



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

Explanation,

using the figure above, when is state high, using action search, with the prob of α we go to state high and get reward r_{search} .

3.15 We know,

$$G_t = \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c)$$

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

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c) \mid S_t = s \right]$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] + E_{\pi} \left[\sum_{k=0}^{\infty} c \gamma^k \mid S_t = s \right]$$

$$\text{let } V_c = E \sum_{k=0}^{\infty} c \gamma^k$$

$$\boxed{V_c = \frac{c}{1-\gamma}}$$

thus adding c doesn't relative values of any state

3.16. If we have episodic task, then the value of V_c would change

$$V_c = \sum_{k=0}^1 c \gamma^k = c \left[\frac{1-\gamma^k}{1-\gamma} \right]$$

hence same value would be added and thus task would be left unchanged.

5. We know,

$$V_*(s) = \max_{\pi} V_{\pi}(s) \quad \forall s$$

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

} optimal values

using the above equations

$$V_*(s) = \max_{a \in A(s)} q_*(s, a)$$