## 1.

Using the definition of Q function and Bellman equation we obtain the following:

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

$$= E_{\pi}\left[ R_{t+1} + \sum_{k=1}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right] \#;$$

$$\# \; E_{\pi}\left[ R_{t+1} \mid S_t = s, A_t = a \right] = E_{\pi} r(s,a);$$

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

$$= E_{\pi}\left[ \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_{t+1} = s', A_{t+1} = a' \right] = E_{\pi}\left[ \gamma Q_{\pi}(s',a') \right];$$

$$\Rightarrow Q^{\pi}(s,a) = E_{(s',a') \sim P(\cdot \mid s,a)}\left[ r(s,a) + \gamma Q^{\pi}(s',a') \right]$$

## 2.

$$Q^*(s,a) = \max_{\pi} q_{\pi}(s,a) = \max_{\pi}\left[ E_{(s',a') \sim P(\cdot \mid s,a)} \left[ r(s,a) + \gamma a^{\pi}(s',a') \right] \right]$$

$$= E_{(s',a') \sim \pi^*}\left[ r(s,a) + \gamma \max_{a'} Q^{\pi^*}(s',a') \right] =$$

$$= E_{(s',a') \sim \pi^*(\cdot \mid s,a)}\left[ r(s,a) + \gamma \max_{a'} Q^{\pi}(s',a') \right]$$

A plausable objective minimizes the difference between current $Q$-function and the optimal $Q$-function

$$\ell(\theta) = E \| Q^*(\xi, a, \theta) - Q(\xi, a, \theta) \|^2$$

$$Q^*(\xi, a, \theta) = E_{\xi' \sim \mathcal{T}^*(\cdot | \xi, a)} \left[ r(\xi, a) + \gamma \max_{a'} Q^*(\xi', a') \right] =$$

$$= E_{\xi' \sim \mathcal{T}^*(\cdot | \xi, a)} \left\{ r(\xi, a) + \gamma \max_{a'} \max Q(\xi', a') \right\}$$

$$\Rightarrow \ell(\theta) = E_{\xi' \sim \mathcal{T}^*(\cdot | \xi, a)} \| r + \gamma \max_{a'} \max Q(\xi', a', \theta) - Q(\xi, a, \theta) \|^2$$