# Sample Complexity of the Linear Quadratic Regulator

Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu

(Joey) Yihan Zhou

Aug 5th, 2020

The University of British Columbia

## Recap: Optimal Control

**Problem Statement**: We assume a dynamical system with state $x_t \in \mathbb{R}^n$ can be acted on by a control $u_t \in \mathbb{R}^p$ and obeys the stochastic dynamics

$$x_{t+1} = f_t(x_t, u_t, w_t),$$

where $w_t$ is a random process with $w_t$ independent of $w_{t'}$ for all $t \neq t'$. Optimal control then seeks to minimize

$$\min \quad \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T} c_t(x_t, u_t)\right]$$

$$\text{subject to} \quad x_{t+1} = f_t(x_t, u_t, w_t)$$

## Recap: Linear Quadratic Regulator

The simplest optimal control problem is the Linear Quadratic Regulator
(LQR), in which costs are a fixed quadratic function of state and control
and the dynamics are linear and time-invariant:

$$\min \quad \mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T} x_t^* Q x_t + u_{t-1}^* R u_{t-1}\right]$$

$$\text{subject to} \quad x_{t+1} = Ax_t + Bu_t + w_t$$

Here $Q$ is a $n \times n$ positive definite matrix and $R$ is a $p \times p$ positive
definite matrix. A and B are **state transition matrices** and $w_t \in \mathbb{R}^n$ is
Gaussian noise with zero-mean and covariance $\Sigma_w$.

## Infinite Time Horizon Variant of the LQR

The problem we concern is the infinite time horizon variant of LQR:

$$\min \quad \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}\left[x_t^* Q x_t + u_{t-1}^* R u_{t-1}\right]$$

$$\text{subject to} \quad x_{t+1} = A x_t + B u_t + w_t$$

When the dynamics are known, this problem has a celebrated closed form solution based on the solution of matrix **Riccati equations**. Indeed, the optimal solution sets $u_t = K x_t$ for a fixed $p \times n$ matrix $K$, and the corresponding optimal cost will serve as our gold-standard baseline to which we will compare the achieved cost of all algorithms.

## Control Performance Baseline

- We want to have baselines delineating the possible control performance achievable given a fixed amount of data collected from a system.

- Such baselines help us compare different methods and better understand tradeoffs between data collection and actions.

- Ideally, we want to find lower bounds stating the minimum amount of knowledge needed to achieve a particular performance, regardless of method.

- However, in the context of controls, even upper bounds describing the worst-case performance of competing methods are exceptionally rare.

## Unknown Dynamics

In the case when the state transition matrices are unknown, fewer results have been established about what cost is achievable.

**Assumptions:** We can conduct experiments of the following form: given some initial state $x_0$ ($x_0 = 0$), we can evolve the dynamics for $T$ time steps using any control sequence $\{u_0, \cdots, u_{T-1}\}$, measuring the resulting output $\{x_1, \cdots, x_T\}$.

**Question:** How can we solve the optimization problem? Can we get a baseline measuring how good our controller is?

How can we do this?

## Intuitive Solution

How can we do this?

A very intuitive method: conducting experiments to estimate transition matrices as $(\hat{A}, \hat{B})$ and then use $(\hat{A}, \hat{B})$ to solve the optimization problem.

Unfortunately, while this procedure might perform well given sufficient data, it is difficult to determine how many experiments are necessary in practice. Furthermore, it is easy to construct examples where the procedure fails to find a stabilizing controller.

## Coarse-ID control

- Part I: Use supervised learning to learn a coarse model of the dynamical system to be controlled. We refer to the system estimate as the nominal system.

- Part II: Using either prior knowledge or statistical tools like the bootstrap, build probabilistic guarantees about the distance between the nominal system and the true, unknown dynamics.

- Part III: Solve a robust optimization problem over controllers that optimizes performance of the nominal system while penalizing signals with respect to the estimated uncertainty, ensuring stable and robust execution.

## Digression: Robust Optimization

Robust optimization want to minimize the worst case performance given some disturbance, in our particular case:

$$\min_{\substack{\|\Delta_A\|_2 \leq \epsilon_A \\ \|\Delta_B\|_2 \leq \epsilon_B}} \sup \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}\left[ x_t^* Q x_t + u_{t-1}^* R u_{t-1} \right]$$

subject to $x_{t+1} = (\hat{A} + \Delta_A)x_t + (\hat{B} + \Delta_B)u_t + w_t$

## Part I: Model Estimation

We run experiments in which the system starts at $x_0 = 0$ and the dynamics evolve with a given input and record the resulting state observations. The set of inputs and outputs from each such experiment will be called a rollout. For system estimation, we excite the system with Gaussian noise for $N$ rollouts, each of length $T$. The resulting dataset is $\{(x_t^{(\ell)}, u_t^{(\ell)}) : 1 \leq \ell \leq N, 0 \leq t \leq T\}$, where $t$ indexes the time in one rollout and indexes independent rollouts. Therefore, we can estimate the system dynamics by

$$(\hat{A}, \hat{B}) \in \underset{(A,B)}{\operatorname{argmin}} \sum_{\ell=1}^{N} \sum_{t=0}^{T=1-1} \frac{1}{2} \left\| A x_t^{(\ell)} + B u_t^{(\ell)} - x_{t+1}^{(\ell)} \right\|_2^2.$$

## Part I: Model Estimation

First, fixing notation to simplify the presentation, let
$\Theta := [A \quad B]^* \in \mathbb{R}^{(n+p) \times n}$ and let $z_t := \begin{bmatrix} x_t \\ u_t \end{bmatrix} \in \mathbb{R}^{n+p}$. Then system
dynamics can be rewritten, for all $t \geq 0$,

$$x_{t+1}^* = z_t^* \Theta + w_t^*.$$

Then in a single rollout, we will collect

$$X := \begin{bmatrix} x_1^* \\ x_2^* \\ \vdots \\ x_T^* \end{bmatrix}, \ Z := \begin{bmatrix} z_1^* \\ z_2^* \\ \vdots \\ z_T^* \end{bmatrix}, \ W := \begin{bmatrix} w_1^* \\ w_2^* \\ \vdots \\ w_T^* a \end{bmatrix}$$

## Part I: Model Estimation

The system dynamics can be written in

$$X = Z\Theta + W.$$

Denote the data for each rollout as $(X^{(\ell)}, Z^{(\ell)}, W^{(\ell)})$. With slight abuse of notation, let $X_N$ be composed of vertically stacked $X^{(\ell)}$, and similarly for $Z_N$ and $W_N$. Then we have

$$X_N = Z_N\Theta + W_N.$$

The full data least square estimator is that

$$\hat{\Theta} = (Z_N^* Z_N)^{-1} Z_N^* X_N = \Theta + (Z_N^* Z_N)^{-1} Z_N^* W_N.$$

Then the estimation error is given by

$$E = \hat{\Theta} - \Theta = (Z_N^* Z_N)^{-1} Z_N^* W_N.$$

## Part II: Theoretical Bounds on Model Estimation

To get independent data, we only use the last sample of each rollout, $(x_T^{(\ell)}, u_{T-1}^{(\ell)}, w_{T-1}^{(\ell)})$. We excite the system with $w_t^{(\ell)} \sim \mathcal{N}(0, \sigma_w^2 I_n)$ and $u_t^{(\ell)} \sim \mathcal{N}(0, \sigma_p^2 I_p)$ i.i.d. for $1 \leq \ell \leq N$ and $0 \leq t \leq T-1$ and to minimize

$$(\hat{A}, \hat{B}) \in \underset{(A,B)}{\mathrm{argmin}} \sum_{\ell=1}^{N} \frac{1}{2} \left\| A x_{T-1}^{(\ell)} + B u_{T-1}^{(\ell)} - x_T^{(\ell)} \right\|_2^2.$$

Also we need to redefine $X_N = \begin{bmatrix} x_T^{(1)} & x_T^{(2)} & \cdots & x_T^{(N)} \end{bmatrix}^*$.

## Part II: Theoretical Bounds on Model Estimation

If we define $G_T = [A^{T-1}B \quad A^{T-2}B \quad \cdots \quad B]$ and
$F_T = [A^{T-1} \quad A^{T-2} \quad \cdots \quad I_n]$. Unroll the system dynamics and see that

$$
x_T = G_T \left[ \begin{array}{c} u_0 \\ u_1 \\ \cdots \\ u_{T-1} \end{array} \right] + F_T \left[ \begin{array}{c} w_0 \\ w_1 \\ \cdots \\ w_{T-1} \end{array} \right].
$$

Using Gaussian excitation, $u_t \sim \mathcal{N}(0, \sigma_u^2 I_p)$ gives

$$
\left[ \begin{array}{c} x_T \\ u_T \end{array} \right] \sim \mathcal{N} \left[ 0, \left[ \begin{array}{cc} \sigma_u^2 G_T G_T^* + \sigma_w^2 F_T F_T^* & 0 \\ 0 & \sigma_u^2 I_p \end{array} \right] \right].
$$

Therefore, bounding the estimation error can be achieved via proving a result on the error in random design linear regression with vector valued observations.

## Part II: Theoretical Bounds on Model Estimation

**Proposition (1.1)**
*Assume we collect data from the linear, time-invariant system initialized at $x_0 = 0$, using inputs $u_t \sim \mathcal{N}(0, \sigma_u^2 I_p)$ i.i.d. for $t = 1, \cdots, T$. Suppose that the process noise is $w_t \sim \mathcal{N}(0, \sigma_w^2 I_n)$ and that*

$$N \geq 8(n + p) + 16 \log(4/\delta).$$

*Then with probability at least $1 - \delta$, the least squares estimator using only the final sample of each trajectory satisfies both the inequality*

$$\left\| \hat{A} - A \right\|_2 \leq \frac{16\sigma_w}{\sqrt{\lambda_{min}(\sigma_u^2 G_T G_T^* + \sigma_w^2 F_T F_T^*)}} \sqrt{\frac{(n + 2p) \log(36/\delta)}{N}},$$

*and the inequality*

$$\left\| \hat{B} - B \right\|_2 \leq \frac{16\sigma_w}{\sigma_u} \sqrt{\frac{(n + 2p) \log(36/\delta)}{N}}.$$

## Part II: Theoretical Bounds on Model Estimation

- The matrices $G_T G_T^*$ and $F_T F_T^*$ are finite time controllability Gramians for the control and noise inputs, respectively. These are standard objects in control: each eigenvalue/vector pair of such a Gramian characterizes how much input energy is required to move the system in that particular direction of the state-space.

- Therefore $\lambda_{\min}(\sigma_u^2 G_T G_T^* + \sigma_w^2 F_T F_T^*)$ quantifies the least controllable, and hence most difficult to excite and estimate, mode of the system.

- On the other hand, this measure of the excitability of the system has no impact on learning $B$.

- Evaluation of this bound is still dependent on $A$ and $B$.

## Part II: Bootstrap

There are two important limitations to using such guarantees in practice to offer upper bounds on $\epsilon_A = \left\| A - \hat{A} \right\|_2$ and $\epsilon_B = \left\| B - \hat{B} \right\|_2$.

First, using only one sample per system rollout is empirically less efficient than using all available data for estimation.

Second, even optimal statistical analyses often do not recover constant factors that match practice. For purposes of robust control, it is important to obtain upper bounds on $\epsilon_A$ and $\epsilon_B$ that are not too conservative. Thus, we aim to find $\hat{\epsilon}_A$ and $\hat{\epsilon}_B$ such that $\epsilon_A \leq \hat{\epsilon}_A$ and $\epsilon_B \leq \hat{\epsilon}_B$ with high probability.

---

**Algorithm 2** Bootstrap estimation of $\epsilon_A$ and $\epsilon_B$

---

1: **Input:** confidence parameter $\delta$, number of trials $M$, data $\{(x_t^{(i)}, u_t^{(i)})\}_{\substack{1 \le i \le N \\ 1 \le t \le T}}$, and $(\widehat{A}, \widehat{B})$ a minimizer of $\sum_{\ell=1}^{N} \sum_{t=0}^{T-1} \frac{1}{2} \|Ax_t^{(\ell)} + Bu_t^{(\ell)} - x_{t+1}^{(\ell)}\|_2^2$.

2: **for** $M$ trials **do**

3:     **for** $\ell$ from 1 to $N$ **do**

4:         $\widehat{x}_0^{(\ell)} = x_0^{(\ell)}$

5:         **for** $t$ from 0 to $T-1$ **do**

6:             $\widehat{x}_{t+1}^{(\ell)} = \widehat{A}\widehat{x}_t^{(\ell)} + \widehat{B}\widehat{u}_t^{(\ell)} + \widehat{w}_t^{(\ell)}$ with $\widehat{w}_t^{(\ell)} \overset{\text{i.i.d}}{\sim} \mathcal{N}(0, \sigma_w^2 I_n)$ and $\widehat{u}_t^{(\ell)} \overset{\text{i.i.d}}{\sim} \mathcal{N}(0, \sigma_u^2 I_p)$.

7:         **end for**

8:     **end for**

9:     $(\widetilde{A}, \widetilde{B}) \in \arg\min_{(A,B)} \sum_{\ell=1}^{N} \sum_{t=0}^{T-1} \frac{1}{2} \|A\widehat{x}_t^{(\ell)} + B\widehat{u}_t^{(\ell)} - \widehat{x}_{t+1}^{(\ell)}\|_2^2$.

10:     record $\widetilde{\epsilon}_A = \|\widehat{A} - \widetilde{A}\|_2$ and $\widetilde{\epsilon}_B = \|\widehat{B} - \widetilde{B}\|_2$.

11: **end for**

12: **Output:** $\widehat{\epsilon}_A$ and $\widehat{\epsilon}_B$, the $100(1-\delta)$th percentiles of the $\widetilde{\epsilon}_A$'s and the $\widetilde{\epsilon}_B$'s.

---

For $\hat{\epsilon}_A$ and $\hat{\epsilon}_B$ estimated by Algorithm 2 we intuitively have

$$P(\left\|A - \widehat{A}\right\|_2 \le \hat{\epsilon}_A) \approx 1 - \delta, \qquad P(\left\|B - \widehat{B}\right\|_2 \le \hat{\epsilon}_B) \approx 1 - \delta.$$

### Digression: System Level Synthesis

With estimates of the system $(\hat{A}, \hat{B})$ and operator norm error bounds $(\epsilon_A, \epsilon_B)$ in hand, we now turn to control design. We will use some results from system level synthesis (SLS).

Consider linear dynamics under a fixed a static state-feedback control policy $K$, i.e., let $u_k = Kx_k$. Then we have

$$x_k = \sum_{t=1}^{k} (A + BK)^{k-t} w_{t-1}, \qquad u_k = \sum_{t=1}^{k} K(A + BK)^{k-t} w_{t-1}.$$

Let $\Phi_x(k) := (A + BK)^{k-1}$ and $\Phi_u(k) := K(A + BK)^{k-1}$. The pair $\{\Phi_x(k), \Phi_u(k)\}$ are called the closed-loop system response elements induced by the static controller $K$. Then we can rewrite the transition as

$$\begin{bmatrix} x_k \\ u_k \end{bmatrix} = \sum_{t=1}^{k} \begin{bmatrix} \Phi_x(k - t + 1) \\ \Phi_u(k - t + 1) \end{bmatrix} w_{t-1}$$

## Digression: System Level Synthesis

Then the transition can be written as

$$\Phi_x(k+1) = A\Phi_x(k) + B\Phi_u(k), \quad \Phi_x(1) = I, \quad \forall k \geq 1.$$

The previous optimization problem is non-convex in $K$ while the later is affine in $\{\Phi_x(k), \Phi_u(k)\}$.

As we work with infinite horizon problems, it is notationally more convenient to work with transfer function representations of the above objects, which can be obtained by taking a $z$-transform of their time-domain representations. The frequency domain variable $z$ can be informally thought of as the time-shift operator, i.e., $z\{x_k, x_{k+1}, \cdots\} = \{x_{k+1}, x_{k+2}, \cdots\}$, allowing for a compact representation of LTI dynamics. Then we can rewrite the constraint as

$$[zI - A \quad -B] \begin{bmatrix} \boldsymbol{\Phi}_x \\ \boldsymbol{\Phi}_u \end{bmatrix} = I.$$

The main theorem we will use from SLS is the following:

**Theorem 3.1** (State-Feedback Parameterization [59]). *The following are true:*

- *The affine subspace defined by*

$$\begin{bmatrix} zI - A & -B \end{bmatrix} \begin{bmatrix} \mathbf{\Phi}_x \\ \mathbf{\Phi}_u \end{bmatrix} = I, \ \mathbf{\Phi}_x, \mathbf{\Phi}_u \in \frac{1}{z}\mathcal{RH}_\infty \qquad (3.6)$$

  *parameterizes all system responses (3.5) from* $\mathbf{w}$ *to* $(\mathbf{x}, \mathbf{u})$*, achievable by an internally stabilizing state-feedback controller* $\mathbf{K}$.

- *For any transfer matrices* $\{\mathbf{\Phi}_x, \mathbf{\Phi}_u\}$ *satisfying (3.6), the controller* $\mathbf{K} = \mathbf{\Phi}_u \mathbf{\Phi}_x^{-1}$ *is internally stabilizing and achieves the desired system response (3.5).*

## Part III: Robust LQR Synthesis

We return to the problem setting where estimates $(\hat{A}, \hat{B})$ of a true system $(A, B)$ satisfy

$$\|\Delta_A\|_2 \le \epsilon_A, \qquad \|\Delta_A\|_2 \le \epsilon_A$$

where $\Delta_A := \hat{A} - A$ and $\Delta_B := \hat{B} - B$ and where we wish to minimize the LQR cost for the worst instantiation of the parametric uncertainty. From the previous theorem, we can reformulate the LQR problem in terms of $\{\Phi_x(k), \Phi_u(k)\}$ as:

$$\min_{\Phi_x, \Phi_u} \sigma_w^2 \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \Phi_x \\ \Phi_u \end{bmatrix} \right\|^2_{\mathcal{H}_2}$$

subject to (3.6).

## Part III: Robust LQR Synthesis

**Lemma (3.4)**

*Let the controller $\mathbf{K}$ stabilize $(\hat{A}, \hat{B})$. Then if $\mathbf{K}$ stabilizes $(A, B)$, it achieves the following LQR cost*

$$J(A, B, \mathbf{K}) = \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \mathbf{\Phi}_x \\ \mathbf{\Phi}_u \end{bmatrix} \left( I + \begin{bmatrix} \Delta_A & \Delta_B \end{bmatrix} \begin{bmatrix} \mathbf{\Phi}_x \\ \mathbf{\Phi}_u \end{bmatrix} \right)^{-1} \right\|_{\mathcal{H}_2}^2$$

With some complex derivation, we can get the following upper bound

$$\sup_{\substack{\|\Delta_A\|_2 \leq \epsilon_A \\ \|\Delta_B\|_2 \leq \epsilon_B}} J(A, B, \mathbf{K}) \leq \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \mathbf{\Phi}_x \\ \mathbf{\Phi}_u \end{bmatrix} \right\|_{\mathcal{H}_2} \frac{1}{1 - H_\alpha(\mathbf{\Phi}_x, \mathbf{\Phi}_u)} = \frac{J(\hat{A}, \hat{B}, \mathbf{K})}{1 - H_\alpha(\mathbf{\Phi}_x, \mathbf{\Phi}_u)}. \tag{3.17}$$

## Part III: Robust LQR Synthesis

To minimize the right side of (17), we can use a decomposition trick and get the following optimization problem:

$$\text{minimize}_{\gamma \in [0,1)} \frac{1}{1-\gamma} \min_{\Phi_x, \Phi_u} \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \Phi_x \\ \Phi_u \end{bmatrix} \right\|_{\mathcal{H}_2}$$

$$\text{s.t.} \begin{bmatrix} zI - \widehat{A} & -\widehat{B} \end{bmatrix} \begin{bmatrix} \Phi_x \\ \Phi_u \end{bmatrix} = I, \quad \left\| \begin{bmatrix} \frac{\epsilon_A}{\sqrt{\alpha}} \Phi_x \\ \frac{\epsilon_B}{\sqrt{1-\alpha}} \Phi_u \end{bmatrix} \right\|_{\mathcal{H}_\infty} \leq \gamma \qquad (3.18)$$

$$\Phi_x, \Phi_u \in \frac{1}{z} \mathcal{RH}_\infty.$$

We note that this optimization objective is jointly quasi-convex in $(\gamma, \Phi_x, \Phi_u)$. We further remark that any feasible solution $(\Phi_x, \Phi_u)$ generates a stabilizing controller $\mathsf{K} = \Phi_u \Phi_x^{-1}$ of the system $(A, B)$. Therefore, even if the solution is approximated, as long as it is feasible, it will be stabilizing. Two approximation techniques are finite impulse response (FIR) approximation and Lyapunov approximation.

We now return to analyzing the Coarse-ID control problem. We can upper bound the performance of the controller synthesized using the optimization (3.18).

**Theorem 4.1.** Let $J_\star$ denote the minimal LQR cost achievable by any controller for the dynamical system with transition matrices $(A, B)$, and let $K_\star$ denote the optimal controller. Let $(\widehat{A}, \widehat{B})$ be estimates of the transition matrices such that $\|\Delta_A\|_2 \leq \epsilon_A$, $\|\Delta_B\|_2 \leq \epsilon_B$. Then, if $\mathbf{K}$ is synthesized via (3.18) with $\alpha = 1/2$, the relative error in the LQR cost is

$$\frac{J(A, B, \mathbf{K}) - J_\star}{J_\star} \leq 5(\epsilon_A + \epsilon_B\|K_\star\|_2)\|\mathfrak{R}_{A+BK_\star}\|_{\mathcal{H}_\infty}, \tag{4.1}$$

as long as $(\epsilon_A + \epsilon_B\|K_\star\|_2)\|\mathfrak{R}_{A+BK_\star}\|_{\mathcal{H}_\infty} \leq 1/5$.

where $\mathfrak{R}_M := (zI - M)^{-1}$ denote the resolvent of the matrix $M$.

Together with the previous bound of $\Delta_A$ and $\Delta_B$, we have this end-to-end performance guarantee:

**Corollary 4.3.** *Let $\lambda_G = \lambda_{\min}(\sigma_u^2 G_T G_T^* + \sigma_w^2 F_T F_T^*)$, where $F_T, G_T$ are defined in (1.4). Suppose the independent data estimation procedure described in Algorithm 1 is used to produce estimates $(\widehat{A}, \widehat{B})$ and $\mathbf{K}$ is synthesized via (3.18) with $\alpha = 1/2$. Then there are universal constants $C_0$ and $C_1$ such that the relative error in the LQR cost satisfies*

$$\frac{J(A, B, \mathbf{K}) - J_\star}{J_\star} \le C_0 \sigma_w \|\mathfrak{R}_{A+BK_\star}\|_{\mathcal{H}_\infty} \left( \frac{1}{\sqrt{\lambda_G}} + \frac{\|K_\star\|_2}{\sigma_u} \right) \sqrt{\frac{(n+p)\log(1/\delta)}{N}} \tag{4.4}$$

*with probability $1 - \delta$, as long as $N \ge C_1(n+p)\sigma_u^2 \|\mathfrak{R}_{A+BK_\star}\|_{\mathcal{H}_\infty}^2 (1/\lambda_G + \|K_\star\|_2^2/\sigma_u^2)\log(1/\delta)$.*

## Summary

Coarse-ID control provides a straightforward approach to merging nonasymptotic methods from system identification with contemporary Systems Level Synthesis approaches to robust control. Indeed, many of the principles of Coarse-ID control were well established in the 90s, but fusing together an end-to-end result required contemporary analysis of random matrices and a new perspective on controller synthesis.

**Questions?**

Thank you!