1, (a)

① $V_*(s) = \max_\pi V^\pi(s) = \max_\pi \sum_{a \in A} \pi(a|s) Q^\pi(s,a) \leq \sum_{a \in A} \pi(a|s) Q_*(s,a)$

$\Rightarrow V_*(s) \leq \max_a Q_*(s,a)$ ( $\because$ deterministic optimal $\pi_*(a|s) = \begin{cases} 1, & a = \arg\max_{a \in A} Q_*(s,a) \\ 0, & \text{other} \end{cases}$ )

Suppose $V_*(s) < \max_a Q_*(s,a)$, then $\exists \pi'$ s.t. $\exists V^{\pi'}_*(s) = \sum_{a \in A} \pi'(a|s) Q_*(s,a) > V_*(s)$

$\Rightarrow V'_*(s) > V_*(s) \Rightarrow (\rightarrow\!\!\!\times) $ with $V_*(s) = \max_\pi V^\pi(s)$

Thus, $V_*(s) = \max_a Q_*(s,a)$ ✗

② $Q_*(s,a) = \max_\pi Q^\pi(s,a) = \max_\pi \left[ R_s^a + \gamma \sum_{s'} P_{ss'}^a V^\pi(s') \right]$

$= R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_\pi V^\pi(s') = R_s^a + \gamma \sum_{s'} P_{ss'}^a V_*(s')$ ✗

(b)

For any two action-value function $Q, Q'$,

$\| T^*(Q) - T^*(Q') \|_\infty = \max_{(s,a)} | T^*(Q)(s,a) - T^*(Q')(s,a) |$

$= \max_{(s,a)} | \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q(s',a') - \gamma \sum_{s'} P_{ss'}^a \max_{a'} Q'(s',a') |$

$\leq \max_{(s,a)} | \gamma \sum_{s'} P_{ss'}^a \max_{a'} [ Q(s',a') - Q'(s',a') ] |$

$\leq \max_{(s,a)} \max_{a'} | \gamma \sum_{s'} P_{ss'}^a ( Q(s',a') - Q'(s',a') ) |$

$\leq \gamma \| Q - Q' \|_\infty \Rightarrow T^*$ is a $\gamma$-contraction operator

2.

Let $\mu \in \mathbb{R}$, consider Lagrange function,

$$L(\pi) = \sum_{a \in A} \left( \pi(a|s) Q_\Omega^{\pi_k}(s,a) - \pi(a|s) \log \pi(a|s) \right) - \mu \left( \sum_{a \in A} \pi(a|s) - 1 \right)$$

$$\forall a \in A, \quad \frac{\partial L(\pi)}{\partial \pi(a|s)} = Q_\Omega^{\pi_k}(s,a) - \log \pi(a|s) - 1 - \mu \implies \pi(a|s) = \exp\left(Q_\Omega^{\pi_k}(s,a) - 1 - \mu\right)$$

$$\sum_{a \in A} \pi(a|s) = 1 \implies \sum_{a \in A} \exp\left(Q_\Omega^{\pi_k}(s,a) - 1 - \mu\right) = e^{-1-\mu} \sum_{a \in A} \exp\left(Q_\Omega^{\pi_k}(s,a)\right) = 1$$

$$\implies e^{1+\mu} = \sum_{a \in A} \exp\left(Q_\Omega^{\pi_k}(s,a)\right) \implies \mu = \ln \sum_{a \in A} \exp\left(Q_\Omega^{\pi_k}(s,a)\right) - 1$$

$$\pi(a|s) = \exp\left(Q_\Omega^{\pi_k}(s,a) - 1 - \left(\ln \sum_{a \in A} \exp\left(Q_\Omega^{\pi_k}(s,a)\right) - 1\right)\right)$$

$$= \frac{\exp\left(Q_\Omega^{\pi_k}(s,a)\right)}{\sum_{a \in A} \exp\left(Q_\Omega^{\pi_k}(s,a)\right)} \implies \text{optimal}$$

Thus, $\pi_{k+1}(\cdot|s) = \arg\max_\pi \left\{ \langle \pi(\cdot|s), Q_\Omega^{\pi_k}(s,a) \rangle - \Omega(\pi(\cdot|s)) \right\}$

$$= \frac{\exp\left(Q_\Omega^{\pi_k}(s,\cdot)\right)}{\sum_{a \in A} \exp\left(Q_\Omega^{\pi_k}(s,a)\right)}$$