$$V_\pi(s) = E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$T_1 = E_\pi[R_{t+1} | S_t = s] \quad ; \quad T_2 = \gamma E_\pi[G_{t+1} | S_t = s]$$

$$V_\pi(s) = T_1 + T_2.$$

$$T_1 = \sum_r r\, p(r|s) = \sum_{s'} p(s'|s) \sum_r r\, p(r|s, s') \text{ —— } \textcolor{green}{\text{May not be reqd.}}$$

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

$$\sum_a \pi(a|s)\, p(s', r|s, a)$$

$$T_2 = \gamma \sum_{s'} p(s'|s)\, V_\pi(s') \text{ — from last slide of Lec 7.}$$

$$\sum_{s'} \sum_a \sum_r \pi(a|s)\, p(s', r|s, a) \times \gamma V_\pi(s')$$
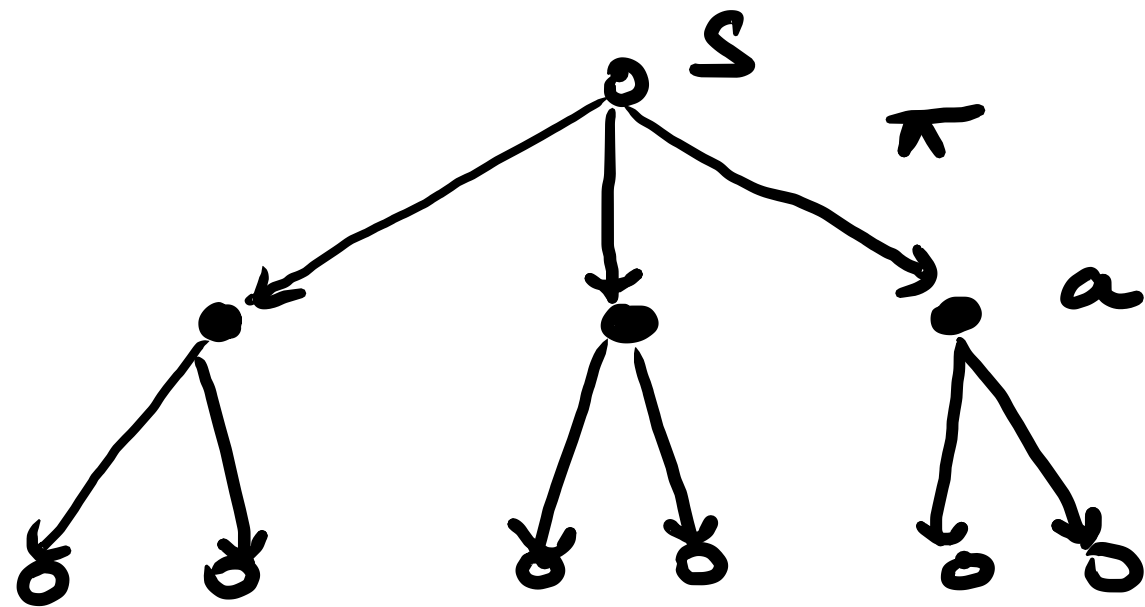
Let take $S = \{ s_1, s_2, s_3, s_4 \}$

$V_\pi(s_1) \rightarrow V_\pi(s_2), V_\pi(s_3), V_\pi(s_4), V_\pi(s_1)$

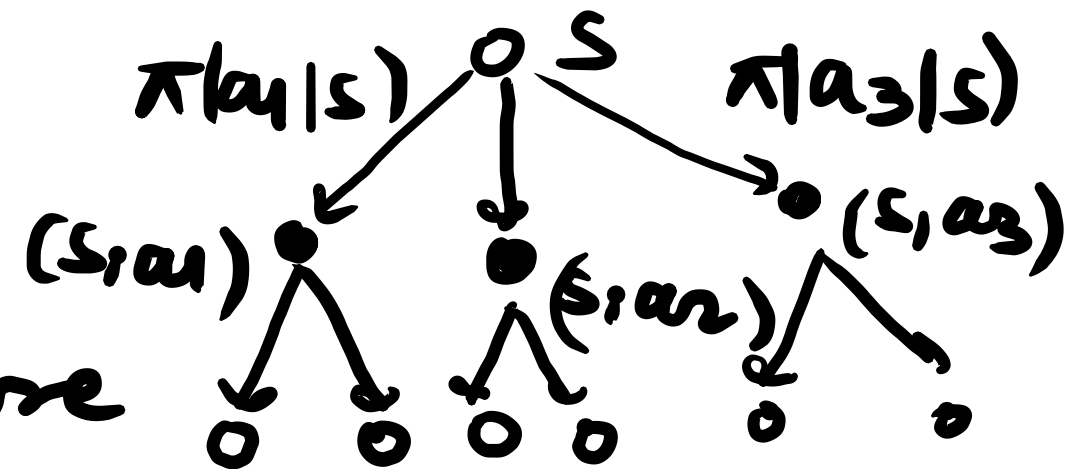$V_\pi(s_2) \rightarrow V_\pi(s_1), V_\pi(s_3), V_\pi(s_4), V_\pi(s_2)$

$\vdots$

because the environment can take the agent from the current state to the same state. (self-loop)

Back up diagram for $V_\pi$



in this example, let the no. of possible actions in $s$ be 3, $a_1, a_2$ & $a_3$.

after taking an action, suppose you can go to 2 possible states

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

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

$\longrightarrow \sum_a \pi(a|s) q_\pi(s,a)$

1. find $V_\pi(s)$ in terms of $q_\pi(s,a)$. $\longrightarrow$ H.W.

2. find $q_\pi(s,a)$  ,,  ,, $V_\pi(s)$ & $p(s', r | s, a)$.

In the backup diagram, each open circle represents a state, while each solid circle rep. a state-action pair.

Bellman equality :- states that the value of the state you are currently in = (discounted) value of the expected next state plus the reward expected along the way.

The state-value function $V_\pi(s)$ is the 'unique solution' to its Bellman equation (BE).

## optimal policies & optimal value functions.

Aim (roughly) is to find a policy that achieves a lot of reward over the long run.

Defn:- A policy $\pi$ is defined to be better than or equal to a policy $\pi'$, if its expected return is $\geq$ that of $\pi'$

<span style="color:red">∀ states</span>

$$\pi \geq \pi' \quad \text{iff} \quad V_\pi(s) \geq V_{\pi'}(s)$$
$$\forall s \in S.$$