

Approximate Value Iteration and Variants

Chen-Yu Wei

Value Iteration

For $k = 1, 2, \dots$

$$\forall s, a, \quad Q^{(k)}(s, a) \leftarrow \underbrace{R(s, a)}_{\text{unknown}} + \gamma \sum_{s'} \underbrace{P(s'|s, a)}_{\text{unknown}} \max_{a'} Q^{(k-1)}(s', a')$$

Idea: In each iteration, use multiple samples to estimate the right-hand side.

Least-Square Value Iteration (LSVI)

For $k = 1, 2, \dots$

 We want these samples to be “exploratory”

Obtain samples $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ where $\mathbb{E}[r_i] = R(s_i, a_i)$, $s'_i \sim P(\cdot | s_i, a_i)$

Find $Q^{(k)}$ such that

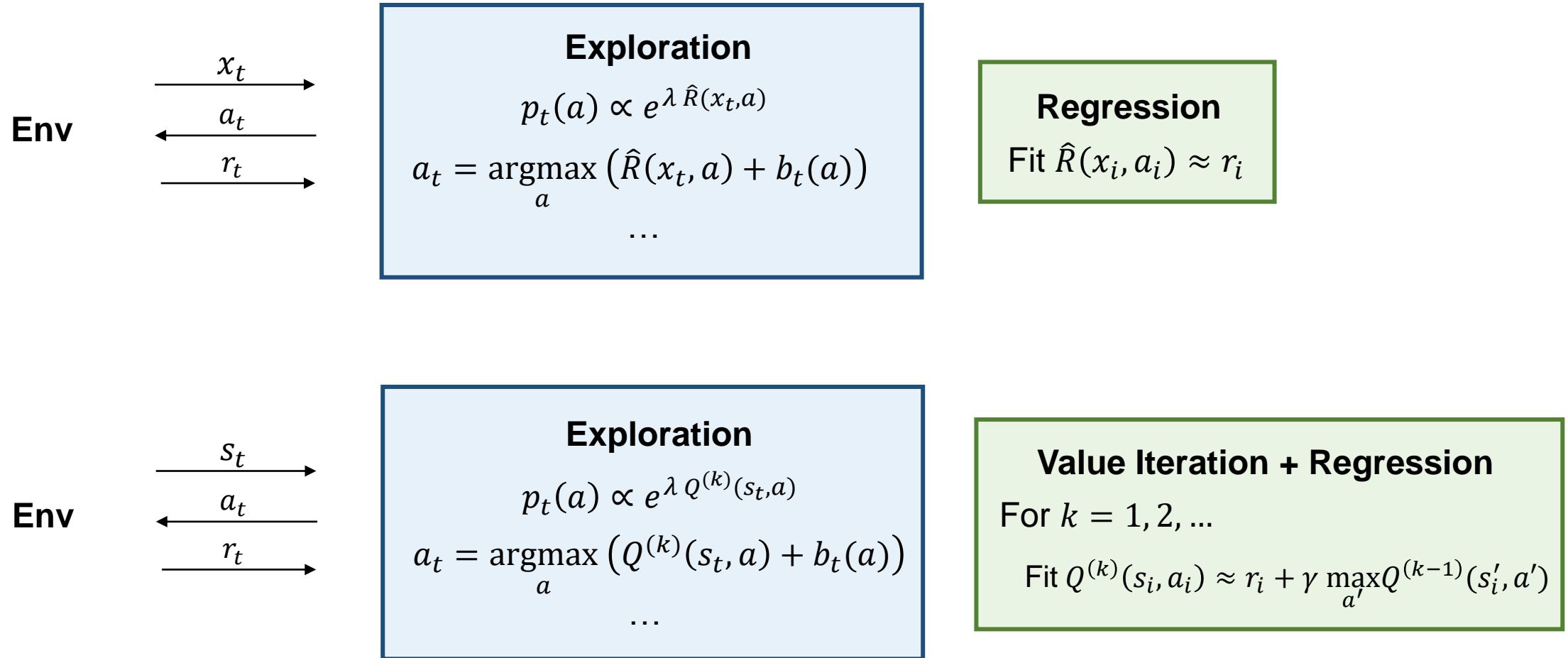
$$Q^{(k)}(s_i, a_i) \approx r_i + \gamma \max_{a'} Q^{(k-1)}(s'_i, a') \quad (\text{regression})$$

Tabular $\forall s, a, \quad Q^{(k)}(s, a) = \frac{\sum_{i=1}^n \mathbb{I}\{(s_i, a_i) = (s, a)\} \left(r_i + \gamma \max_{a'} Q^{(k-1)}(s'_i, a') \right)}{\sum_{i=1}^n \mathbb{I}\{(s_i, a_i) = (s, a)\}}$

General function approximation $\theta_k = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_{k-1}}(s'_i, a') \right)^2$

Linear function approximation $\theta_k = \left(\lambda I + \sum_{i=1}^n \phi(s_i, a_i) \phi(s_i, a_i)^{\top} \right)^{-1} \left(\sum_{i=1}^n \phi(s_i, a_i) \left(r_i + \gamma \max_{a'} \phi(s'_i, a')^{\top} \theta_{k-1} \right) \right)$

Comparison with Contextual Bandits



Analysis of LSVI under Certain Assumptions

To theoretically show that LSVI converges to the optimal value function, we will make some assumptions to ensure the following holds for all iteration k :

$$Q^{(k)}(s, a) \approx R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a'} Q^{(k-1)}(s', a') \right]$$

Linear case:

$$\phi(s, a)^\top \theta^{(k)} \approx R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a'} \phi(s', a')^\top \theta^{(k-1)} \right]$$

Analysis of LSVI under Certain Assumptions

1. Approximate Bellman Completeness Assumption: For any $\theta \in \mathbb{R}^d$, there exists a $\theta' \in \mathbb{R}^d$ such that

$$\left| \phi(s, a)^\top \theta' - \left(R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a'} \phi(s', a')^\top \theta \right] \right) \right| \leq \epsilon_{\text{fa}}$$

This ensures that no matter what θ_{k-1} is, there always exists a θ_k^* such that

$$\phi(s, a)^\top \theta_k^* \approx R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a'} \phi(s', a')^\top \theta_{k-1} \right]$$

This is similar to the linear assumption $|\phi(s, a)^\top \theta^* - R(s, a)| \leq \epsilon_{\text{fa}}$ in contextual bandits, but is qualitatively stronger because the assumption require “for any θ ”.

Analysis of LSVI under Certain Assumptions

2. Coverage Assumption: The dataset $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ collected in each iteration allows us to find θ_k so that for any s, a ,

$$|\phi(s, a)^\top \theta_k - \phi(s, a)^\top \theta_k^*| \leq \epsilon_{\text{stat}}$$

Recall from linear contextual bandit analysis, with

$$\theta_k = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n \left(\phi_i^\top \theta - \underbrace{\left(r_i + \gamma \max_{a'} \phi(s'_i, a')^\top \theta_{k-1} \right)}_{\text{Expectation} = \phi_i^\top \theta_k^*} \right)^2 + \lambda \|\theta\|^2$$

we have $|\phi(s, a)^\top (\theta_k - \theta_k^*)| \lesssim \sqrt{\beta} \|\phi(s, a)\|_{\Lambda^{-1}}$ where $\Lambda = \lambda I + \sum_{i=1}^n \phi_i \phi_i^\top$

In linear CB, we did not make such an assumption. What we did there is adding $\sqrt{\beta} \|\phi(s, a)\|_{\Lambda^{-1}}$ as **exploration bonus**, which aims to make $\sqrt{\beta} \|\phi(s, a)\|_{\Lambda^{-1}}$ small for all s, a .

Analysis of LSVI under Certain Assumptions

Under approximate Bellman completeness and coverage assumptions, LSVI ensures

$$\|Q^{(k)} - Q^*\|_{\infty} \leq O\left(\gamma^k \|Q^{(0)} - Q^*\|_{\infty} + \frac{\epsilon_{\text{fa}} + \epsilon_{\text{stat}}}{1 - \gamma}\right)$$

where $\|Q^{(k)} - Q^*\|_{\infty} := \max_{s,a} |Q^{(k)}(s, a) - Q^*(s, a)|$

Also, the greedy policy $\pi^{(k)}(s) = \operatorname{argmax}_a Q^{(k)}(s, a)$ satisfies for all s ,

$$V^*(s) - V^{\pi^{(k)}}(s) \leq O\left(\gamma^k \|Q^{(0)} - Q^*\|_{\infty} + \frac{\epsilon_{\text{fa}} + \epsilon_{\text{stat}}}{1 - \gamma}\right)$$