

1. To show this expression holds $\forall i \in \{1, \dots, H\}$, we proceed by mathematical induction:

Base case ($i=H$): As $\hat{V}_{H+1}(s) = 0 \forall s \in S$,

$$\hat{V}_H(s) - V_H(s) = \sum_{t=H}^H E[\hat{r}(s_t, a_t) - r(s_t, a_t) | S_H = s]$$

$$= E[\hat{r}(s_H, a_H) - r(s_H, a_H)]$$

Hence, this checks out trivially.

Inductive step: Suppose the expression holds for $i = T+1$ (inductive hypothesis). We now show it holds for $i = T$:

$$\begin{aligned} \hat{V}_T(s) - V_T(s) &= \sum_{t=T}^H E[\hat{r}(s_t, a_t) - r(s_t, a_t) | S_T = s] \\ &= E[\hat{r}(s_T, a_T) - r(s_T, a_T)] + \sum_{t=T+1}^H E[\hat{r}(s_t, a_t) - r(s_t, a_t) | S_T = s] \end{aligned}$$

Now, by Bellman equation, this becomes

$$E[\hat{r}(s_t, a_t) - r(s_t, a_t)] + \sum_{s'} (P(s'|s_t, a_t) \hat{V}_{T+1}(s') - P(s'|s_t, a_t) V_{T+1}(s'))$$

2nd term can be re-written as

$$\sum_{s'} (\hat{P} - P) \hat{V}_{T+1}(s') + \sum_{s'} P (\hat{V}_{T+1}(s') - V_{T+1}(s'))$$

Now making the substitution via inductive hypothesis, the 2nd term above becomes

$$\sum_{s'} P(s'|s_t, a_t) \sum_{t=T+1}^H E[\hat{r}(s_t, a_t) - r(s_t, a_t) + \sum_{s'} (\hat{P} - P) \hat{V}_{t+1}(s') | s = s']$$

$$= \sum_{t=T+1}^H E[\hat{r}(s_t, a_t) - r(s_t, a_t) + \sum_{s'} (\hat{P} - P) \hat{V}_{t+1}(s') | s = s']$$

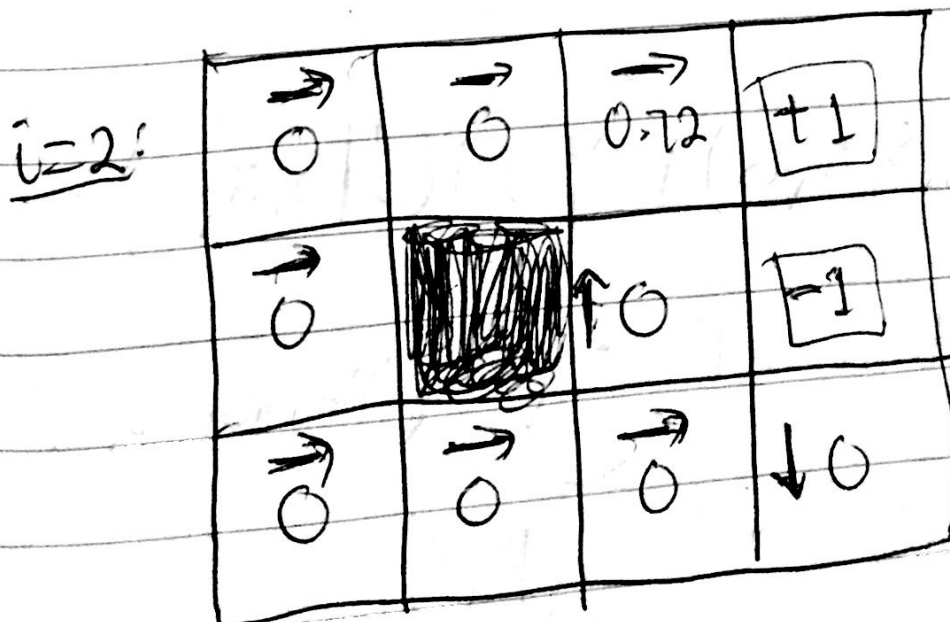
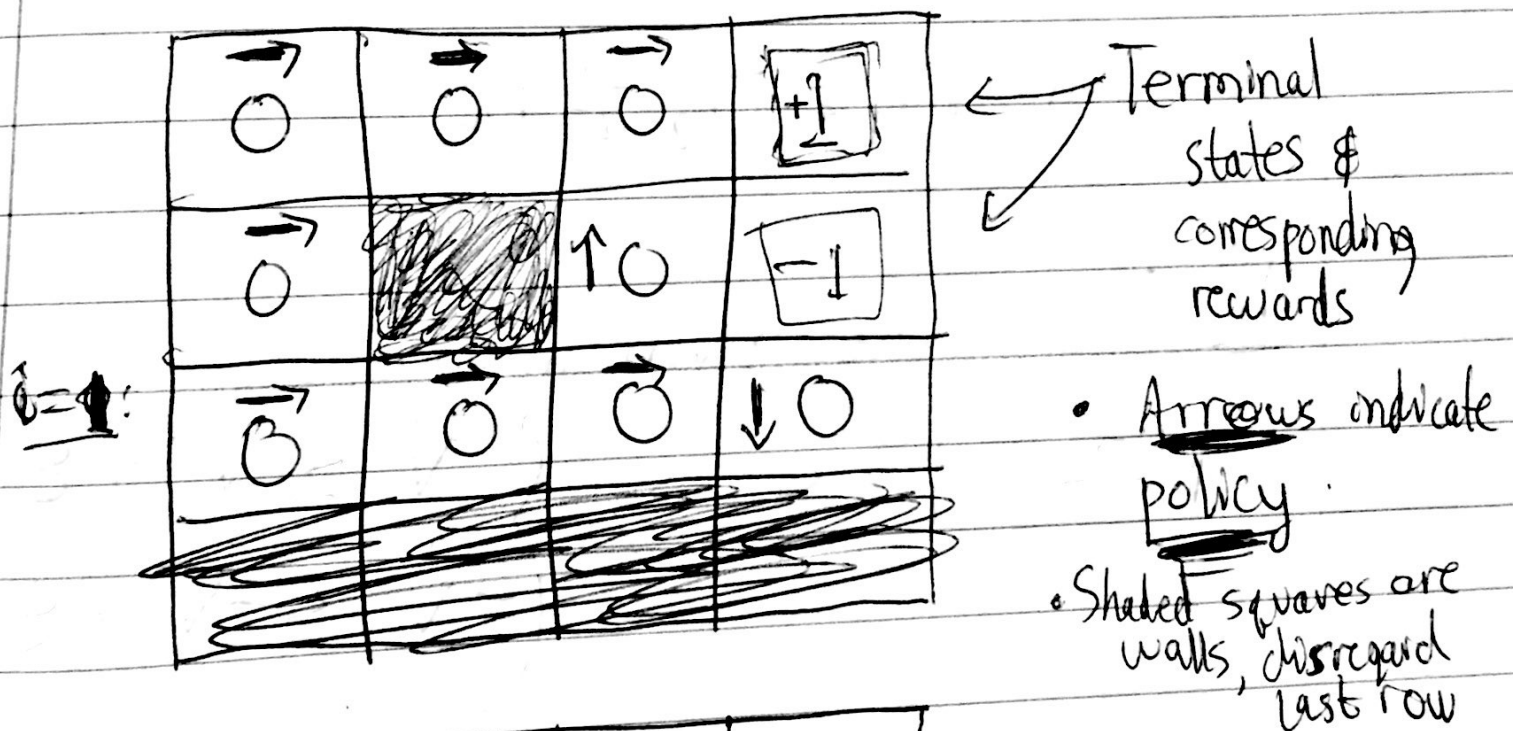
Adding this to $\sum_{s'} (\hat{P} - P) \hat{V}_{T+1}(s') + E[\hat{r}(s_T, a_T) - r(s_T, a_T)]$

we obtain the desired result

∴ By principle of mathematical induction, we are done

2. a) This situation is certainly a possibility if the value function hasn't converged and extraction of the current policy yields a suboptimal policy that is stuck for the moment

Ex) Consider the following gridworld setup:



$i=3$:

$\vec{0}$	$\vec{0}$	$\vec{0}$	$\boxed{1}$
$\vec{0}$	$\vec{0}$	$\vec{0}$	$\boxed{-1}$
$\vec{0}$	$\vec{0}$	$\vec{0}$	$\downarrow 0$

Clearly $\pi_1 = \pi_2 \neq \pi_3$

b) From the given information,

$$|V^* - \tilde{V}| \leq \epsilon \Rightarrow \begin{cases} \tilde{V} \leq V^* + \epsilon \\ -\tilde{V} \leq -V^* + \epsilon \end{cases} \quad \text{--- ①}$$

So,

$$V_{\pi_{\tilde{V}}}(s) = R(s, \pi_{\tilde{V}}(s)) + \gamma \sum_{s'} P(s'|s, \pi_{\tilde{V}}(s)) V_{\pi_{\tilde{V}}}(s')$$

By definition of $\pi_{\tilde{V}}(s)$ (Bellman update of \tilde{V})

$$\begin{aligned} &= \left\{ \max_a \left[R(s, a) + \gamma \sum_{s'} P(s'|s, a) \cdot \tilde{V}(s') \right] \right. \\ &\quad \left. + \gamma \sum_{s'} P(s'|s, a) (V_{\pi_{\tilde{V}}}(s') - \tilde{V}(s')) \right\} \quad \text{--- ②} \end{aligned}$$

Substituting the terms from (1) into (2), we have

$$\begin{aligned}
 &= \max_a \left[R(s, a) + \gamma \sum_{s'} P(s'|s, a) (V^*(s') + \epsilon) \right] \\
 &\quad + \gamma \sum_{s'} P(s'|s, \pi_{\tilde{V}}(s)) (V_{\pi_{\tilde{V}}}(s') - V^*(s') + \epsilon) \\
 &= \max_a \left[R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s') \right] + \gamma \epsilon \sum_{s'} P(s'|s, a) \\
 &\quad + \gamma \sum_{s'} P(s'|s, \pi_{\tilde{V}}(s)) (V_{\pi_{\tilde{V}}}(s') - V^*(s')) + \gamma \epsilon \sum_{s'} P(s'|s, a)
 \end{aligned}$$

As V^* is optimal, $B V^* = V^*$ and $\sum_{s'} P(s'|s, a) = 1$ (\because probability distribution)

$$\begin{aligned}
 &= V^*(s) + 2\gamma\epsilon + \gamma \sum_{s'} P(s'|s, \pi_{\tilde{V}}(s)) (V_{\pi_{\tilde{V}}}(s') - V^*(s')) \\
 &\leq V^*(s) + 2\gamma\epsilon + \gamma (V_{\pi_{\tilde{V}}}(s) - V^*(s))
 \end{aligned}$$

This is because of the contraction property of the Bellman operator i.e.

$$B(V_{\pi_{\tilde{V}}}(s) - V^*(s)) \leq \gamma (V_{\pi_{\tilde{V}}}(s) - V^*(s))$$

Solving for $L_{\tilde{V}}(s)$ after re-arranging terms, we obtain

$$L_{\tilde{V}}(s) \leq \frac{2\gamma\epsilon}{1-\gamma}$$

3. a.) In order to ensure the optimal policy leads to shortest path, we define $R(s) = -1 \forall s \in S$ where $S = \{1, \dots, 24\} \setminus \{5, 11\}$. This is forced because the agent is penalized for every move taken to reach $s=5$.

b.) As $\gamma = 1$, $V_i(s) = E_{\pi} [R(s) + V(s')]$
backpropagating from $s=5$ yields

$$V(s) = \begin{cases} 5 & s = 5 \\ -5 & s = 11 \\ 4 & s = 0, 6, 10 \\ 3 & s = 1, 7, 15 \\ 2 & s = 2, 8, 12, 16, 20 \\ 1 & s = 4, 14, 18, 22 \\ 0 & s = 19, 23 \\ -1 & s = 3 \\ -2 & s = 24 \end{cases}$$

c.) This changes state values but not the optimal policy because discounting only adds further incentive to reduce number of moves required as rewards will still be negative $\forall s \neq 5$. Thus, the best strategy is to still take the shortest path.

d.) Now we have $R(s) = 0 \forall s \in S$. So ~~the~~ agent is no longer penalized for moving to any $s \neq 5$. However, as the discount factor is no longer 1, the agent is still incentivized to reach state 5 in the fewest number of moves. Hence, once again, shortest path is the optimal strategy to maximize reward. The value of each state will therefore be equal to that of the previous part for the same corresponding state plus an increment of

$$\sum_{i=0}^{L-1} \gamma^i$$

where $L =$ length of shortest path to state 5
 $\forall s \in S$

4

c)

Stochasticity in the second environment increases number of required iterations for convergence. The policies in both environments are identical.