HW4 SUNID: BEN MA

1a. $V_{opt}^{(t)}(s) \leftarrow$ MAX $\sum\limits_{s'} T(s,a,s')[$ Reward $(s,a,s') + r V_{opt}^{(t-1)}(s')]$
       over Actions(s)

Iter = 0,  $V_{opt} = \begin{bmatrix} 0, & 0, & 0, & 0, & 0 \end{bmatrix}$
                       $-2$  $-1$  $0$  $+1$  $+2$

Iter = 1,  $V_{opt}(S=-2) = 0$
$V_{opt}(S=-1, -1) = 0.8 \times (20+0) + 0.2 \times (-5+0) = 15$ } $\rightarrow = 15$
$V_{opt}(S=-1, +1) = 0.7 \times (20+0) + 0.3 \times (-5+0) = 12.5$ } max

$V_{opt}(S=0, -1) = 0.8 \times (-5 + 1\times 0) + 0.2 \times (-5 + 1\times 0) = -5$
$V_{opt}(S=0, +1) = 0.7 \times (-5) + 0.3 \times (-5) = -5$

$V_{opt}(S=1, -1) = 0.8 \times (-5+0) + 0.2 \times (100+0) = 16$ } $\rightarrow = 26.5$
$V_{opt}(S=1, +1) = 0.7 \times (-5) + 0.3 \times 100 = 26.5$ } max

$V_{opt}(S=+2) = 0$
$V_{opt}^{(1)}(s) = \begin{bmatrix} 0, & 15, & -5, & 26.5, & 0 \end{bmatrix}$
                    $-2$  $-1$  $0$  $+1$  $+2$

Iter = 2,  $V_{opt}(S=-2) = 0$
$V_{opt}(S=-1, -1) = \cancel{0.3 \times} \;\; 0.8 \times (20+0) + 0.2 \times (-5-5) = 14$
$V_{opt}(S=-1, +1) = 0.7 \times (20+0) + 0.3 \times (-5-5) = 11$ } $\rightarrow$ max $= 14$

$V_{opt}(S=0, -1) = 0.8 \times (-5 + 15) + 0.2 \times (-5 + 26.5 \times 1) = 12.3$ } $\rightarrow = 13.45$
$V_{opt}(S=0, +1) = 0.3 \times (-5 + 26.5 \times 1) + 0.7 \times (-5 + 15 \times 1) = 13.45$ } max

$V_{opt}(S=+1, -1) = 0.8 \times (-5 + (-5) \times 1) + 0.2 \times (100 + 0) = 12$ } $\rightarrow = 23$
$V_{opt}(S=+1, +1) = 0.7 \times (-5 + (-5) \times 1) + 0.3 \times (100+0) = 23$ } max
$V_{opt}(S=2) = 0$
$V_{opt}^{(2)}(s) = \begin{bmatrix} 0, & 14, & 13.45, & 23, & 0 \end{bmatrix}$
                    $-2$  $-1$  $0$  $+1$  $+2$

1 b. $\pi_{opt}(s) = \underset{a \in Action(s)}{arg\,max}\ Q_{opt}(s,a)$

So base on $V_{opt}^{(2)}(s)$, $\pi_{opt} = [\ -1, +1, +1\ ]$
$\qquad\qquad\qquad\qquad\qquad\quad s=-1\quad s=0\quad s=1$


2 b. Because we have an acyclic MDP, so it ~~will~~ will
not go back to a previous state, we can use recursion
to find the optimal value via one pass like below:
def recursive (s):
$\qquad V_{opt} = \underset{a \in Actions(s)}{Max}\ \sum_{s'} T(s,a,s')[\ Rewards(s,a,s') + recursive\ V(s)\ ]$


2 C. $V_{opt}^{(t)} = \underset{a \in Actions(s)}{MAX}\ \sum_{s'} T(s,a,s')[\ Reward(s,a,s') + \gamma V_{opt}^{(t+1)}(s')\ ]$

since $\gamma' = 1$. ~~Then we can define Reward(s,a,s')~~ $\frac{1}{T(s,a,s')}$

$V_{opt}' = \underset{a \in Actions(s)}{Max}\ \sum_{s'} T'(s,a,s')[\ Reward'(s,a,s') + V_{opt}^{(t+1)}(s')\ ]$

$\qquad = \underset{a \in Actions(s)}{Max}\ \sum_{s'} (T'\,R' + T'\,V_{opt}^{(t+1)})$

$V_{opt}' = \underset{a \in Actions(s)}{Max}\ \sum_{s'} (T\,R + \gamma T\,V_{opt}^{(t+1)})$

If $V_{opt}' = V_{opt}$ for all $s \in States$ Then $Reward'(s,a,s') = \cancel{\frac{1}{T}} Reward = 0$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad T'(s,a,s') = (1-\gamma)\ T(s,a,s')$

4b.

After running multiple value iterations between smallMDP and largeMDP, I find that largerMDP has larger discrepancies than the smallMDP. The reason is that the largeMDP has more state to explore than the smallMDP. The difference for the largeMDP is 1102 and smallMDP is 4. The percentage difference is 99.23%. Another issue is that the feature extractor extracted only state and action pairs.

4d.

The average reward from originalMDP is 12 and the average reward from new thresholdMDP is 10.

When value iteration find a policy of the originalMDP, it is to find the optimal policy according to the original threshold of 10 and if the threshold change to 15, it's still trying to avoid busting. The reason is that Q-Learning is offline policy which can find the new optimal policy to collect the maximum rewards but fixRLalgorithm is stateful, partial feedback, which cannot find optimal policy.

Therefore Q Learning is adaptable and can handle changes of thresholds.