

Overview

- Planning with a simulator: estimate $V(s)$ with accuracy ε .
- Sample complexity = # calls to the simulator to get accuracy ε .
- Non-regularized case: there are no known algorithms with polynomial guarantees in a general setting.
- Regularization \rightarrow smooth Bellman operator.
- **SmoothCruiser**: exploits smoothness \rightarrow sample complexity that is always **polynomial**.

Setting & Assumptions

- MDPs and two-player games: $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$.
- \mathcal{S}, \mathcal{A} : state and action spaces.
- $P \triangleq \{P(\cdot|s, a)\}_{s, a \in \mathcal{S} \times \mathcal{A}}$ transition probabilities.
- $R: \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ reward function.
- $\gamma \in [0, 1]$ discount factor.
- **Assumption 1**: $|\mathcal{S}|$ arbitrary, $|\mathcal{A}| = K < \infty$.
- **Assumption 2**: we have a generative model (oracle, simulator) that can sample rewards and transitions: $R, Z \leftarrow \text{oracle}(s, a)$.

Regularized Value Functions

- Given a function $F_s: \mathbb{R}^K \rightarrow \mathbb{R}$, we define the value function as

$$V(s) \triangleq F_s(Q_s),$$

$$Q_s(a) \triangleq \mathbb{E}_{z \sim P(\cdot|s, a)}[R(s, a) + \gamma V(z)].$$

- In MDPs: $F_s = \max$ gives Bellman equations.
- In games: $F_s = \max$ for player 1 and $F_s = \min$ for player 2.
- Let $\text{LogSumExp}_\lambda(x) \triangleq \lambda \log \sum_{i=1}^K \exp(x_i/\lambda)$.
- **Entropy regularization** with parameter λ :

$$\max \text{ becomes } \text{LogSumExp}_\lambda$$

$$\min \text{ becomes } -\text{LogSumExp}_{-\lambda}.$$

- Property: with regularization, F_s becomes L -smooth

$$|F_s(Q) - F_s(Q') - (Q - Q')^\top \nabla F_s(Q')| \leq L \|Q - Q'\|_2^2$$

with $L = 1/\lambda$, $\nabla F_s(Q) \succeq 0$ and $\|\nabla F_s(Q)\|_1 = 1$.

Motivation

- **Strong regularization** $\lambda \rightarrow \infty$ and $L = 0$, that is, F_s is linear for all s : $F_s(x) = w_s^\top x$, with $\|w_s\|_1 = 1$, $w_s \in \mathbb{R}^k$ and $w_s \succeq 0$.
 \rightarrow Monte Carlo sampling, $\mathcal{O}(1/\varepsilon^2)$ calls to the oracle.
- **No regularization** $\lambda = 0$ and $L \rightarrow \infty$, that is, F_s cannot be well approximated by a linear function.
 \rightarrow Sparse Sampling algorithm (Kearns et al., 1999), $\mathcal{O}((1/\varepsilon)^{\log(1/\varepsilon)})$ calls to the oracle.
- For $0 < \lambda < \infty$: interpolate between the two extreme cases using linear approximations of F_s .

Theoretical Guarantees

Theorem 1. Let $n(\varepsilon, \delta')$ be the number of calls to the generative model before the algorithm terminates. For any state $s \in \mathcal{S}$ and $\varepsilon, \delta' > 0$,

$$n(\varepsilon, \delta') = \tilde{\mathcal{O}}(1/\varepsilon^4).$$

Theorem 2. For any $s \in \mathcal{S}$, $\varepsilon > 0$ and $\delta > 0$, there exists a choice of δ' that depends on ε and δ such that the output $\hat{V}(s)$ of **SmoothCruiser**(s, ε, δ') satisfies

$$\mathbb{P}[|\hat{V}(s) - V(s)| > \varepsilon] \leq \delta.$$

and such that $n(\varepsilon, \delta') = \mathcal{O}(1/\varepsilon^{4+c})$ for any $c > 0$.

Intuition

- Linear approximation of F_s around $\hat{Q}_s = \text{estimateQ}(s, \sqrt{\varepsilon/L})$:
$$F_s(Q_s) \approx F_s(\hat{Q}_s) + (Q_s - \hat{Q}_s)^\top \nabla F_s(\hat{Q}_s) + \varepsilon$$
- One of the terms can be written as an expected value:
$$Q_s^\top \nabla F_s(\hat{Q}_s) = \mathbb{E}[Q_s(A) | \hat{Q}_s], \text{ with } A \sim \nabla F_s(\hat{Q}_s)$$
- **Problem**: Q_s is unknown.
- **Solution**: use $\tilde{Q}(A) = R_{s,A} + \gamma \text{sampleV}(Z_{s,A}, \varepsilon/\sqrt{\gamma})$.
- Finally, we have (in expectation):
$$F_s(Q_s) \approx F_s(\hat{Q}_s) - \hat{Q}_s^\top \nabla F_s(\hat{Q}_s) + \tilde{Q}(A) + \varepsilon.$$

SmoothCruiser

Algorithm 1 SmoothCruiser

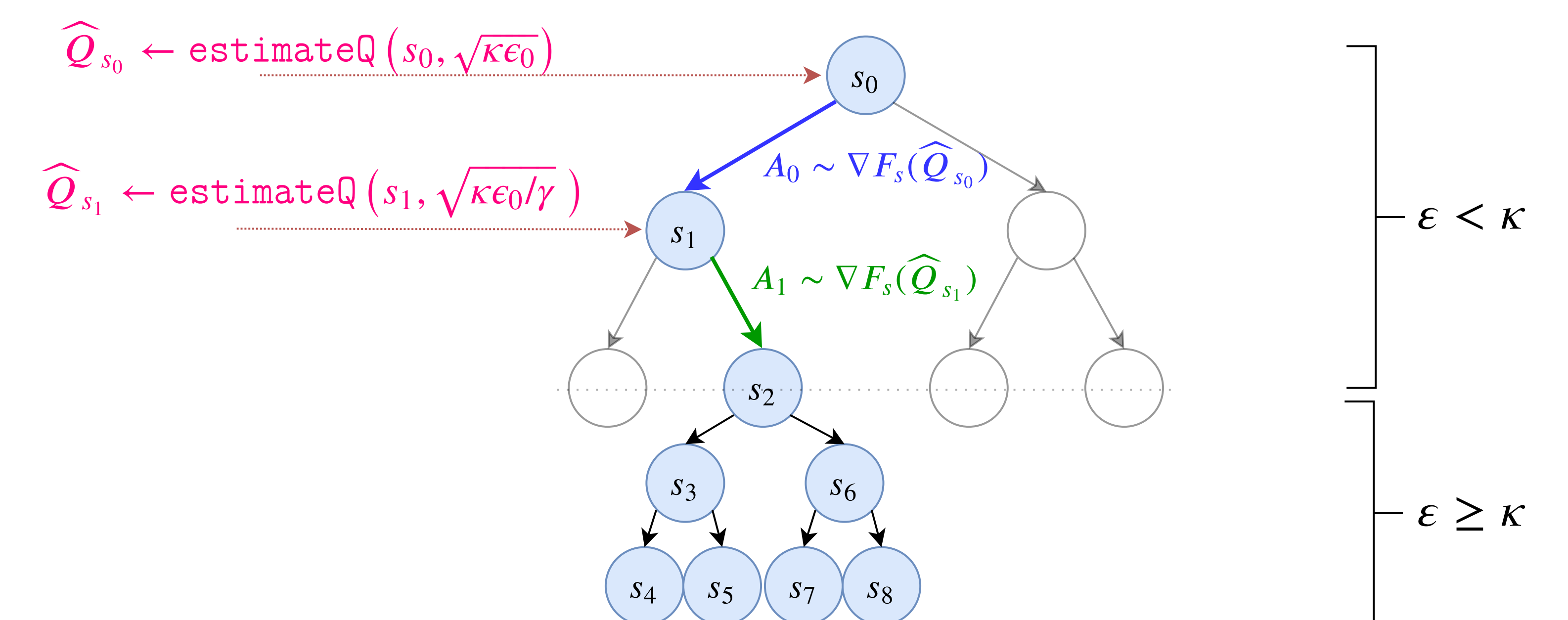
Input: $(s, \varepsilon, \delta') \in \mathcal{S} \times \mathbb{R}_+ \times \mathbb{R}_+$
 $M_\lambda \leftarrow \sup_{s \in \mathcal{S}} |F_s(0)| = \lambda \log K$
 $\kappa \leftarrow (1 - \sqrt{\gamma})/(KL)$
 Set δ', κ and M_λ as a global parameters
 $\hat{Q}_s \leftarrow \text{estimateQ}(s, \varepsilon)$
Output: $F_s(\hat{Q}_s)$

Algorithm 2 estimateQ

1: **Input:** (s, ε)
 2: $N(\varepsilon) \leftarrow \mathcal{O}(\log(2K/\delta')/\varepsilon^2)$
 3: **for** $a \in \mathcal{A}$ **do**
 4: $\hat{Q}_s(a) \leftarrow 0$
 5: **for** $i \in 1, \dots, N(\varepsilon)$ **do**
 6: $(R_i, Z_i) \leftarrow \text{oracle}(s, a)$.
 7: $\hat{V}_i \leftarrow \text{sampleV}(Z_i, \varepsilon/\sqrt{\gamma})$
 8: **end for**
 9: $\hat{Q}_s(a) = \text{average of } \{R_i + \gamma \hat{V}_i\}_{i=1}^{N(\varepsilon)}$
 10: **end for**
 11: **Output:** \hat{Q}_s clipped to $[0, Q_{\max}]$

Algorithm 3 sampleV

1: **Input:** $(s, \varepsilon) \in \mathcal{S} \times \mathbb{R}_+$
 2: **if** $\varepsilon \geq (1 + M_\lambda)/(1 - \gamma)$ **then**
 3: **Output:** 0
 4: **else if** $\varepsilon \geq \kappa$ **then**
 5: $\hat{Q}_s \leftarrow \text{estimateQ}(s, \varepsilon)$
 6: **Output:** $F_s(\hat{Q}_s)$
 7: **else if** $\varepsilon < \kappa$ **then**
 8: $\hat{Q}_s \leftarrow \text{estimateQ}(s, \sqrt{\kappa\varepsilon})$
 9: $A \leftarrow \text{action drawn from } \nabla F_s(\hat{Q}_s)$
 10: $(R, Z) \leftarrow \text{oracle}(s, A)$
 11: $\hat{V} \leftarrow \text{sampleV}(Z, \varepsilon/\sqrt{\gamma})$
 12: **Output:** $F_s(\hat{Q}_s) - \hat{Q}_s^\top \nabla F_s(\hat{Q}_s) + (R + \gamma \hat{V})$
 13: **end if**



Sample Complexity Simulation

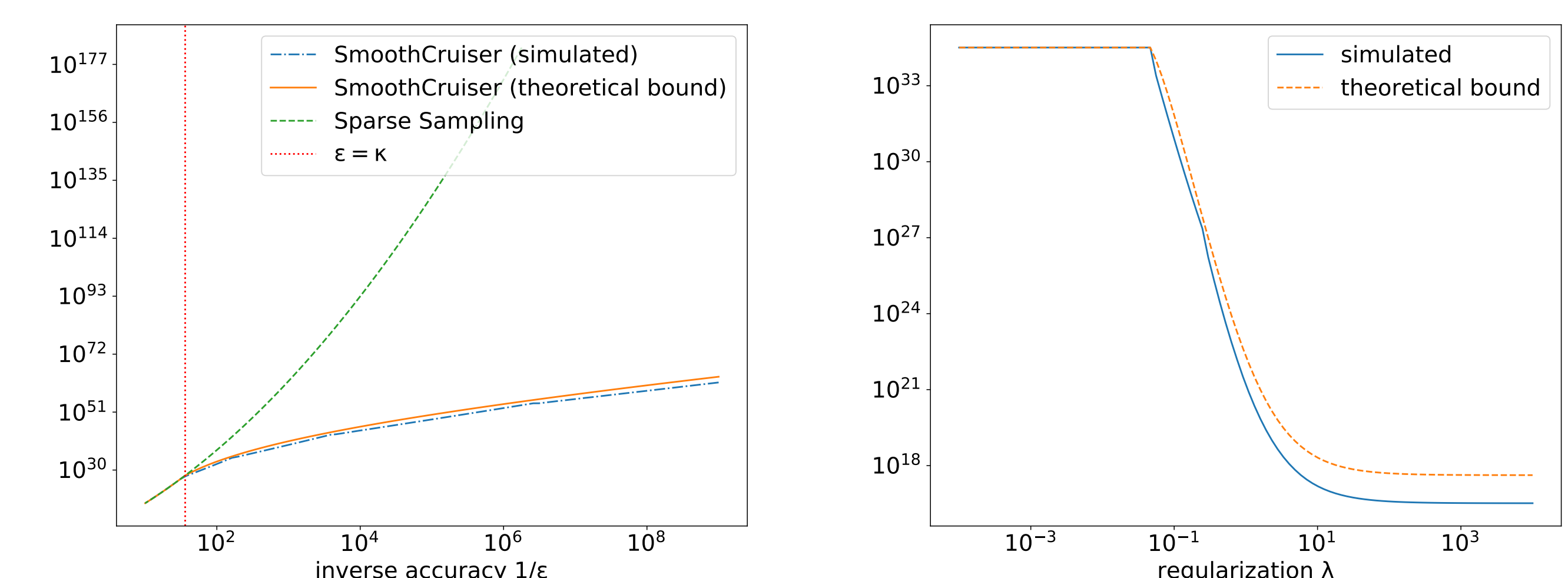


Figure 1: Left: number of calls to **sampleV** as a function of $1/\varepsilon$. Right: number of calls to **sampleV** required to achieve a 0.01 relative error as a function of the regularization λ .