

$A$  = Tensors,  $\otimes$  = Tensor Product,  $\times_k$  = mode- $k$  product,  $\odot$  = Hadamard product  
 $\pi$  = policy,  $(s, a)$  = state-action pair

$$Q(\pi|s, a) = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid s_t = s, a_t = a \right]$$

By Bellman equation we get,

$$Q(\pi|s, a) = R_{s, s'}^a + \gamma \sum_{a'} \pi(s', a') Q(\pi|s', a').$$

Expanding recursion we get,

$$Q(\pi|s_t, a_t) = R_{s_t, s_{t+1}}^{a_t} + \gamma \sum_{a_{t+1}} R_{s_{t+1}, s_{t+2}}^{a_{t+1}} \pi_{t+1} + \gamma^2 \sum_{a_{t+1}, a_{t+2}} R_{s_{t+2}, s_{t+3}}^{a_{t+2}} \pi_{t+1} \pi_{t+2} + \dots$$

we assume that a transition,  $s_t \xrightarrow{a_t} s_{t+1}$ , exists

Combining terms and considering  $C^{(k)} \in \mathbb{C}^{|S \times A|^k}$  is a sum of terms, we get

$$\sum_{s_t, a_t}^{S \times A} Q(\pi|s_t, a_t) = \left( \sum_{k=0}^{\infty} \sum_{\mu_{t+1}, \dots, \mu_{t+k}}^{S \times A} C^{(k)} \pi(\mu_{t+1}) \dots \pi(\mu_{t+k}) \right)$$

has the same form as a Hamiltonian. We can derive the optimal policy,  $\pi^*$ , by minimizing that function.

This approach aims to combat the issues surrounding linearity in Q-learning by representing the problem as a Hamiltonian & solving to find the optimal policy. This allows for the leveraging of quantum properties to consider multiple  $(s, a)$  pairs simultaneously, thereby deriving a Q-value faster than traditional Q-learning techniques.