

Init $\omega = 1$

$$V_{new}(S_1) = 0 + \delta \cdot V(S_2) = 1.8$$

$$V_{new}(S_2) = 0 + \delta \cdot V(S_2) = 1.8$$

Init $\omega = 1$.

LEAST SQUARES

$$V(S_1) = \omega \cdot 1$$

$$V(S_2) = \omega \cdot 2$$

$$ERROR = 1 \cdot (\omega \cdot 1 - 1.8)^2 + 9 \cdot (\omega \cdot 2 - 1.8)^2$$

$$2(\omega - 1.8) + 2 \cdot 9(2\omega - 1.8) = 0$$

$$\omega - 1.8 + 4\omega - 3.6 = 0$$

$$\omega = \frac{1.8 + 3.6}{5} = 1.08$$

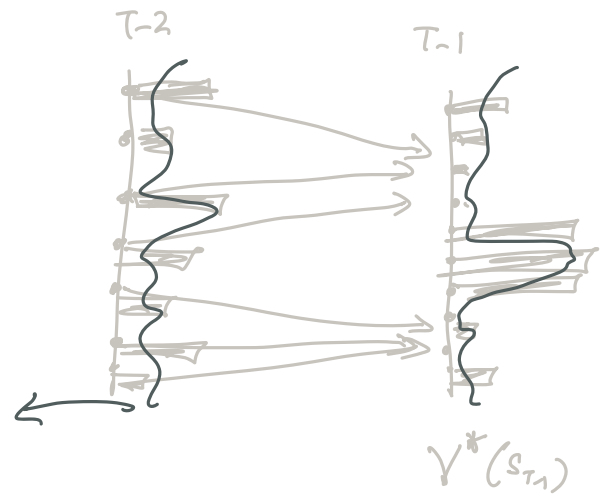
$$2(\omega - 1.8) + 2 \times 9(2\omega - 1.8) = 0$$

$$\omega - 1.8 + 18(2\omega - 1.8) = 0$$

$$\omega = \frac{1.8(1 + 18)}{37} \approx 0.92$$

DYNAMIC PROGRAMMING

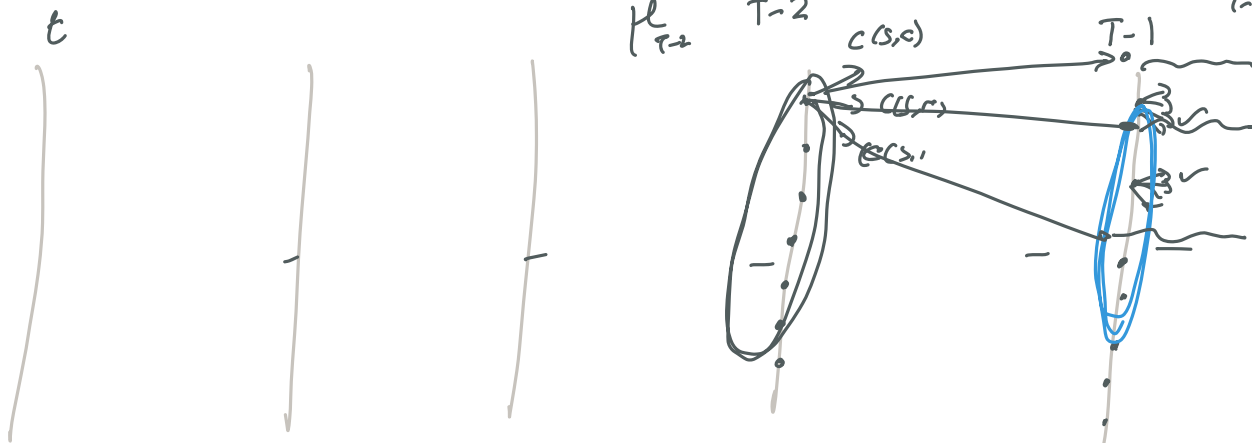
FITTED DYNAMIC PROGRAMMING



POLICY SEARCH VIA DYNAMIC PROGRAMMING

First

$$\pi_{T-1}^* \approx \underset{a}{\operatorname{argmax}} C(s, a)$$



$$\pi_{T-2}^*$$

$$= \underset{a_{t-2}}{\operatorname{argmax}} C(s_{t-2}, a_{t-2}) +$$

$$\cancel{V^* \pi_{T-1}^*(s_{t-1})}$$

Restart
 π_{T-1}^*

$$\pi_{T-1}^*(s) = \underset{a}{\operatorname{argmax}} C(s, a)$$

RESTART DISTRIBUTION $H_{T-1}^*(s)$