

Optimal Control, MAP w/ known dynamics,

$$V^*(s) = \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma V^*(s'))$$

$$\pi^*(s) = \operatorname{argmax}_a \text{---}$$

① Value Iteration

$$V_k^*(s) = \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma V_{k-1}^*(s'))$$

$$\pi_k^*(s) = \operatorname{argmax}_a \text{---}$$

$$Q^*(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma \max_{a'} Q^*(s', a'))$$

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

② Q-Value Iteration

$$Q_k^*(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma \max_{a'} Q_{k-1}^*(s', a'))$$

$$\pi_k^*(s) = \operatorname{argmax}_a Q_k^*(s, a)$$

Policy evaluation for a given $\pi(s)$

$$V^\pi(s) = \sum_{s', r} p(s', r | s, \pi(s)) (r + \gamma V^\pi(s)) \rightarrow \text{solve system of } \beta \text{ linear equations}$$

OR

$$V_k^\pi(s) = \sum_{s', r} p(s', r | s, \pi(s)) (r + \gamma V_{k-1}^\pi(s)) \rightarrow \text{run until convergence.}$$

③ Policy Iteration

1. Policy evaluation for current policy π_j

2. Policy improvement (best action using one-step lookahead)

$$\pi_{j+1} = \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) (r + \gamma V_j^\pi(s))$$

repeat until convergence.