

Handling Heterogeneous Curvatures in Bandit LQR Control

Yu-Hu Yan, Jing Wang, Peng Zhao (**Spotlight**)

LAMDA Group, Nanjing University, China



ICML
International Conference
On Machine Learning

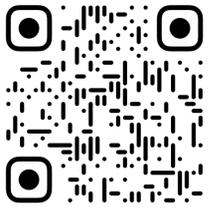


南京大學
NANJING UNIVERSITY

LAMDA
Learning And Mining from Data



Our Paper (OpenReview)



Homepage (Yu-Hu Yan)

Problem Setup

Linear Quadratic Regulator (LQR):

LQR aims to control the following **linear dynamical system**:

$$\begin{cases} \mathbf{x}_{t+1} = A\mathbf{x}_t + B\mathbf{u}_t + \boldsymbol{\xi}_t \\ \mathbf{y}_t = C\mathbf{x}_t + \mathbf{e}_t \end{cases} \quad \begin{cases} A, B: \text{known system transition matrices} \\ \mathbf{x}_t, \mathbf{u}_t: \text{state and action} \\ \boldsymbol{\xi}_t: \text{disturbances (noises)} \end{cases}$$

(partial observation)

LQR aims to minimize the **cumulative cost**:

$$\min_{\pi \in \Pi} \sum_{t=1}^T c_t(\mathbf{x}_t^\pi, \mathbf{u}_t^\pi) \triangleq \sum_{t=1}^T \mathbf{x}_t^\pi \top Q_t \mathbf{x}_t^\pi + \sum_{t=1}^T \mathbf{u}_t^\pi \top R_t \mathbf{u}_t^\pi$$

(time-varying PSD Q_t, R_t)

In **online** LQR problem, we focus on **policy regret**: [Dekel et al., ICML 2012]

$$\sum_{t=1}^T c_t(\mathbf{x}_t, \mathbf{u}_t) - \min_{\pi \in \Pi} \sum_{t=1}^T c_t(\mathbf{x}_t^\pi, \mathbf{u}_t^\pi) \quad \pi: \text{a specific policy chosen from a predefined policy class } \Pi$$

Bandit cost: Learner only observes a scalar value of $c_t(\mathbf{y}_t, \mathbf{u}_t)$.

Oblivious adversary: Costs and disturbances are chosen *in advance*.

Assumptions:

① **Stability:** The system is stable, i.e., the spectral radius $\rho(A) < 1$.
(A standard assumption that can be extended to strongly stabilizable systems.)

② **Disturbances:** **semi-adversarial** disturbances [Simchowitz et al., COLT 2020]

$$\boldsymbol{\xi}_t = \boldsymbol{\xi}_t^{\text{adv}} + \boldsymbol{\xi}_t^{\text{sto}}, \quad \mathbf{e}_t = \mathbf{e}_t^{\text{adv}} + \mathbf{e}_t^{\text{sto}}$$

Adversarial condition: $\|\boldsymbol{\xi}_t\|_2, \|\mathbf{e}_t\|_2 \leq W$.

Stochastic condition: $\text{semi-adversarial} = \text{adversarial} + \text{stochastic}$

$$\mathbb{E}[\boldsymbol{\xi}_t^{\text{sto}}] = \mathbb{E}[\mathbf{e}_t^{\text{sto}}] = \mathbf{0}, \quad \mathbb{E}[\boldsymbol{\xi}_t^{\text{sto}} \boldsymbol{\xi}_t^{\text{sto} \top}] \succeq \text{Var}_\xi \cdot I, \quad \mathbb{E}[\mathbf{e}_t^{\text{sto}} \mathbf{e}_t^{\text{sto} \top}] \succeq \text{Var}_e \cdot I.$$

(A standard assumption when dealing with **strongly convex** cost functions.)

③ **Costs:** **heterogeneous** quadratic costs

quadratic, **heterogeneous**: $Q_t, R_t \succeq \alpha_t I$ smooth: $\nabla^2 c_t(\cdot, \cdot) \preceq \beta_c I$

Lipschitz: $|c_t(\mathbf{y}, \mathbf{u}) - c_t(\mathbf{y}', \mathbf{u}')| \leq L_c R_c \|(\mathbf{y} - \mathbf{y}', \mathbf{u} - \mathbf{u}')\|_2$

Motivation: Why heterogeneous curvatures?

(i) **adaptive to true curvatures** of costs: capable of using more curvature information (since $\alpha_t \geq \alpha_{\min}$) for better performance. α_{\min} is hard to obtain before algorithms start.

(ii) **robust to corrupted quadratic costs**: allowing $\alpha_t = 0$ for some rounds where algorithm using only $\alpha_{\min} = 0$ will perform poorly.

Key Question: Is it possible design an algorithm that is

- ① **adaptive to heterogeneous curvatures for better performance** and
- ② **robust to (possibly) corrupted cost functions?**

Our Results

A general theorem for bandit LQR with heterogeneous curvatures

Theorem 1. By choosing the regularization coefficients $\{\lambda_t\}_{t=0}^T$ and step sizes $\{\eta_t\}_{t=1}^T$ properly, for Lipschitz and α_t -quadratic costs, we obtain

$$\mathbb{E}[\text{REG}_T] \leq \tilde{\mathcal{O}} \left(\inf_{\{\lambda_1^*, \dots, \lambda_T^*\}} \left\{ T^{1/3} + \lambda_{1:T-1}^* + \sum_{t=1}^{T-1} (\alpha_{1:t} + \lambda_{0:t}^*)^{-1/3} \right\} \right).$$

For Lipschitz, smooth, and α_t -quadratic costs, we obtain

$$\mathbb{E}[\text{REG}_T] \leq \tilde{\mathcal{O}} \left(\inf_{\{\lambda_1^*, \dots, \lambda_T^*\}} \left\{ \sqrt{T} + \lambda_{1:T-1}^* + \sum_{t=1}^{T-1} (\alpha_{1:t} + \lambda_{0:t}^*)^{-1/2} \right\} \right),$$

where $\{\lambda_t^*\}_{t=1}^T$ represent the optimal regularization coefficients.

➤ Coefficients $\{\lambda_t^*\}_{t=1}^T$ are only for analysis, so one may plug in any feasible $\{\lambda_t^\dagger\}_{t=1}^T$.

➤ By choosing different $\{\lambda_t^\dagger\}_{t=1}^T$ **in analysis**, we achieve different guarantees (pure convexity, pure quadraticity, corrupted quadraticity, decaying quadraticity).

	Convexity ($\alpha_t = 0$)	Quadraticity ($\alpha_t = \alpha > 0$)	Corrupted Quadraticity ($\alpha_t = \alpha \cdot \mathbb{1}_{t \notin \mathcal{T}}$)	Decaying Quadraticity ($\alpha_t = t^{-\gamma}$)
	$\tilde{\mathcal{O}}(T^{3/4})$ (Grady et al., 2020)	N/A	N/A	N/A
Lipschitz Functions	$\tilde{\mathcal{O}}(T^{3/4})$	$\tilde{\mathcal{O}}(T^{2/3})$	(when $ \mathcal{T} = T^{3/4}$) $\tilde{\mathcal{O}}(T^{2/3})$	$\begin{cases} \tilde{\mathcal{O}}(T^{2/3+\gamma/3}), & \gamma \in [0, 1/4] \\ \tilde{\mathcal{O}}(T^{2/3}), & \gamma \in (1/4, 1] \end{cases}$
	$\tilde{\mathcal{O}}(T^{2/3})$ (Cassel and Koren, 2020)	$\tilde{\mathcal{O}}(\sqrt{T})$ (Sun et al., 2023)	N/A	N/A
Smooth Functions	$\tilde{\mathcal{O}}(T^{2/3})$	$\tilde{\mathcal{O}}(\sqrt{T})$	(when $ \mathcal{T} = T^{3/4}$) $\tilde{\mathcal{O}}(\sqrt{T})$	$\begin{cases} \tilde{\mathcal{O}}(T^{1/2+\gamma/2}), & \gamma \in [0, 1/3] \\ \tilde{\mathcal{O}}(T^{2/3}), & \gamma \in (1/3, 1] \end{cases}$

Why not consider more general **strongly convex** functions?

Our results rely on the **with-history reduction** scheme [Sun et al., NeurIPS 2023], which needs quadraticity.

Reduction to BCO with Memory

Technique I: With-History Reduction to BCO with Memory

A common parametrization: **Nature's \mathbf{y}** and **disturbance-response policy**

Nature's \mathbf{y} : $\mathbf{y}_t^{\text{nat}} \triangleq \mathbf{e}_t + \sum_{i=1}^{t-1} C A^{i-1} \boldsymbol{\xi}_{t-i}$

Disturbance-Response Policy: $\mathbf{u}_t(M) = \sum_{i=0}^{M-1} M^{[i]} \mathbf{y}_{t-i}^{\text{nat}}$ [Simchowitz et al., COLT 2020]

Previous works use **truncation** to define a with-memory function:

⇒ **Drawback:** **truncation error** is hard to control in this problem.

Ours: The **with-history loss function** is defined as [Sun et al., NeurIPS 2023].

$$f_t(N_0, \dots, N_H) \triangleq c_t \left(\mathbf{y}_t^{\text{nat}} + \sum_{i=1}^H G^{[i]} \mathbf{u}_{t-i}(N_{H-i}) + \sum_{i=H+1}^{t-1} G^{[i]} \mathbf{u}_{t-i}(M_{t-i}), \mathbf{u}_t(N_H) \right).$$

↑ history more than H steps before

- **Essential idea:** **absorbing** history into c_t to form the “with-history” function, leading to a **lossless** reduction.

- **Benefit with price:** leading to **non-oblivious** $f_t(\cdot)$ for BCO problems, but fortunately can be handled thanks to **obliviousness** of $c_t(\cdot)$.

Handling Heterogeneous Curvatures

Technique II: Stability Analysis using Newton Decrement

BCO with heterogeneous curvatures: [Luo et al., COLT 2022]

Consider a general setup with online functions $\{h_t(\cdot)\}_{t=1}^T$ and decisions $\{w_t\}_{t=1}^T$.

$$\text{REG}_T \triangleq \mathbb{E} \left[\sum_{t=1}^T h_t(w_t) - \min_{w \in \mathcal{W}} \sum_{t=1}^T h_t(w) \right] \quad \begin{array}{l} \text{heterogeneous curvatures:} \\ \nabla^2 h_t(\cdot) \succeq \sigma_t I \end{array}$$

- For heterogeneous curvatures, this work follows AOGD (OGD with carefully designed **time-varying step sizes**) [Bartlett et al., NIPS 2008].

- For bandit feedback with strongly convex functions, **FTRL** is adopted.

⇒ **FTRL** with **time-varying step sizes** for BCO with heterogeneous curvatures.

Issue: A crucial term in regret analysis is $\|\nabla \psi(w)\|_{\nabla^{-2} \psi(w)}$. ($\psi(\cdot)$: FTRL regularizer)

Solution: **lifting** the domain to $\mathcal{W} \triangleq \{w = (w, 1) \mid w \in \mathcal{W}\}$ and $\|\nabla \psi(w)\|_{\nabla^{-2} \psi(w)}$ can be bounded by normal barrier (a special self-concordant barrier).

Algorithm 1 Subroutine of Luo et al. (2022)

Input: bandit value $h_t(w_t)$, curvature σ_t , last-round decision \bar{w}_t , step size η_{t+1} , regularization coefficient λ_t

- 1: Estimate gradient $\mathbf{g}_t \triangleq d(h_t(w_t) + \frac{\lambda_t}{2} \|w_t\|_2^2) / H_t^{1/2} s_t$ ↗ **gradient estimator**
- 2: Update as $\bar{w}_{t+1} = \arg \min_{w \in \mathcal{W}} F_{t+1}(w)$ ↘ **key update step**
- 3: Compute $H_{t+1} \triangleq \nabla^2 \Psi(\bar{w}_{t+1}) + \eta_{t+1} (\sigma_{t+1} I + \lambda_{0,t} I)$ ↗ **perturbation step**
- 4: Draw s_{t+1} randomly from $S^{d+1} \cap (H_{t+1}^{-1/2} \mathbf{e}_{d+1})^\perp$
- 5: Perturb $w_{t+1} = (w_{t+1}, 1) = \bar{w}_{t+1} + H_{t+1}^{-1/2} s_{t+1}$

$$F_{t+1}(w) \triangleq \sum_{s=1}^t \left(g_s^\top w + \frac{\sigma_s}{2} \|w - \bar{w}_s\|_2^2 \right) \quad (\text{APPROX})$$

$$+ \sum_{s=1}^t \frac{\lambda_s}{2} \|w - \bar{w}_s\|_2^2 + \frac{\lambda_0}{2} \|w\|_2^2 \quad (\text{REGLR-I})$$

$$+ \frac{1}{\eta_{t+1}} \Psi(w). \quad (\text{REGLR-II})$$

(notions in the lifted domain use bold italic symbols)

BCO with switching costs and heterogeneous curvatures:

Memory terms can be further transformed into **switching costs**:

$$\|w_t - w_{t+1}\|_2 \text{ and } \|w_t - w_{t+1}\|_2^2$$

We identify a **key stability** term in both **regret** and **switching cost** analysis.

$$\|\bar{w}_t - \bar{w}_{t+1}\|_{H_t} \quad \begin{array}{l} - \{\bar{w}_t\}_{t=1}^T \text{ are online decisions} \\ - H_t \text{ is the perturbation matrix, i.e., } w_t = \bar{w}_t + H_t^{1/2} s_t \end{array}$$

- Important in **regret**: closely related to $\langle g_t, \bar{w}_t - \bar{w}_{t+1} \rangle$, a crucial term in FTRL analysis.

- Important in **memory**: closely related to **switching cost** terms of $\|\bar{w}_t - \bar{w}_{t+1}\|_2$.

Newton decrement: $\lambda(w, f) \triangleq \|\nabla f(w)\|_{\nabla^{-2} f(w)}$ [Nesterov and Nemirovskii, 1994]

Proposition 1. For a self-concordant barrier f , if the Newton decrement $\lambda(w, f) \leq 1/2$, then it follows that $\|w - w^*\|_{\nabla^2 f(w)} \leq 2\lambda(w, f)$, where $w^* = \arg \min_w f(w)$.

Proof sketch: As \bar{w}_{t+1} is a minimizer of $F_{t+1}(\cdot)$, we transform H_t to $\nabla^2 F_{t+1}(\cdot)$.

Using Newton decrement, we are able to:

(i) **greatly simplify** the proof for BCO with heterogeneous curvatures.

(ii) handle heterogeneous curvatures in BCO **with memory**.

Note that our analysis for heterogeneous curvatures in BCO-M holds for **general strongly convex functions**.

Key References

- [1] Luo et al., Adaptive Bandit Convex Optimization with Heterogeneous Curvature, COLT 2022
- [2] Sun et al., Optimal Rates for Bandit Nonstochastic Control, NeurIPS 2023