

HW14

1. LSTD Idea of using LSTD for approximate policy evaluation in PI

Start with random weights w (i.e. value function)

Repeat until Convergence

$\pi(s) = \text{greedy}(\Phi w)$

Evaluate π using LSTD

- Generate sample trajectories of P_π
- Use LSTD to produce new weights w (w gives an approximated value function of π)

Linear model approximation $\mathbf{W} = \left(I - \gamma (\Phi^T \Phi)^{-1} \Phi^T \Phi' \right)^{-1} (\Phi^T \Phi)^{-1} \Phi^T R = (\Phi^T \Phi - \gamma \Phi^T \Phi')^{-1} \Phi^T R$
aligns with LSTD solution $\mathbf{W} = \left(\Phi^T \Phi - \gamma \Phi^T \Phi' \right)^{-1} \Phi^T R$.

2. LSPI LSPI is similar to previous loop by replaces LSTD with a new algorithm LSTDQ