# Ex 4.3

$$q_\pi(s,a) = E_\pi\left[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a\right]$$

$$= E_\pi\left[R_{t+1} + \gamma V_\pi(s') \mid S_t = s, A_t = a\right]$$

$$(4.3) \quad = \sum_{s',r} P(s',r \mid s,a)\left[r + \gamma V_\pi(s')\right]$$

$$(4.4) \quad = \sum_{s',r} P(s',r \mid s,a)\left[r + \gamma\left[\sum_{a'} \pi(a' \mid s') \cdot q_\pi(s',a')\right]\right]$$

$$q_{k+1}(s,a) = \sum_{s',r} P(s',r \mid s,a)\left[r + \gamma\left[\sum_{a'} \pi(a' \mid s') \cdot q_\pi(s',a')\right]\right]$$

# Ex 4.5

1. Initialization.

   $q(s,a) \in \mathbb{R}$ and $\pi(s) \in A(s)$ arbitrarily for all $s \in S$

2. Policy Evaluation

   Loop:

   $\Delta \leftarrow 0$

   Loop for each $s \in S$

   $q \leftarrow q(s, \pi(s))$

   $q(s,\pi(s)) \leftarrow \sum_{s',r} P(s',r \mid s,a)\left[r + \gamma\left[\sum_{a'} \pi(a' \mid s') \cdot q_\pi(s',a')\right]\right]$

   $\Delta \leftarrow \max\left(\Delta, \mid q - q(s, \pi(s))\mid\right)$

3. Policy Improvement

   Policy $-$ stable $\leftarrow$ true

   for each $s \in S$:

   old$-$action $\leftarrow \pi(s)$

   $\pi(s) \leftarrow \max\left(\sum_{s',r} P(s',r \mid s,a)\left[r + \gamma\left[\sum_{a'} \pi(a' \mid s') \cdot q_\pi(s',a')\right]\right.\right.$

   If old$-$action $\neq \pi(s)$, then policy$-$stable $\leftarrow$ false

   If policy$-$stable, then stop and return $q \approx q_*$ and $\pi \approx \pi_*$; else go to 2.

Question 2

Part 3 : The line in the DP method reach nearly 1.0 when state
is about 25, then keep steady in the rest of states. And
it is very smooth. However, the line in the Monte Carlo
method which is in episode 8000 increases very slowly compare
to the line in DP method, and it is not smooth. Because
the Monte Carle method of learning is depending on experience,
it trys the random policy. In gambler's problem, all
possibilities at states and action are exposed to agent,
and the possibilities are finite. So DP is more suitable