

Model-based (tabular) RL:

For each (s, a) , we have $\{(r_i, s'_i)\}_{i=1}^n$
 $\left. \begin{array}{l} R(s, a) \\ P(\cdot | s, a) \end{array} \right\}$

Build empirical model \hat{M} .

$$\hat{R}(s, a) = \frac{1}{n} \sum_{i=1}^n \hat{r}_i \approx R(s, a)$$

$$\hat{P}(s' | s, a) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[s'_i = s'] \approx P(\cdot | s, a)$$

Run Value-iteration in \hat{M} .

$$\underline{(\hat{T}f)}(s, a) = \underline{\hat{R}(s, a)} + \gamma \underline{\mathbb{E}_{s' \sim \hat{P}(\cdot | s, a)} [\max_{a'} f(s', a')]}$$

$$= \frac{1}{n} \sum_{i=1}^n r_i + \gamma \sum_{s' \in S} \hat{P}(s' | s, a) \cdot (\max_{a'} f(s', a'))$$

$$= \frac{1}{n} \sum_{i=1}^n r_i + \gamma \sum_{s' \in S} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{I}[s'_i = s'] \right) \cdot \max_{a'} f(s', a')$$

$$= \frac{1}{n} \left(\sum_{i=1}^n r_i + \gamma \sum_{s' \in S} \sum_{i=1}^n \mathbb{I}[s'_i = s'] \max_{a'} f(s', a') \right)$$

$$= \frac{1}{n} \left(\sum_{i=1}^n r_i + \gamma \sum_{i=1}^n \max_{a'} f(s'_i, a') \right)$$

$$= \frac{1}{n} \sum_{i=1}^n (r_i + \gamma \max_{a'} f(s'_i, a')) \approx (\hat{T}f)(s, a)$$

For fixed f . empirical Bellman update.

$$\begin{aligned} \mathbb{E}_{\substack{r_i \sim R(\cdot|s,a) \\ s_i \sim P(\cdot|s,a)}} \left[r_i + \gamma \max_{a'} f(s_i, a') \right] &= \underbrace{R(s,a)} + \gamma \mathbb{E}_{s_i \sim P(\cdot|s,a)} \left[\max_{a'} f(s_i, a') \right] \\ &= (Tf)(s,a). \end{aligned}$$

