Oracle Complexity Reduction for Model-free LQR: A Stochastic Variance-Reduced Policy Gradient Approach

Leonardo F. Toso, Han Wang, James Anderson

Department of Electrical Engineering, Columbia University

February 28, 2025



American Control Conference (ACC), 2024









Cost queries are expensive

What is a cost query?

Precision Robotic Arm



Control: Joint angle, speed

Cost metric: Positioning error of the silicon wafer

Cost query: Configure the arm and measure the positioning error associated with a specific set of system parameters.

Cost queries are expensive



- Expensive
- Time consuming
- May incur a high risk for human being interaction

Cost queries are expensive



We need to reduce the oracle complexity for large scale optimal control via policy gradient methods

$$\min_{x\in\mathcal{X}}f(x)=\frac{1}{n}\sum_{i=1}^{n}f_{i}(x)$$

where $f_i(x)$ are convex functions.

$$\min_{x\in\mathcal{X}}f(x)=\frac{1}{n}\sum_{i=1}^{n}f_{i}(x)$$

where $f_i(x)$ are convex functions.

Gradient Descent (GD): $x_{k+1} = x_k - \eta \nabla_x f(x_k)$

Stochastic GD: $x_{k+1} = x_k - \eta \nabla_x f_z(x_k)$, where $z \sim \{1, 2, \dots, n\}$

$$\min_{x\in\mathcal{X}}f(x)=\frac{1}{n}\sum_{i=1}^{n}f_{i}(x)$$

where $f_i(x)$ are convex functions.

Gradient Descent (GD): $x_{k+1} = x_k - \eta \nabla_x f(x_k)$

Stochastic GD: $x_{k+1} = x_k - \eta \underbrace{\nabla_x f_z(x_k)}_{\approx \nabla_x f(x_k)?}$, where $z \sim \{1, 2, \dots, n\}$

$$\min_{x\in\mathcal{X}}f(x)=\frac{1}{n}\sum_{i=1}^{n}f_{i}(x)$$

where $f_i(x)$ are convex functions.

Gradient Descent (GD): $x_{k+1} = x_k - \eta \nabla_x f(x_k)$

Stochastic GD: $x_{k+1} = x_k - \eta \underbrace{\nabla_x f_z(x_k)}_{\approx \nabla_x f(x_k)?}$, where $z \sim \{1, 2, \dots, n\}$

Stochastic Variance-Reduced:

$$x_{k+1} = x_k - \eta (\nabla_x f_z(x_k) \underbrace{-\nabla_x f_z(y) + \nabla_x f(y)}_{\text{zero mean}})$$

where y is generic point.

Johnson and Zhang, 2013, Reddi et al., 2016.

Stochastic Variance-Reduced:

$$x_{k+1} = x_k - \eta \underbrace{\left(\nabla_x f_z(x_k) - \nabla_x f_z(y) + \nabla_x f(y) \right)}_{\text{control variates } u}$$

Mean: $\mathbb{E}(v_k) = \mathbb{E}(\nabla_x f_z(x_k))$

Variance: $X = \nabla_x f_z(x_k), \ Y = -\nabla_x f_z(y) + \nabla_x f(y)$

$$\mathsf{var}(v_k) = \mathsf{var}(X) + \mathsf{var}(Y) - 2\mathsf{cov}(X, Y)$$

Johnson and Zhang, 2013, Reddi et al., 2016.

Stochastic Variance-Reduced:

$$x_{k+1} = x_k - \eta \underbrace{\left(\nabla_x f_z(x_k) - \nabla_x f_z(y) + \nabla_x f(y)\right)}_{\text{control variates } u}$$

Mean: $\mathbb{E}(v_k) = \mathbb{E}(\nabla_x f_z(x_k))$

Variance: $X = \nabla_x f_z(x_k), Y = -\nabla_x f_z(y) + \nabla_x f(y)$

$$var(v_k) = var(X) + var(Y) - 2cov(X, Y)$$

Question: Can we design an oracle-efficient solution to address the model-free LQR problem by building upon the success of stochastic variance-reduced approaches?

Linear Quadratic Regulator (LQR)

Consider a discrete LTI dynamical system

$$x_{t+1} = Ax_t + Bu_t, \quad t = 0, 1, 2, \dots$$
 (sys-dyn)

where $x_t \in \mathbb{R}^{d_x}$ and $u_t \in \mathbb{R}^{d_u}$.

Linear Quadratic Regulator (LQR)

Consider a discrete LTI dynamical system

$$x_{t+1} = Ax_t + Bu_t, \quad t = 0, 1, 2, \dots$$
 (sys-dyn)

where $x_t \in \mathbb{R}^{d_x}$ and $u_t \in \mathbb{R}^{d_u}$.

LQR Objective: Design a controller K^* ($u_t = -K^* x_t$) that solves

$$\mathcal{K}^{\star} = \operatorname{argmin}_{\mathcal{K} \in \mathcal{K}} \mathcal{C}(\mathcal{K}) = \mathbb{E}\left[\sum_{t=0}^{\infty} x_t \left(Q + \mathcal{K}^{\top} \mathcal{R} \mathcal{K}\right) x_t\right],$$

subject to (sys-dyn).

Stabilizing set: $\mathcal{K} = \{ \mathcal{K} \mid \rho(\mathcal{A} - \mathcal{B}\mathcal{K}) < 1 \}.$

Linear Quadratic Regulator (LQR) Formulation

LQR Objective: Design a controller K^* ($u_t = -K^* x_t$) that solves

$$\mathcal{K}^{\star} = \operatorname{argmin}_{\mathcal{K} \in \mathcal{K}} \mathcal{C}(\mathcal{K}) = \mathbb{E}\left[\sum_{t=0}^{\infty} x_t \left(Q + \mathcal{K}^{\top} \mathcal{R} \mathcal{K}\right) x_t\right],$$

subject to (sys-dyn).

Model-based LQR: Given (A, B, Q, R),

 $K^* = \text{DARE}(A, B, Q, R) \rightarrow \text{Riccati Equation}$

Linear Quadratic Regulator (LQR) Formulation

LQR Objective: Design a controller K^* ($u_t = -K^* x_t$) that solves

$$\mathcal{K}^{\star} = \operatorname{argmin}_{\mathcal{K} \in \mathcal{K}} \mathcal{C}(\mathcal{K}) = \mathbb{E}\left[\sum_{t=0}^{\infty} x_t \left(Q + \mathcal{K}^{\top} \mathcal{R} \mathcal{K}\right) x_t\right],$$

subject to (sys-dyn).

Model-based LQR: Given (A, B, Q, R),

 $K^{\star} = \text{DARE}(A, B, Q, R) \rightarrow \text{Riccati Equation}$

Q: How to design K^* when (A, B, Q, R) is unknown?

Linear Quadratic Regulator (LQR)

Policy Gradient LQR

Fazel et al., ICML 2018, proved that despite of the non-convexity^{*} of C(K), PG methods globally converge to K^* , i.e., given

Initial Stabilizing Controller: $\textit{K}_{0} \in \mathcal{K}$

Controllability: (A, B) is controllable

Linear Quadratic Regulator (LQR) Policy Gradient LQR

Fazel et al., ICML 2018, proved that despite of the non-convexity^{*} of C(K), PG methods globally converge to K^* , i.e., given

Initial Stabilizing Controller: $\textit{K}_{0} \in \mathcal{K}$

Controllability: (A, B) is controllable

$$K_{n+1} = K_n - \eta \widehat{\nabla} C(K_n), \text{ for } n \in \{0, 1, \dots, N-1\}$$

$$\downarrow$$

$$C(K_N) - C(K^*) \leq \epsilon,$$

after $N \geq \mathcal{O}(\log(1/\epsilon))$.

Gradient Dominance: $C(K) - C(K^*) \leq \lambda \|\nabla C(K)\|_F^2$ for any $K \in \mathcal{K}$.

* Non-convex for $d_X \geq 3$.

Linear Quadratic Regulator (LQR)

Policy Gradient LQR

Other nice properties of the LQR cost:

- Uniform bounds for the gradient
- Lipschitz of the cost
- Lipschitz of the gradient
- Bounded controller difference

Linear Quadratic Regulator (LQR)

Policy Gradient LQR

Other nice properties of the LQR cost:

- Uniform bounds for the gradient
- Lipschitz of the cost
- Lipschitz of the gradient
- Bounded controller difference

Stabilizing sub-level set: Given (sys-dyn), the stabilizing sub-level set $\mathcal{G}\subseteq \mathcal{K}$ is

$$\mathcal{G} = \left\{ \mathcal{K} \mid \mathcal{C}(\mathcal{K}) - \mathcal{C}(\mathcal{K}^{\star}) \leq \gamma(\mathcal{C}(\mathcal{K}_0) - \mathcal{C}(\mathcal{K}^{\star})) \right\},$$

for some $\gamma > 0$.

Fazel et al. 2018, Bu et al., 2019, Gravell et al. 2020.

Zeroth-order Gradient Estimation

$$K_{n+1} = K_n - \eta \widehat{\nabla} C(K_n), \text{ for } n \in \{0, 1, \dots, N-1\}$$

Zeroth-order Gradient Estimation

$$K_{n+1} = K_n - \eta \widehat{\nabla} C(K_n)$$
, for $n \in \{0, 1, \dots, N-1\}$

One-point Zeroth-order Estimation (Z01P):

$$(m,r) \rightarrow \text{ZO1P}: \overline{\nabla}C(K) := \sum_{i=1}^{m} \frac{d_{x}d_{u}C(K+U_{i})U_{i}}{mr^{2}},$$

m (number of trajectories), *r* (smoothing radius) and $||U_i||_F = r$.

Zeroth-order Gradient Estimation

$$K_{n+1} = K_n - \eta \widehat{\nabla} C(K_n)$$
, for $n \in \{0, 1, \dots, N-1\}$

One-point Zeroth-order Estimation (Z01P):

$$(m,r) \rightarrow \text{ZO1P}: \overline{\nabla}C(K) := \sum_{i=1}^{m} \frac{d_{x}d_{u}C(K+U_{i})U_{i}}{mr^{2}},$$

m (number of trajectories), *r* (smoothing radius) and $||U_i||_F = r$.

Zeroth-order Gradient Estimation

$$K_{n+1} = K_n - \eta \widehat{\nabla} C(K_n)$$
, for $n \in \{0, 1, \dots, N-1\}$

One-point Zeroth-order Estimation (Z01P):

$$(m,r) \rightarrow \text{ZO1P}: \overline{\nabla}C(K) := \sum_{i=1}^{m} \frac{d_{X}d_{u}C(K+U_{i})U_{i}}{mr^{2}},$$

m (number of trajectories), *r* (smoothing radius) and $||U_i||_F = r$.

Two-points Zeroth-order Estimation (Z02P):

$$(m,r) \rightarrow \text{ZO2P}: \widetilde{\nabla}C(K) := \sum_{i=1}^{m} \frac{d_{x}d_{u}(C(K+U_{i})-C(K-U_{i}))U_{i}}{2mr^{2}},$$

Illustrative Example - ZO Bias and Variance

$$A = \begin{bmatrix} 1.20 & 0.50 & 0.40 \\ 0.01 & 0.75 & 0.30 \\ 0.10 & 0.02 & 1.50 \end{bmatrix}, B = \begin{bmatrix} \frac{1}{2} \\ 1 \\ \frac{1}{2} \end{bmatrix}, Q = 2I_3, R = \frac{1}{2},$$

Improvements on the Oracle Complexity

Comparison on the sample complexity (\mathbb{S}_c), and two-point oracle complexity ($\mathcal{N}_{\text{ZO2P}}$) required to achieve $\mathbb{E}\left(C(K_N) - C(K^*)\right) \leq \epsilon$.

Methods	\mathbb{S}_{c}	$\mathcal{N}_{\mathrm{ZO2P}}$
PG - ZO1P (Fazel et al (2018))	$\mathcal{O}(1/\epsilon^4 \cdot \log{(1/\epsilon)})$	-
PG - ZO1P (Gravell et al (2019)	$\mathcal{O}(1/\epsilon^4 \cdot \log{(1/\epsilon)})$	-
PG - ZO1P (Malik et al. (2019)	$\mathcal{O}(1/\epsilon^2 \cdot \log{(1/\epsilon)})$	-
PG - ZO2P (Malik et al. (2019)	$\mathcal{O}(1/\epsilon \cdot \log{(1/\epsilon)})$	$\mathcal{O}(1/\epsilon \cdot \log{(1/\epsilon)})$
PG - ZO2P (Mohammadi et al. (2020))	$\mathcal{O}(\log{(1/\epsilon)})$	$\mathcal{O}(\log{(1/\epsilon)})$

Improvements on the Oracle Complexity

Comparison on the sample complexity (\mathbb{S}_c), and two-point oracle complexity ($\mathcal{N}_{\text{ZO2P}}$) required to achieve $\mathbb{E}\left(C(K_N) - C(K^*)\right) \leq \epsilon$.

Methods	\mathbb{S}_{c}	$\mathcal{N}_{\mathrm{ZO2P}}$
PG - ZO1P (Fazel et al (2018))	$\mathcal{O}(1/\epsilon^4 \cdot \log{(1/\epsilon)})$	-
PG - ZO1P (Gravell et al (2019)	$\mathcal{O}(1/\epsilon^4 \cdot \log{(1/\epsilon)})$	-
PG - ZO1P (Malik et al. (2019)	$\mathcal{O}(1/\epsilon^2 \cdot \log{(1/\epsilon)})$	-
PG - ZO2P (Malik et al. (2019)	$\mathcal{O}(1/\epsilon \cdot \log{(1/\epsilon)})$	$\mathcal{O}(1/\epsilon \cdot \log{(1/\epsilon)})$
PG - ZO2P (Mohammadi et al. (2020))	$\mathcal{O}(\log{(1/\epsilon)})$	$\mathcal{O}(\log{(1/\epsilon)})$

Question: Can we harness synergy between Z01P (i.e., cheap cost queries) and Z02P (i.e., low variance) to reduce two-point oracle complexity for the model-free LQR problem?

Stochastic Variance-Reduced Policy Gradient (SVRPG) Difficulties

SVRPG and Reinforcement Learning: Xu et al., 2020, demonstrate a sample complexity reduction from $\mathcal{O}(\epsilon^2)$ to $\mathcal{O}(\epsilon^{5/3})$ in the local convergence analysis, i.e., $\|\nabla C(K_N)\|_F^2 \leq \epsilon$

They assume: $\mathbb{E}(\widehat{\nabla}C(K)) = \nabla C(K)$,

Stochastic Variance-Reduced Policy Gradient (SVRPG) Difficulties

SVRPG and Reinforcement Learning: Xu et al., 2020, demonstrate a sample complexity reduction from $\mathcal{O}(\epsilon^2)$ to $\mathcal{O}(\epsilon^{5/3})$ in the local convergence analysis, i.e., $\|\nabla C(K_N)\|_F^2 \leq \epsilon$

They assume: $\mathbb{E}(\widehat{\nabla}C(K)) = \nabla C(K)$,

Difficulties of SVRPG for the model-free LQR:

- K_n needs to stay stabilizing for all iterations
- Z0 gradient estimation is biased, i.e., $\mathcal{O}(r^2)$
- Z02P estimation is cost-query expensive
- Z01P estimation has a large variance

Stochastic Variance-Reduced Policy Gradient (SVRPG) Difficulties

SVRPG and Reinforcement Learning: Xu et al., 2020, demonstrate a sample complexity reduction from $\mathcal{O}(\epsilon^2)$ to $\mathcal{O}(\epsilon^{5/3})$ in the local convergence analysis, i.e., $\|\nabla C(K_N)\|_F^2 \leq \epsilon$

They assume: $\mathbb{E}(\widehat{\nabla}C(K)) = \nabla C(K)$,

Difficulties of SVRPG for the model-free LQR:

- K_n needs to stay stabilizing for all iterations
- Z0 gradient estimation is biased, i.e., $\mathcal{O}(r^2)$
- Z02P estimation is cost-query expensive
- Z01P estimation has a large variance

We are the first to combine Z01P and Z02P in a SVRPG approach

A Mixed ZO1P and ZO2P Estimation with Control Variates

Initialization: *N*, *T*, η , n_{out} , n_{in} , n_{out} , r_{in} , $K_0 \in G$

- N epochs with length T
- \bullet step-size η
- (*n*_{out}, *r*_{out}) Z02P's parameters
- (n_{in}, r_{in}) Z01P's parameters
- K₀ initial stabilizing controller

A Mixed ZO1P and ZO2P Estimation with Control Variates

Initialization: N, T, η , n_{out} , n_{in} , r_{out} , r_{in} , $K_0 \in G$

Step 1: For all $n \in \{0, 1, \dots, N-1\}$ set $K_0^{n+1} = \widetilde{K}^n = K_T^n$ and compute

$$ilde{\mu} = \widetilde{
abla} C(ilde{K}^n) ext{ with } (n_{ ext{out}}, r_{ ext{out}}) o ext{ Z02P},$$

Step 2: Within each epoch repeat for $t \in \{0, 1, \dots, T-1\}$

$$\overline{\nabla}C(K_t^{n+1}), \overline{\nabla}C(\tilde{K}^n) \text{ with } (n_{\text{in}}, r_{\text{in}}) \to \text{ Z01P},$$
$$K_{t+1}^{n+1} = K_t^{n+1} - \eta \left(\overline{\nabla}C(K_t^{n+1}) + \tilde{\mu} - \overline{\nabla}C(\tilde{K}^n)\right)$$

Repeat steps 1 and 2 and return $K_{NT} = K_T^N$.

Why Oracle Complexity Reduction?

Step 1: For all $n \in \{0, 1, ..., N-1\}$ set $K_0^{n+1} = \widetilde{K}^n = K_T^n$ and compute $\widetilde{\mu} = \widetilde{\nabla} C(\widetilde{K}^n)$ with $(n_{\text{out}}, r_{\text{out}}) \rightarrow \mathbb{Z}02P$,

Step 2: Within each epoch repeat for $t \in \{0, 1, \dots, T-1\}$

$$\overline{\nabla}C(\mathcal{K}_{t}^{n+1}), \overline{\nabla}C(\widetilde{\mathcal{K}}^{n}) \text{ with } (n_{\text{in}}, r_{\text{in}}) \to \text{ Z01P},$$
$$\mathcal{K}_{t+1}^{n+1} = \mathcal{K}_{t}^{n+1} - \eta \left(\overline{\nabla}C(\mathcal{K}_{t}^{n+1}) + \widetilde{\mu} - \overline{\nabla}C(\widetilde{\mathcal{K}}^{n})\right),$$

Why Oracle Complexity Reduction?

Step 1: For all $n \in \{0, 1, ..., N-1\}$ set $K_0^{n+1} = \widetilde{K}^n = K_T^n$ and compute $\widetilde{\mu} = \widetilde{\nabla} C(\widetilde{K}^n)$ with $(n_{\text{out}}, r_{\text{out}}) \rightarrow \mathbb{Z}02P$,

Step 2: Within each epoch repeat for $t \in \{0, 1, \dots, T-1\}$

$$\overline{\nabla}C(K_t^{n+1}), \overline{\nabla}C(\tilde{K}^n) \text{ with } (n_{\text{in}}, r_{\text{in}}) \to \text{ Z01P},$$

$$K_{t+1}^{n+1} = K_t^{n+1} - \eta \left(\overline{\nabla}C(K_t^{n+1}) + \tilde{\mu} - \overline{\nabla}C(\tilde{K}^n)\right),$$

Idea: Use **ZO2P less often** and control the inner loop gradient variance with **more ZO1P** + control variates.

Convergence: Given $K_0 \in \mathcal{G}$. Suppose that $r_{in} = r_{out} = r$,

 $n_{out} \geq \mathcal{O}(1), \ n_{in} \geq \mathcal{O}(T^2), \text{ and } \eta \text{ sufficiently small},$

then it holds that

$$\mathbb{E}\left(\mathit{C}(\mathit{K}_{\mathsf{NT}})-\mathit{C}(\mathit{K}^{\star})\right) \leq \Delta_{0}\rho^{\mathit{NT}} + \frac{\mathit{C}_{\mathsf{bias}}r^{2}}{\mathit{n}_{\mathsf{in}}}$$

where $\rho \in (0,1)$ and $\Delta_0 = C(K_0) - C(K^{\star})$.

$$\mathbb{E}\left(C(\mathcal{K}_{\mathsf{NT}}) - C(\mathcal{K}^{\star})\right) \leq \Delta_0 \rho^{\mathsf{NT}} + \frac{\mathsf{C}_{\mathsf{bias}}r^2}{\mathsf{n}_{\mathsf{in}}}$$

$$\mathbb{E}\left(C(K_{\mathsf{NT}}) - C(K^{\star})\right) \leq \Delta_0 \rho^{\mathsf{NT}} + \frac{C_{\mathsf{bias}}r^2}{n_{\mathsf{in}}}$$

$$\mathbb{E}\left(\mathcal{C}(\mathcal{K}_{\mathsf{NT}}) - \mathcal{C}(\mathcal{K}^{\star})\right) \leq \epsilon$$

for some small tolerance ϵ .

$$\mathbb{E}\left(C(K_{\mathsf{NT}}) - C(K^{\star})\right) \leq \Delta_0 \rho^{\mathsf{NT}} + \frac{C_{\mathsf{bias}}r^2}{n_{\mathsf{in}}}$$

$$\mathbb{E}\left(C(K_{\mathsf{NT}})-C(K^{\star})\right)\leq\epsilon$$

for some small tolerance ϵ .

Sample complexity: $S_c = 2Nn_{out} + NTn_{in}$

$$\mathbb{E}\left(C(K_{\mathsf{NT}}) - C(K^{\star})\right) \leq \Delta_0 \rho^{\mathsf{NT}} + \frac{C_{\mathsf{bias}}r^2}{n_{\mathsf{in}}}$$

$$\mathbb{E}\left(C(K_{\mathsf{NT}}) - C(K^{\star})\right) \leq \epsilon$$

for some small tolerance ϵ .

Sample complexity: $S_c = 2Nn_{out} + NTn_{in}$

$$egin{aligned} & \mathsf{NT} \geq \mathcal{O}(\log(1/\epsilon)) o \mathsf{N} = \mathcal{O}(\log(1/\epsilon))^{eta}, \ \ T = \mathcal{O}(\log(1/\epsilon))^{1-eta} \ & \mathsf{n}_{\mathsf{out}} = \mathcal{O}(1), \ \ \mathsf{n}_{\mathsf{in}} = \mathcal{O}(T^2) = \mathcal{O}(\log(1/\epsilon))^{2-2eta}, \ \ eta \in (0,1), \end{aligned}$$

$$\mathbb{E}\left(C(\mathcal{K}_{\mathsf{NT}}) - C(\mathcal{K}^{\star})\right) \leq \Delta_0 \rho^{\mathsf{NT}} + \frac{C_{\mathsf{bias}}r^2}{n_{\mathsf{in}}}$$

$$\mathbb{E}\left(C(K_{\mathsf{NT}}) - C(K^{\star})\right) \leq \epsilon$$

for some small tolerance ϵ .

Sample complexity: $S_c = 2Nn_{out} + NTn_{in} = O(\log(1/\epsilon))^{3-2\beta}$

$$egin{aligned} \mathsf{NT} &\geq \mathcal{O}(\log(1/\epsilon)) o \mathsf{N} = \mathcal{O}(\log(1/\epsilon))^{eta}, \ \ \mathsf{T} = \mathcal{O}(\log(1/\epsilon))^{1-eta}, \ \ \mathsf{n}_{\mathsf{out}} = \mathcal{O}(1), \ \ n_{\mathsf{in}} = \mathcal{O}(\mathsf{T}^2) = \mathcal{O}(\log(1/\epsilon))^{2-2eta}, \ \ eta \in (0,1), \end{aligned}$$

Two-point oracle complexity: $N_{ZO2P} = 2Nn_{out} = O(\log(1/\epsilon))^{\beta}$

Stability: Given $K_0 \in \mathcal{G}$. Suppose that

 n_{out}, n_{in} sufficiently large, and r, η sufficiently small

then $K_t^n \in \mathcal{G}$, with high probability, \forall epochs $n \in [N]$ of length $t \in [T]$.

Stability: Given $K_0 \in \mathcal{G}$. Suppose that

 $n_{\text{out}}, n_{\text{in}}$ sufficiently large, and r, η sufficiently small

then $K_t^n \in \mathcal{G}$, with high probability, \forall epochs $n \in [N]$ of length $t \in [T]$.

Main takeaway: By carefully controlling the quality of the inner and outer gradient estimations and not taking larger steps the learned controller provably stays within the stabilizing sub-level set.



Numerical Validation

Example 1 - $n_x = 3$, $n_u = 1$



Numerical Validation

Example 2 - $n_x = 4$, $n_u = 2$



- We propose a stochastic variance-reduced policy gradient approach for the model-free LQR problem.
- Our approach combines the benefits of one-point Z0 estimation (i.e., cheap in cost queries) and two-point Z0 estimation (i.e., lower variance) with the help of a mixed SVRPG approach.
- We prove that our approach achieves an ϵ -approximate solution with $\mathcal{O}\left(\log\left(1/\epsilon\right)^{3-2\beta}\right)$ queries, with only $\mathcal{O}\left(\log\left(1/\epsilon\right)^{\beta}\right)$ two-point query information for $\beta \in (0, 1)$.
- We prove that (sys-dyn) is stable under the learned controller.

Collaborators



To find out more:





My website

Full paper

Happy to take questions!

