

$$\Theta = A^{-1} b$$

$$A = E [\phi_t (\phi_t - \gamma \phi_{t+1})^T]$$

$$b = E [r_{t+1} \phi_t]$$

$$\hat{A}_t = \sum_{k=0}^{t-1} \phi_k (\phi_k - \gamma \phi_{k+1})^T + \epsilon I$$

$$\hat{b}_t = \sum_{k=0}^{t-1} r_{k+1} \phi_k$$

$$\hat{\Theta}_t = \hat{A}_t^{-1} \hat{b}_t$$

$$\hat{A}_t^{-1} = \left(\hat{A}_{t-1} + \phi_t (\phi_t - \gamma \phi_{t+1})^T \right)^{-1}$$

Sherman - Morrison formula: v_{t-1}^T

$$\hat{A}_t^{-1} = \hat{A}_{t-1}^{-1} + \frac{\hat{A}_{t-1}^{-1} \phi_t^T (\phi_t - \gamma \phi_{t+1})^T \hat{A}_{t-1}^{-1}}{1 + (\phi_t - \gamma \phi_{t+1})^T \hat{A}_{t-1}^{-1} \phi_t}$$

$$v_{t-1} = (\hat{A}_{t-1}^{-1})^T (\phi_t - \gamma \phi_{t+1})$$

$$\hat{A}_t^{-1} = \hat{A}_{t-1}^{-1} - \frac{\hat{A}_{t-1}^{-1} \phi_t v_{t-1}^T}{1 + v_{t-1}^T \phi_t}$$

$$J(\omega) = \frac{V(s_0)}{\pi(\omega)} \quad \nabla_{\omega} J(\omega) = ?$$

$$\begin{aligned} \nabla_{\omega} J(s) &= \nabla_{\omega} \left[\sum_a \pi(a, s, \omega) Q_{\pi}(s, a) \right] = \\ &= \sum_a \left[\nabla_{\omega} \pi(a, s, \omega) Q_{\pi}(s, a) + \pi(a, s, \omega) \nabla_{\omega} Q_{\pi}(s, a) \right] = \\ &= \sum_a \left[-11- + \pi(a, s, \omega) \nabla_{\omega} \left[\sum_{s', r'} p(s', r' | s, a) (r' + \gamma V_{\pi}(s')) \right] \right] \\ &= \sum_a \left[-11- + \pi(a, s, \omega) \sum_{s', r'} p(s', r' | s, a) \gamma \nabla_{\omega} V_{\pi}(s') \right] \\ &= \sum_a \left[\nabla_{\omega} \pi(a, s, \omega) Q_{\pi}(s, a) + \pi(a, s, \omega) \sum_{s', r'} p(s', r' | s, a) \gamma \cdot \right. \\ &\quad \left. \left(\sum_{a'} \nabla_{\omega} \pi(a', s', \omega) Q_{\pi}(s', a') + \pi(a', s', \omega) \sum_{s'', r''} p(s'', r'' | s', a') \nabla_{\omega} V_{\pi}(s'') \right) \right] \\ &= \dots = \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \gamma^k p(s \rightarrow x, k, \pi) \sum_a \nabla_{\omega} \pi(a, x, \omega) Q_{\pi}(x, a) \end{aligned}$$

$$\nabla_{\omega} J(\omega) = \nabla_{\omega} V_{\pi}(s_0)$$

$$= \sum_s \nabla_{\omega} \pi(s) \sum_a \nabla_{\omega} \pi(a, s, \omega) Q_{\pi}(s, a)$$

$$\nabla_{\omega} J(\omega) = E_{\pi} \left[\gamma^t \sum_a \nabla_{\omega} \pi(a, s, \omega) Q_{\pi}(s, a) \right]$$

$$= E_{\pi} \left[\gamma^t \sum_a \pi(a, s, \omega) \left(\frac{\nabla_{\omega} \pi(a, s, \omega)}{\pi(a, s, \omega)} \right) Q_{\pi}(s, a) \right]$$

$$= E_{\pi} \left[\gamma^t \sum_a \pi(a, s, \omega) \nabla_{\omega} \log \pi(a, s, \omega) Q_{\pi}(s, a) \right]$$

$$= E_{\pi} \left[\gamma^t \nabla_{\omega} \log \pi(a, s, \omega) Q_{\pi}(s, a) \right]$$