Q1.

$$q_\lambda(y,b) \doteq E_\lambda\left[\sum_{k=t+1}^\infty (R_k - r(\lambda)) \mid S_t = y, A_t = b\right]$$

$$= E_\lambda\left[R_{k+1} - r(\lambda) + \sum_{k=t+2}^\infty (R_k - r(\lambda)) \mid S_t = y, A_t = b\right]$$

$$= E_\lambda\left[R_{k+1} - r(\lambda) + E_\lambda\left[\sum_{k=t+2}^\infty (R_k - r(\lambda)) \mid S_{t+1}, S_t = y\right]^{A_{t+1}, A_t = b} \mid S_t = y, A_t = b\right]$$

$$= E_\lambda\left[R_{k+1} - r(\lambda) + E_\lambda\left[\sum_{k=t+2}^\infty (R_k - r(\lambda)) \mid S_{t+1}, A_{t+1}\right] \mid S_t = y, A_t = b\right]$$

$$\text{(MP)}$$

$$= E_\lambda\left[R_{k+1} - r(\lambda) + q_\lambda(y_{t+1}, b_{t+1}) \mid S_t = y, A_t = b\right]$$

$$\text{($q_\lambda$ definition)}$$

$$= \sum_{b,y',r} P_\lambda(A_t = b, S_t = y', R_{t+1} = r \mid S_t = y)[r - r(\lambda) + q_\lambda(y', b)]$$

$$\text{(LOTUS)}$$

$$= \sum_{y',r} P(y', r \mid y, b)\left[r - r(\lambda) + \sum_b \lambda(b' \mid y') q_\lambda(y', b')\right]$$

Q2

For the transition S, A, R, S', where the action $A \sim b$ is drawn from policy distribution $b$, update the action value like $b = \pi$ (behavior policy = target policy)

following way:

$$Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma \sum_{a'} \pi(a'|S')Q(S',a') - Q(S,A)]$$

using greedification such as $\epsilon$-greedy.

Q3

1. Switch line 5 and line 6.

2. Add "Choose A' from s' using policy derived from Q" before line 8.

3. Change line 8 with
   "$Q(S,A) \leftarrow Q(S,A) + \lambda [Rt + r Q(S',A') - Q(S,A)]$"

4. Add "$A \leftarrow A'$" to line 9.

Q4

$$E_{A \sim b}\left[\rho(A|S)\left(R + \gamma V(S') - V(S)\right) \mid S=s\right]$$

$$= \sum_{a,s',r} P_b(A=a, S'=s', R=r \mid S=s)\; \rho(a|s)\left(r + \gamma V(s') - V(s)\right)$$

$$= \sum_{a,s',r} b(a|s)\, P(s',r \mid s, a)\; \frac{\pi(a|s)}{b(a|s)}\left(r + \gamma V(s') - V(s)\right)$$

$$= \sum_{a,s',r} P(s',r \mid s, a)\; \pi(a|s)\left(r + \gamma V(s') - V(s)\right)$$

$$= E_{A \sim \pi}\left[R + \gamma V(S') - V(S) \mid S=s\right]$$

On the other hand,

$$E_{A \sim b}\left[V(S) \mid S=s\right] = V(s)$$

Therefore , $E_{A \sim b}\left[V(S) \mid S=s\right] = S_{A \sim \pi}\left[V(S) | S=s\right]$.

Q5

$$\frac{\partial g(a)}{\partial a} = g(a)(1-g(a))$$

$$= \frac{1}{1+e^{-a}}\left(1- \frac{1}{1+e^{-a}}\right)$$

$$= \frac{1}{1+e^{-a}} \times \frac{e^{-a}}{1+e^{-a}} = \frac{e^{-a}}{(1+e^{-a})^2}$$

$$g(0) = \frac{1}{1+e^0} = \frac{1}{2}$$

$$\frac{\partial \hat{y}_k}{\partial B_{kj}} = X_j = g(\psi_j) = \frac{1}{1+e^{-\psi_j}}$$

$$\frac{\partial \hat{y}_k}{\partial A_{i,j}} = B_{k,i} \frac{\partial X_i}{\partial A_{i,j}} = B_{k,i} \frac{\partial g(\psi_i)}{\partial \psi_i} S_j$$

$$= B_{k,i}\left( g(\psi_i)(1-g(\psi_i))\right) S_j$$

$$= B_{k,i}\left(\frac{1}{1+e^{-\psi_i}} \times \frac{e^{-\psi_j}}{1+e^{-\psi_j}}\right) S_j$$

$$= B_{k,i} \cdot \frac{e^{-\psi_j}}{(1+e^{-\psi_j})^2} \cdot S_j$$

$$\psi_i = \sum_L A_{i,L} S_L = 0$$

$$\frac{\partial \hat{y}}{\partial B_{kj}} = g(\psi_i) = g(0) = \frac{1}{2}$$

And $B_{k,i} = 0$

Vicky Zhao

$$\frac{\partial \hat{y}_k}{\partial A_{ij}} = B_{k,i} \, g'(\psi_i) S_j = 0 \times \left( g(\psi_i) \times (1 - g(\psi_i)) \right) S_i$$

$$= 0 \times \left( \frac{1}{2} \times (1 - \frac{1}{2}) \right) S_j$$

$$= 0 \times \frac{1}{4} \times S_j = 0$$

Q6.

a)  Sarsa :

$$Q(S,A) \leftarrow Q(S,A) + \partial[R + \gamma Q(S',A') - Q(S,A)]$$

B: $Q(S,A) \leftarrow 1.5 + 0.2[0 + 0.9 \times 1.5 - 1.5]$

$$= 1.5 + 0.2[1.35 - 1.5]$$
$$= 1.5 + 0.2 \times (-0.15)$$
$$= 1.47.$$

D: $Q(S,A) \leftarrow 1.5 + 0.2[4 + 0.9 \times 0 - 1.5]$

$$= 1.5 + 0.2[4 + 0 - 1.5]$$
$$= 1.5 + 0.2 \times 2.5$$
$$= 1.5 + 0.5$$
$$= 2.$$

Vicky Zhao

b)

B: $Q(s,A) \leftarrow 1.47 + 0.2[0 + 0.9 \times 2 - 1.47]$

$= 1.47 + 0.2[1.8 - 1.47]$

$= 1.536$

D: $Q(s,A) \leftarrow 2 + 0.2[4 + 0.9 \times 0 - 2]$

$= 2 + 0.2[4 - 2]$

$= 2 + 0.2 \times 2$

$= 2 + 0.4$

$= 2.4$