==Question 1==

| s | a | s' | r | $r(s,a,s')$ | $p(s'\circ|s,a)$ | $p(s',\sigma|s,a)$ |
|---|---|---|---|---|---|---|
| high | search | high | 10 | $r_s$ | $\alpha$ | $+0.1\,\alpha\cdot r_s$ |
| high | search | low | -10 | $r_s$ | $1-\alpha$ | $-0.1(1-\alpha)\,r_s$ |
| low | search | high | 0 | -3 | $1-\beta$ | 0 |
| low | search | low | -10 | $r_s$ | $\beta$ | $-0.1\beta\,r_s$ |
| high | wait | high | 10 | $r_w$ | 1 | $0.1\,r_w$ |
| high | wait | low | 10 | $\cos$ - | 0 | - |
| low | wait | high | 0 | - | 0 | - |
| low | wait | low | 10 | $r_{wait}$ | 1 | $0.1\,r_w$ |
| low | sech. | high | 0 | 0 | 1 | - |
| low | sech. | low | 0 | - | 0 | - |

Since $r(s,a,s') = \dfrac{\sum r\; p(s',r|s,a)}{p(s'|s,a)}$

① $r_s = \dfrac{10\cdot p(s',r|s,a)}{p(high|high, search)} \to \alpha$

$p(s',r|s,a) = \alpha\cdot 0.1\cdot r_s$

② $r_s = \dfrac{-10}{(1-\alpha)}\, p(s',r|s,a)$

$p(low, -10|high, search) = -0.1\cdot(1-\alpha)\,r_s.$

**Question 2 -** Given the gridworld environment,

Solving the Bellman equation, $v_\pi(s)$ :

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

Solving for state 1 for illustration -

$$v_\pi(1) = 0.25 \times 1 \ [-1 + \gamma v_\pi(1)] + $$
$$0.25 \times 1 \ [-1 + \gamma v_\pi(1)] + $$
$$0.25 \times 1 \ [0 + \gamma v_\pi(2)] + $$
$$0.25 \times 1 \ [0 + \gamma v_\pi(4)]$$

Similarly $v_\pi(2) = 0.25 \times 1 \ [-1 + \gamma v_\pi(2)] +$
$$0.25 \times 1 \ [0 + \gamma v_\pi(1)] +$$
$$0.25 \times 1 \ [0 + \gamma v_\pi(3)] +$$
$$0.25 \times 1 \ [0 + \gamma v_\pi(4)]$$

Thus
$$\begin{bmatrix} v_\pi(1) \\ v_\pi(2) \\ \vdots \\ \vdots \\ v_\pi(25) \end{bmatrix}_{25 \times 1} = \begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ \vdots \\ R_{25} \end{bmatrix}_{25 \times 1} + \gamma \underbrace{P_\pi}_{25 \times 25} \underbrace{v_\pi}_{25 \times 1}$$

$\downarrow$ $\downarrow$
$25 \times 1$

$\downarrow$
Total no.
of states

$$\therefore \quad v_\pi = R_\pi + \gamma P_\pi v_\pi$$
$$v_\pi - \gamma P_\pi v_\pi = R_\pi$$
$$\boxed{v_\pi = (I - \gamma P_\pi)^{-1} R_\pi} \qquad \forall s \in S.$$

**Question 3**

**Exercise 3.15**
$$v_\pi(s) = E_\pi \left[ R_t \mid S_t = s \right]$$

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

Adding a constant to each reward—

let $\hat{R}_{t+k+1} = R_{t+k+1} + c$

$$\hat{v}_\pi(s) = E_\pi \left[ \gamma^k \hat{R}_{t+k+1} \mid S_t = s \right]$$

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

$$= v_\pi(s) + \frac{c}{1-\gamma}$$

Thus we see the best policy doesn't change with addition of constant $c$ at each reward.

**Exercise 3.16** Adding a constants at episodic tasks:

$$v_\pi(s) = E \left[ R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{k-1} R_{t+k} \right]$$

$$\hat{v}_\pi(s) = E \left[ (R_{t+1} + c) + \gamma (R_{t+2} + c) + \cdots + \gamma^{k-1} (R_{t+k} + c) \right]$$

$$= c \left[ 1 + \gamma + \gamma^2 + \cdots + \gamma^{k-1} \right] + v_\pi(s)$$

$$= \underbrace{\frac{c(1-\gamma^k)}{1-\gamma}}_{(A)} + v_\pi(s)$$

$k \rightarrow$ small value, (A) is closer to terminal state, thus adding a constant only ~~adds~~ prolongs / reduces the

steps to terminal states.

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

Thus,

$$v^*(s) \geq R(s) + \gamma \max_{a \in A} \sum p(s' | s, a) \, v(s')$$

thus thinking $|A|$ as the linear constraints;

$$v(s) \geq R(s) + \gamma \sum_{s' \in S} p(s' | s, a) \, v(s') \quad \forall a \in A$$

Thus, minimise $\sum_s v(s)$
$v$

subject to $v(s) \geq R(s) + \gamma \sum_{s' \in S} p(s' | s, a) \, v(s') \quad \forall a \in A$
$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad s \in A$

Solving for optimal policy,

if we optimise,

$$\min_v \sum_s d(s) \, v(s)$$

Subject to $v(s) \geq R(s) + \gamma \sum_{s' \in S} p(s' | s, a) \, v(s') \quad \forall a \in A,$
$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad s \in S$

Here $d(s)$ is the distribution over states.

Adding $\mu(s, a)$ for each constraint,

$$\max_{\mu(s, a)} \sum_{s \in S} R(s) \sum_{a \in A} \mu(s, a)$$

Subject to $\sum_{a \in A} \mu(s', a) = d(s') + \sum_{s \in S} \sum_{a \in A} p(s' | s, a) \, \mu(s, a) \quad \forall s' \in S.$

where $\mu(s, a) = \sum_{t=0}^{\infty} \gamma^t \, P(S_t = s, A_t = a)$

Thus $\pi^*(s) = \max_{a \in A} \mu(s, a)$

**Notes**             **Appointment**

## Question 5

$$v_\pi(s) = E_\pi \left[ R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s \right]$$

$$v_*(s) = E_{\pi_*} \left[ R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s \right]$$

$$= \max_a E \left[ R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s \right]$$

$$= \max_a \sum_\lambda \sum_{s'} \left( \lambda + \gamma v_*(s') \right) p(s', \lambda \mid s, a)$$

Thus,

$$\boxed{v_*(s) = \max_a q_*(s,a)}$$

value of each state when taken an
action a.

## Question 6

The bug described in question 4 $\downarrow$
to find a way such that policy doesn't keep
on switching in case of finding multiple policies

In this case,
pseudo code

    if $v_\pi(s) << \theta$:
            don't update the $v_\pi(s)$
            $\theta$ increase by a small value
    to break continuous update of states.