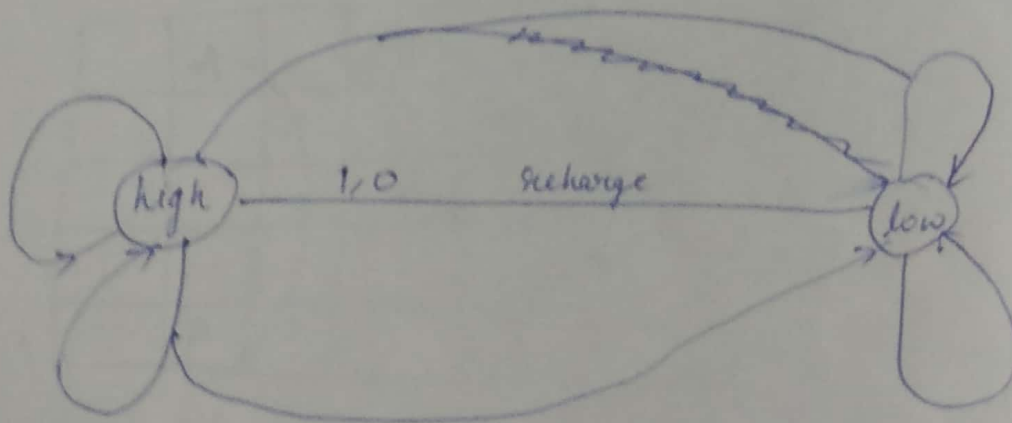
$$\max_s |V_n(s) - V_{n-1}(s)| < 0.1$$

Output  $V_n(s) \quad \forall s.$



# Recycling Robot:

$s$	$a$	$s'$	$p(s' s,a)$	$r(s,a,s')$
high	search	high	$\alpha$	$R_{search}$
high	search	low	$1-\alpha$	$R_{search}$
low	search	high	$1-\beta$	$-3$
low	search	low	$\beta$	$R_{search}$
high	wait	high	1	$R_{wait}$
high	wait	low	0	$R_{wait}$
low	wait	high	0	$R_{wait}$
low	wait	low	1	$R_{wait}$
low	recharge	high	1	0
low	recharge	low	0	0



Bellman optimality eqn for the robot

$$high \equiv h, \quad low \equiv l$$

$$search \equiv s, \quad wait \equiv w, \quad recharge \equiv r$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$V_*(h) = \max \left\{ \begin{aligned} & p(h|h,s) [r(h,s,h) + \gamma V_*(h)] + \\ & p(l|h,s) [r(h,s,l) + \gamma V_*(l)] + \\ & p(h|h,w) [r(h,w,h) + \gamma V_*(h)] + \\ & p(l|h,w) [r(h,w,l) + \gamma V_*(l)] \end{aligned} \right\}$$

$$= \max \left( \begin{aligned} & \alpha [r_s + \gamma V_*(h)] + (1-\alpha) [r_s + \gamma V_*(l)], \\ & 1 [r_w + \gamma V_*(h)] + 0 [r_w + \gamma V_*(l)] \end{aligned} \right)$$

$$= \max \left\{ \begin{aligned} & r_s + \gamma [\alpha V_*(h) + (1-\alpha) V_*(l)], \\ & r_w + \gamma V_*(h) \end{aligned} \right\}$$

$$V_*(l) = \max \left\{ \begin{aligned} & \beta r_s - \beta(1-\beta) + \gamma [(1-\beta) V_*(h) + \beta V_*(l)], \quad (\text{search}) \\ & r_w + \gamma V_*(l), \quad (\text{wait}) \\ & \gamma r_s(h) \quad (\text{recharge}) \end{aligned} \right\}$$