

$$A3) R = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \gamma^3 R_{t+3} + \dots$$

$\therefore C \rightarrow \text{Constant}$

$$\tilde{R} = (R_{t+1} + C) + \gamma(R_{t+2} + C) + \dots$$

Then  $V_{\pi}(s) = E[R|s]$

$$\tilde{V}_{\pi}(s) = E[\tilde{R}|s]$$

$$= E[(R_{t+1} + C) + \gamma(R_{t+2} + C) + \gamma^2(R_{t+3} + C) | s]$$

$$= E\left[\sum_{k=1}^{\infty} \gamma^k R_{t+k} | s\right] + E\left[C \sum_{k=0}^{\infty} \gamma^k | s\right] \quad \text{where } 0 < \gamma < 1$$

$\uparrow$  constant

$$= E[R|s] + C \sum_{k=0}^{\infty} \gamma^k$$

$$\boxed{\tilde{V}_{\pi}(s) = V_{\pi}(s) + \frac{C}{1-\gamma}}$$

$$\therefore \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

So relative evaluation does not change.

for episodic Task,

$$\tilde{V}_{\pi}(s) = V_{\pi}(s) + C \sum_{k=0}^T \gamma^k$$

$$\boxed{\tilde{V}_{\pi}(s) = V_{\pi}(s) + \frac{C(1-\gamma^T)}{1-\gamma}}$$

hence as we go closer to terminal state the change in  $\tilde{V}_{\pi}(s)$  increases

$$A5) V_*(s) = \max_{\pi} V(s) \quad \forall s \quad \text{--- ①} \quad q_*(s) = \max_{\pi} q_{\pi}(s, a) \quad \text{--- ②}$$

as

$$V_{\pi}(s) = q_{\pi}(s, \pi(s)).$$

$$V_*(s) = \max_{\pi} q_{\pi}(s, \pi(s)). \quad V_*(s) = \max_{\pi} q_{\pi}(s, \pi(s)).$$

$$V_*(s) = \max_{\pi} q_{\pi}(s, \pi(s)).$$

$$V_*(s) = \max_{\pi} q_{\pi}(s, \pi(s)).$$

$$So \quad V_*(s) = \max_a q_*(s, a) \quad \max_{\pi} \max_a q_{\pi}(s, a)$$

$$V_*(s) = \max_a q_*(s, a)$$

considered them to be constants  
from function T.

A1)

S	a	S'	a	$P(S', a   S, a)$
high	search	high	$a_{search}$	$\alpha$
high	search	low	$a_{search}$	$1 - \alpha$
low	search	high	-3	$1 - \beta$
low	search	low	$a_{search}$	$\beta$
high	wait	high	$a_{wait}$	1
low	wait	low	$a_{wait}$	1
low	recharge	high	0	1

\* Eq. As  $P(a | S, S') = 0$  for each tuple  $(S, a, S')$ , only 1 reward exists  
 $\downarrow$   
 T

$$\sum_a P(S', a | S, a) = P(S' | a, S)$$

so

$$P(S', T(S, a, S') | S, a) = P(S' | a, S)$$