

Q3. Ex 3.15

$$G_t = \sum_{k=1}^{\infty} \gamma^{k-1} \cdot R_{t+k}$$

If we add  $c$  to each reward then the new  $G_t$  will be:

$$G'_t = \sum_{k=1}^{\infty} \gamma^{k-1} (R_{t+k} + c)$$

$$= \sum_{k=1}^{\infty} \gamma^{k-1} R_{t+k} + \sum_{k=1}^{\infty} \gamma^{k-1} \cdot c$$

$$= G_t + c \left( \frac{1}{1-\gamma} \right) \quad (\text{GP formula})$$

Taking expectation on both sides ~~and~~ conditioned on state  $s$

$$E[G'_t | s] = E\left[G_t + \frac{c}{1-\gamma} \mid s\right]$$

$$\Rightarrow V_{\pi}(s) = E[G_t | s] + \frac{c}{1-\gamma} \quad (\text{as } c \text{ \& } \gamma \text{ are constants})$$

$$V'_{\pi}(s) = V_{\pi}(s) + \frac{c}{1-\gamma} \quad \forall s \text{ in states.}$$

Thus, adding constant  $c$  adds a value of  $V_c = \frac{c}{1-\gamma}$  to ~~all~~ values of all states.

Q3. Ex 3.16

considers that a state  $s$  ends in  $T$  steps so:

$$V_{\pi}(s) = \sum_{k=1}^T \gamma^k R_{t+k}$$

$$V'_{\pi}(s) = \sum_{k=1}^T \gamma^k (R_{t+k} + c) = \sum_{k=1}^T \gamma^k \cdot c + V_{\pi}(s)$$

$$\boxed{V'_{\pi}(s) = V_{\pi}(s) + \frac{c(1-\gamma^T)}{(1-\gamma)}}$$

Now here the constant term added for each state  $s$  depends on the no. of timesteps  $T$  after which the terminal state is reached.

Since these timesteps can be diff for each state hence the value func<sup>n</sup> of each state will increase by ~~the~~ different amount and thus adding  $c$  to all rewards will have an effect & it will give preference to states ~~which~~ whose episodes last longer (as  $1-r^T$  will be more for these states).

Q1. Ex. 3.4

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

The table is obtained using the formula<sup>s</sup>

$$p(s' | s, a) = \sum_{r \in \text{rewards}} p(s', r | s, a)$$

Since we have been provided with no prob<sup>s</sup> distribution but only given expected values of rewards for each state so we'll assume that that reward comes with prob = 1 & other with prob = 0



$$\begin{aligned}
 R(s, a, s') &= E[R | s, a, s'] \\
 &= \sum_r R \cdot p(r | s, a, s')
 \end{aligned}$$

So for low, search, low,  $R_{\text{search}}$ :

$$R(\text{low}, \text{search}, \text{low}) = \sum_r R \cdot p(r | \text{low}, \text{search}, \text{low})$$

$$R_{\text{search}} = \frac{\sum_r R \cdot p(r, \text{low} | \text{low}, \text{search})}{p(\text{low} | \text{low}, \text{search})}$$

$$R_{\text{search}} = \frac{R_{\text{search}} \cdot p(r, \text{low} | \text{low}, \text{search})}{\beta}$$

$$p(r, \text{low} | \text{low}, \text{search}) = \beta$$

Q5. we know that by def<sup>n</sup>  $V^*$  is the optimal state value func<sup>n</sup> and  $q^*$  is the optimal state action value func<sup>n</sup>, so from a given state,  $V^*$  is the max of all possible  $q^*$  possible from that state (based on all diff actions).

hence 
$$V^*(s) = \max_a q^*(s, a)$$