

On the Robustness and Convergence of Policy Optimization in Continuous-Time Mixed $\mathcal{H}_2/\mathcal{H}_\infty$ Stochastic Control

Lekan Molu

Microsoft Research

New York City, NY 10012

Presented by **Lekan Molu** (Lay-con Mo-lu)

November 2, 2023

Talk Outline and Overview

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- Policy Optimization and Stochastic Linear Control
 - Connections to risk-sensitive control;
 - Mixed $\mathcal{H}_2/\mathcal{H}_\infty$ control theory.
- The case for convergence analysis in stochastic PO.
 - Kleinman's algorithm, *redux*.
 - Kleinman's algorithm in an iterative best response setting;
 - PO Convergence in best response settings.
- Robustness margins in model- and sampling- settings.
 - PO as a discrete-time nonlinear system;
 - Kleinman and input-to-state-stability;
 - Robust policy optimization as a small-input stable state optimization algorithm

Research Significance

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- (Deep) RL and modern AI
 - Robotic manipulation (Levine et al., 2016), text-to-visual processing (DALL-E), Atari games (Mnih et al., 2013), e.t.c.
 - Policy optimization (PO) is fundamental to modern AI algorithms' success.
 - Major success story: functional mapping of observations to policies.
 - But how does it work?

Policy Optimization – General Framework

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- PO encapsulates policy gradients (Kakade, 2001) or PG, actor-critic methods (Vrabie and Lewis, 2011), trust region PO Schulman et al. (2015), and proximal PO methods (Schulman et al., 2017).
- PG particularly suitable for complex systems.

$$\begin{aligned} & \min J(K) \\ & \text{subject to } K \in \mathcal{K} \end{aligned} \quad (1)$$

where $\mathcal{K} = \{K_1, K_2, \dots, K_n\}$.

- $J(K)$ could be tracking error, safety assurance, goal-reaching measure of performance e.t.c. required to be satisfied.

Continuous-time RL control applications

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- A little randomness in a system's mathematical model coefficients?
 - Population growth model: $dN/dt = a(t)N(t)$, $N(0) = N_0$; growth rate $a(t)$ subject to random effects e.g. $a(t) = r(t) + \text{"noise"}$.
 - We only know the distribution of "noise".
- Filtering and state estimation problems where the nature of the noise is unknown, but it is observed via sensor measurements.
 - Kalman + Bucy Filters – aerospace (Apollo, Mariner etc.).

Continuous-time RL control applications

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- Semielliptic P.D.E.s with Dirichlet boundary value problems e.g. slender flexible rods, Cosserat dynamics etc:

$$\Delta q = \sum_{i=1}^n \frac{\partial^2 q}{\partial \xi_i^2} = 0 \in \Omega, \quad q = q_{\rightarrow} \text{ on } \partial\Omega, \quad \Omega \subset \mathbb{R}^n$$

- An economic portfolio problem where the price, $p(t)$, of a stock satisfies a stochastic differential equation e.g. $dp/dt = (a + \alpha \cdot \text{"noise"})p$ for $a > 0$, $\alpha \in \text{reline}$.
- Call options pricing: The *Black-Scholes option price formula*.

Policy Optimization – Open questions

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- Gradient-based data-driven methods: prone to divergence from true system gradients.
 - Challenge I: Optimization occurs in non-convex objective landscapes.
 - Get performance certificates as a mainstay for control design: Coerciveness property (Hu et al., 2023).
 - Challenge II: Taming PG's characteristic high-variance gradient estimates (REINFORCE, NPG, Zeroth-order approx.).
 - Hello, (linear) robust (\mathcal{H}_∞ -synthesis) control!

Policy Optimization – Open questions

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- Challenge III: Under what circumstances do we have convergence to a desired equilibrium in RL settings?
- Challenge IV: Stochastic control, not deterministic control settings.
 - models involving round-off error computations in floating point arithmetic calculations; the stock market; protein kinetics.
- Challenge V: Continuous-time RL control.
 - Very little theory. Lots of potential applications encompassing rigid and soft robotics, aerospace or finance engineering, protein kinetics.

\mathcal{H}_∞ -Control Under Model Mismatch

$$\begin{aligned} dx(t) &= Ax(t)dt + Bu(t)dt + Ddw(t), \\ z(t) &= Cx(t) + Eu(t), \quad \alpha > 0; \end{aligned}$$

Algorithm 1 Search for the closed-loop \mathcal{H}_∞ -norm

```
1: Given a user-defined step size  $\eta > 0$ 
2: Set the initial upper bound on  $\gamma$  as  $\gamma_{ub} = \infty$ .
3: Initialize a buffer for possible  $\mathcal{H}_\infty$  norms for each  $K_1$ 
   to be found,  $\Gamma_{buf} = \{\}$ .
4: Initialize ordered poles  $\mathcal{P} = \{p_i \in \text{Re}(s) < 0 \mid i =$ 
    $1, 2, \dots\}$   $\triangleright p_1 < p_2 < \dots$ 
5: for  $p_i \in \mathcal{P}$  do
6:   Place  $p_i$  on (2);  $\triangleright$  (Tits and Yang, 1996)
7:   Compute stabilizing  $K_1^{p_i}$ 
8:   Find lower bound  $\gamma_{lb}$  for  $H(\gamma, K_1^{p_i})$ ;  $\triangleright$  using (22)
9:    $\Gamma_{buf}(i) = \text{get\_hinf\_norm}(T_{zw}, \gamma_{lb}, K_1^{p_i})$ .
10: end for
11: function  $\text{get\_hinf\_norm}(T_{zw}, \gamma_{lb}, K_1^{p_i})$ 
12:   while  $\gamma_{ub} = \infty$  do
13:      $\gamma := (1 + 2\eta)\gamma_{lb}$ ;
14:     Get  $\lambda_i(H(\gamma, K_1^{p_i}))$   $\triangleright$  c.f. (14)
15:     if  $\text{Re}(\Lambda) \neq \emptyset$  for  $\Lambda = \{\lambda_1, \dots, \lambda_n\}$  then
16:       Set  $\gamma_{ub} = \gamma$ ; exit
17:     else
18:       Set buffer  $\Gamma_{lb} = \{\}$ 
19:       for  $\lambda_k \in \{\text{Imag}(\Lambda)_{p-1}\}$  do  $\triangleright k = 1$  to  $K$ 
20:         Set  $m_k = \frac{1}{2}(\omega_k + \omega_{k+1})$ 
21:         Set  $\Gamma_{lb}(k) = \max\{\sigma [T_{zw}(jm_k)]\}$ ;
22:       end for
23:        $\gamma_{lb} = \max(\Gamma_{lb})$ 
24:     end if
25:     Set  $\gamma_{ub} = \frac{1}{2}(\gamma_{lb} + \gamma_{ub})$ .
26:   end while
27:   return  $\gamma_{ub}$ 
28: end function
```

Tools: Complexity, Convergence, Robustness.

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- Risk-sensitive \mathcal{H}_∞ -control (Glover, 1989) and discrete- and continuous-time mixed $\mathcal{H}_2/\mathcal{H}_\infty$ design (Khargonekar et al., 1988; Hu et al., 2023):
 - min. upper bound on \mathcal{H}_2 cost subject to satisfying a set of risk-sensitive (often \mathcal{H}_∞) constraints (Basar, 1990):

$$\min_{K \in \mathcal{K}} J(K) := \text{Tr}(P_K D D^\top) \quad (2)$$

$$\text{subject to } \mathcal{K} := \{K \mid \rho(A - BK) < 1, \|T_{zw}(K)\|_\infty < \gamma\}$$

- P_K : solution to the generalized algebraic Riccati equation (GARE);
- A, B, D, K : standard closed-loop system matrices;
- $\|T_{zw}(K)\|_\infty$: \mathcal{H}_∞ -norm of the closed-loop transfer function from a disturbance input w to output z .

Tools: Complexity, Convergence, Robustness.

Infinite-horizon

- discrete-time deterministic LQR settings (Fazel et al., 2018):

$$\min_{K \in \mathcal{K}} \mathbb{E} \sum_{t=0}^{\infty} (x_t^\top Q x_t + u_t^\top R u_t) \text{ s.t. } x_{t+1} = A x_t + B u_t, x_0 \sim \mathcal{P}_0$$

- discrete-time LQ problems under multiplicative noise (Gravell et al., 2021):

$$\begin{aligned} \min_{\pi \in \Pi} \mathbb{E}_{x_0, \{\delta_i\}, \{\gamma_i\}} \sum_{t=0}^{\infty} (x_t^\top Q x_t + u_t^\top R u_t) \\ \text{subject to } x_{t+1} = (A + \sum_{i=1}^p \delta_{ti} A_i) x_t + (B + \sum_{i=1}^q \gamma_{ti} B_i) u_t; \end{aligned}$$

(Non-exhaustive) Lit. Landscape on PO Theory

Literature landscape	Cont. time (Kalman '61, Luenberger '63)	Stochastic. LQR (Kalman '60)	Cont. Phase	LEQG or Mixed H_2/H_∞	Finite/Infinite Horizon
Fazel (2018)	No	No	Yes	No	Finite-horizon
Mohammadi (TAC -- 2020)	Yes	No	Yes	No	Finite-Horizon
Zhang (2019)	Yes	Yes (Gaussian)	Yes	Yes	Inf-horizon
Gravell (2021)	No	Multiplicative	Yes	No	Inf-horizon
Zhang (2020)	No	No	Yes	Yes	Rand-horizon
Molu (2022)	Yes	Yes (Brownian)	Yes	Yes	Inf-Horizon
Cui & Molu (2023)	Yes	Yes (Brownian)	Yes	Yes	Inf-Horizon

- Continuous-time infinite-dimensional linear systems.
 - Disturbances enter additively as random stochastic Wiener processes.
 - Many natural systems admit uncertain additive Brownian noise as diffusion processes.
 - Theoretical analysis machinery: Itô's stochastic calculus.
- Goal: keep controlled process, z , small i.e.

$$\|z\|_2 = \left(\int |z(t)|^2 dt \right)^{1/2},$$

- Under a minimizing $u(x(t)) \in \mathcal{U}$ in spite of unforeseen $w(t) \in \mathcal{W} \subseteq \mathbb{R}^q$.

Minimization Objective and Risk-Sensitive Control

- Risk-sensitive linear exponential quadratic Gaussian objective functional (Jacobson, 1973):

$$\min_{u \in \mathcal{U}} \mathcal{J}_{\text{exp}}(x_0, u, w) = \mathbb{E} \Big|_{x_0 \in \mathcal{P}_0} \exp \left[\frac{\alpha}{2} \int_0^{\infty} z^\top(t) z(t) dt \right],$$

subject to $dx(t) = Ax(t)dt + Bu(t)dt + Ddw(t)$,

$$z(t) = Cx(t) + Eu(t), \quad \alpha > 0; \quad (3)$$

- where $dw/dt = \mathcal{N}(0, W)$, $x_0 = \mathcal{N}(0, \mu)$, and $(x_0, w(t)) \subseteq (\Omega, \mathcal{F}, \mathcal{P})$.

Minimization Objective and Risk-Sensitive Control

- A Taylor series expansion of (3) reveals:

$$\mathcal{J}_{exp}(x_0, u, w) = \lim_{T \rightarrow \infty} \mathbb{E} \Big|_{x_0 \in \mathcal{P}_0} \left[\frac{\alpha}{2} \sum_{t=0}^T z^\top(t)z(t) \right] + \frac{\alpha^2}{4} \text{var} \left[\sum_{t=0}^T z^\top(t)z(t) \right]. \quad (4)$$

- Consider the variance term $\frac{\alpha^2}{4} \text{var} \left[\sum_{t=0}^T z^\top(t)z(t) \right] \rightarrow \epsilon$.
 - α a measure of risk-propensity if $\alpha > 0$;
 - α a measure of risk-aversion if $\alpha < 0$;
 - $\alpha = 0$ implies solving a classic LQP.

RL PO as a Risk-Sensitive Control Problem

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- RL (via PG) computes high-variance gradient estimates from Monte-Carlo trajectory roll-outs and bootstrapping.
- If we set $\alpha > 0$ in the LEQG problem (3), we have a controlled setting where we can study the theoretical properties of RL-based PO.
- Framework: an ADP policy iteration (PI) in a continuous PO setting.
- LEQG also interprets as a risk-attenuation algorithm.

Contributions

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- A two-loop iterative alternating best-response procedure for computing the optimal mixed-design policy;
- Rigorous convergence analyses follow for the model-based loop updates;
- In the absence of exact system models, we provide an input-to-state-stable hybrid robust stabilization scheme.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

This page is left blank intentionally.

Problem Setup

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

For $\alpha > 0$, the cost

$\mathcal{J}_{exp}(x_0, u) = \mathbb{E} \Big|_{x_0 \in \mathcal{P}_0} \exp \left[\frac{\alpha}{2} \int_0^\infty z^\top(t) z(t) dt \right]$, becomes

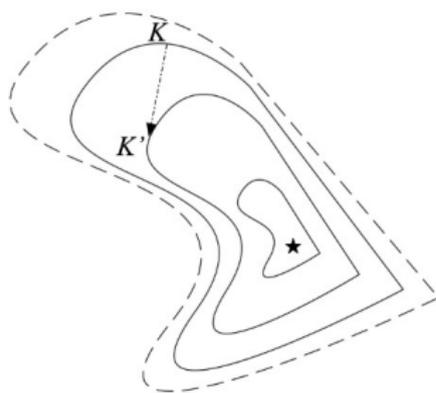
$$\mathbb{E} \Big|_{x_0 \in \mathcal{P}_0} \exp \left\{ \frac{\alpha}{2} \int_0^\infty \left[x^\top(t) Q x(t) + u^\top(t) R u(t) \right] dt \right\}, \quad (5)$$

with the associated closed loop transfer function,

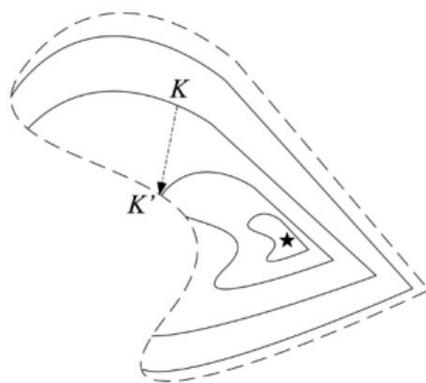
$$T_{zw}(K) = (C - EK)(sI - A + BK)^{-1}D. \quad (6)$$

Nonconvexity and Coercivity in PG

- Coercivity: iterates remain feasible and strictly separated from the infeasible set as the cost decreases.



(a) Landscape of LQR



(b) Landscape of Mixed $\mathcal{H}_2/\mathcal{H}_\infty$ Control

Figure: Coercivity property of PG on LQR and in mixed-design settings.
Credit: (Zhang et al., 2019).

Assumptions

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

- $C^T C = Q \succ 0$, $E^T (C, E) = (0, R)$ for some $R \succ 0$.
- Coercivity satisfaction: (A, B) is stabilizable;
- Optimization satisfaction: (\sqrt{Q}, A) is detectable.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

PO and Dynamic Games: Finite-horizon Gain

- Coercivity: feasibility set of optimization iterates

$$\mathcal{K} = \{ K : \lambda_i(A - B_1 K) < 0, \|T_{zw}(K)\|_\infty < \gamma \}. \quad (7)$$

- Finite-horizon optimization $u^*(t) = -K_{leqg}^* \hat{x}(t)$.
- $K_{leqg}^* = R^{-1} B^\top P_\tau$, and P_τ is the unique, symmetric, positive definite solution to the algebraic Riccati equation (ARE)

$$A^\top P_\tau + P_\tau A - P_\tau (B R^{-1} B^\top - \alpha^{-2} D D^\top) P_\tau = -Q. \quad (8)$$

(Cui and Molu, 2023, Proposition I), (Duncan, 2013) .

- ∞ -horizon case: $P^* \triangleq P_\infty = \lim_{\tau \rightarrow \infty} P_\tau$, and $K_{leqg}^* \triangleq K_\infty = \lim_{\tau \rightarrow \infty} K_\tau$ [Theorem on limit of monotonic operators (Kan, 1964)].

Solving the LEQG Problem

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- Directly solving the LEQG problem (3) in policy-gradient frameworks incurs biased gradient estimates during iterations;
- Affects risk-sensitivity preservation in infinite-horizon LTI settings (see (Zhang et al., 2021; Zhang et al., 2019));
- Workaround: an equivalent dynamic game formulation to the stochastic LQ PO problem.

Two-Player Zero-Sum Game and LEQG

- An equivalent closed-loop two-player game connection (Cui and Molu, 2023, Lemma 1):

$$\begin{aligned} \min_{u \in \mathcal{U}} \max_{\xi \in \mathcal{W}} \bar{\mathcal{J}}_{\gamma}(x_0, u, \xi) \\ \text{subject to } dx(t) = Ax(t)dt + Bu(t)dt + Ddw(t), \\ z(t) = Cx(t) + Eu(t) \end{aligned} \quad (9)$$

$$\begin{aligned} \bar{\mathcal{J}}_{\gamma}(x_0, u, \xi) = \mathbb{E}_{x_0 \sim \mathcal{P}_0, \xi(t)} \int_0^{\infty} \left[x^{\top}(t)Qx(t) + u^{\top}(t)Ru(t) \right] dt \\ - \mathbb{E}_{x_0 \sim \mathcal{P}_0, \xi(t)} \int_0^{\infty} \left[\gamma^2 \xi^{\top}(t)\xi(t) \right] dt \end{aligned}$$

, $\xi(\equiv dw) \sim \mathcal{N}(0, \Sigma)$, and $\gamma \equiv \alpha$.

Proof Sketch (Cui and Molu, 2023, Lemma 1)

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

- If a non-negative definite (n.n.d) GARE (8)'s solution exists, then a minimal realization P^* must exist.
 - Existence: the bounded real Lemma (Zhou et al., 1996).
- If $(A, Q^{\frac{1}{2}})$ is observable, then every n.n.d solution of (8), *i.e.* P^* , is positive definite.
- For a n.n.d P^* , we essentially have a Nash (equivalently a Saddle) equilibrium with $\bar{\mathcal{J}}_\gamma = \underline{\mathcal{J}}_\gamma$.

Proof Sketch (Cui and Molu, 2023, Lemma 1)

- If $\bar{\mathcal{J}}_\gamma$ is finite for some $\gamma = \hat{\gamma} > 0$, then $\bar{\mathcal{J}}_\gamma$ is bounded (if and only if the pair (A, B) is stabilizable).
- For a bounded $\bar{\mathcal{J}}_\gamma$ for some $\gamma = \hat{\gamma}$ and for optimal $K^* = R^{-1}B^\top P_{K,L}$, $L^* = \gamma^{-2}D^\top P_{K,L}$ and all $\gamma > \hat{\gamma}$, $\bar{\mathcal{J}}_\gamma$ admits the closed-loop matrices

$$A_K^* = A - BK^*, A_{K,L}^* = A_K^* + DL^*. \quad (10)$$

- Whence, the saddle-point optimal controllers are

$$u^*(x(t)) = -K^*x(t), \quad \xi^*(x(t)) = L^*x(t). \quad (11)$$

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Model-based PO

- Define $\{p, q\}_{p=1, q=1}^{\bar{p}, \bar{q}}$ where $(\bar{p}, \bar{q}) \in \mathbb{N}_+$ as nested iteration indices for a gain K_p (in an outer loop) and an alternating gain $L_q(K_p)$ (in an inner-loop).

Problem 1 (Model-Based Policy Iteration)

Given system matrices A, B, C, D, E , find the optimal controller gains $K_p, L_q(K_p)$ that robustly stabilizes (3) such that the controller gains do not leave the set of all suboptimal controllers denoted by

$$\check{\mathcal{K}} = \{(K_p, L_q(K_p)) : \lambda_i(A_K^p) < 0, \lambda_i(A_{K,L}^{p,q}) < 0, \\ \|T_{zw}(K_p, L_q(K_p))\|_\infty < \gamma \text{ for all } (p, q) \in \mathbb{N}\}. \quad (12)$$

Model-based Policy Optimization

- Further, define the following closed-loop matrix identities

$$\begin{aligned} A_K^p &= A - BK_p, & A_{K,L}^{p,q} &= A_K^p + DL_q(K_p), \\ Q_K^p &= Q + K_p^\top RK_p, & A_K^\gamma &= A_K^p + \gamma^{-2}DD^\top P_K^p. \end{aligned} \quad (13)$$

- Equation (13) informs the value iterations of the Riccati equations for the outer and inner loops.

$$A_K^{p\top} P_K^p + P_K^p A_K^p + Q_K^p + \gamma^{-2} P_K^p DD^\top P_K^p = 0, \quad (14a)$$

$$K_{p+1} = R^{-1} B^\top P_K^p. \quad (14b)$$

$$A_{K,L}^{(p,q)\top} P_{K,L}^{p,q} + P_{K,L}^{p,q} A_{K,L}^{p,q} + Q_K^p - \gamma^2 L_q^\top(K_p) L_q(K_p) = 0 \quad (15a)$$

$$K_{p+1} = R^{-1} B^\top P_{K,L}^{p,q}, \quad L_{q+1}(K_p) = \gamma^{-2} D^\top P_{K,L}^{p,q}. \quad (15b)$$

Kleinman's Algorithm

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

- An iterative algorithm for solving infinite-time Riccati equations (Kleinman, 1968).
- Based on a successive substitution method.
- For a *deterministic LTI system's* cost matrix P_d , the value iterations of P_d^k are monotonically convergent to P_d^* .
- Kleinman's algorithm as policy iteration
 - Choose a stabilizing control gain K_0 , and let $p = 0$.
 - (Policy evaluation) Evaluate the performance of K_p from the GARE's solution.
 - (Policy improvement) Improve the policy:
$$K_p = -R^{-1}B^T P_d^p.$$
 - Advance iteration $p \leftarrow p + 1$.

Model-based Policy Iteration

Algorithm 1: (Model-Based) PO via Policy Iteration

Input: Max. outer iteration \bar{p} , $q = 0$, and an $\epsilon > 0$;

Input: Desired risk attenuation level $\gamma > 0$;

Input: Minimizing player's control matrix $R \succ 0$.

- 1 Compute $(K_0, L_0) \in \mathcal{K}$; \triangleright From [24, Alg. 1];
 - 2 Set $P_{K,L}^{0,0} = Q_K^0$; \triangleright See equation (9);
 - 3 **for** $p = 0, \dots, \bar{p}$ **do**
 - 4 Compute Q_K^p and A_K^p \triangleright See equation (9);
 - 5 Obtain P_K^p by evaluating K_p on (10);
 - 6 **while** $\|P_K^p - P_{K,L}^{p,q}\|_F \leq \epsilon$ **do**
 - 7 Compute $L_{q+1}(K_p) := \gamma^{-2} D^\top P_{K,L}^{p,q}$;
 - 8 Solve (11) until $\|P_K^p - P_{K,L}^{p,q}\|_F \leq \epsilon$;
 - 9 $\bar{q} \leftarrow q + 1$
 - 10 **end**
 - 11 Compute $K_{p+1} = R^{-1} B^\top P_{K,L}^{p,\bar{q}}$ \triangleright See (11b);
 - 12 **end**
-

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Convergence Analyses: Outer Loops

Lemma 1

Under our assumptions and for the ARE (14), if $K_0 \in \mathcal{K}$, then for any $p \in \mathbb{N}_+$, we must have the following conditions for the optimal K^ and P^* ,*

- (1) $K_p \in \mathcal{K}$;
- (2) $P_K^0 \succeq P_K^1 \succeq \dots \succeq P_K^p \succeq \dots \succeq P^*$;
- (3) $\lim_{p \rightarrow \infty} \|K_p - K^*\|_F = 0$, $\lim_{p \rightarrow \infty} \|P_K^p - P^*\|_F = 0$.

Proof Sketch: The Bounded Real Lemma

Under our standard stabilizability and observability assumptions, for a stabilizing gain K , the following conditions are equivalent



$$\|\mathcal{T}(K)\|_{\infty} < \gamma;$$

- The Riccati equation

$$A_K^{\top} P_K + P_K A_K + C^{\top} C + K^{\top} R K + \gamma^{-2} P_K D D^{\top} P_K = 0, \quad (16)$$

admits a unique positive definite solution $P_K \succeq 0$ for a Hurwitz matrix $(A_K + \gamma^{-2} D D^{\top} P_K)$;

- There exists $P_K \succ 0$ such that

$$A_K^{\top} P_K + P_K A_K + Q + K^{\top} R K + \gamma^{-2} P_K D D^{\top} P_K \prec 0. \quad (17)$$

Stabilizing Proof Sketch

- At an iteration 0, find a K_0 that is stabilizing (Molu, 2023, Alg. 1), so that $K_0 \in \mathcal{K}$ by the bounded real Lemma.
- For $p > 0$, set $Q_K^{p+1} = C^\top C + K_{p+1}^\top R K_{p+1}$, the outer loop GARE is

$$A_K^{(p+1)\top} P_K^p + P_K^p A_K^{(p+1)} + \gamma^{-2} P_K^p D D^\top P_K^p + C^\top C \quad (\text{A.2}) \\ + K_{p+1}^\top R K_{p+1} + (K_{p+1} - K_p)^\top R (K_{p+1} - K_p) = 0.$$

Thus, for a stabilizing $K_{p+1} (\neq K_p)$ we must have $(K_{p+1} - K_p)^\top R (K_{p+1} - K_p) \succ 0$ so that

$$A_K^{(p+1)\top} P_K^p + P_K^p A_K^{(p+1)} + \gamma^{-2} P_K^p D D^\top P_K^p + Q_K^{p+1} \prec 0. \quad (\text{A.3})$$

- For $p > 1$, $K_p \in \mathcal{K}$. Rest: completion of squares, the bounded real Lemma, and the theorem on the “limit of monotonic operators.” (Kan, 1964).

Convergence Analysis

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

- In (Zhang et al., 2019, Theorem A.7 and A.8), the authors showed that this controller update in the outer-loop has a global sub-linear and local quadratic convergence rates.
- We now show that the outer-loop iteration has a global linear convergence rate.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Convergence Analysis: Outer Loop

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

Lemma 2

Let $\Psi = (K_{p+1} - K_p)^\top R (K_{p+1} - K_p)$; and $\Psi = \Psi^\top \succeq 0$.

Furthermore, let $\Phi \in \mathbb{R}^{n \times n}$ be Hurwitz so that

$\Theta = \int_0^\infty e^{(\Phi^\top t)} \Psi e^{(\Phi t)} dt$ and define $c(\Phi) = \log(5/4) \|\Phi\|^{-1}$.

Then, $\|\Theta\| \geq \frac{1}{2} c(\Phi) \|\Psi\|$.

Convergence Analysis: Outer Loop

Remark 1

For $A_K = A - BK$, we know from the bounded real Lemma (Zhang et al., 2019, Lemma A.1) that the Riccati equation

$$A_K^\top P_K + P_K A_K + Q_K + \gamma^{-2} P_K D D^\top P_K = 0 \quad (18)$$

admits a unique positive definite solution $P_K \succ 0$ with a Hurwitz $(A_K + \gamma^{-2} D D^\top P_K)$.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Optimality of the Iteration

Lemma 3 (Optimality of the iteration)

Consider any $K \in \mathcal{K}$, let $K' = R^{-1}B^\top P_K$ (where P_K is the solution to (18)), and $\Psi_K = (K - K')^\top R(K - K')$. If $\Psi_K = 0$, then $K = K^*$.

Proof.

Since $R \succ 0$, $\Psi_K = 0$ implies $K = K'$. Therefore at $\Psi_K = 0$, we must have $K = K'$ which implies that $P_K = P'_K$. If $K = K'$ and $P_K = P'_K$, it suffices to conclude that $K' = K \triangleq K^*$ where $K^* = R^{-1}B^\top P^*$. Hence, $\Psi_K = 0$ is tantamount to $P_K = P^*$ and $K = K^*$. \square

Bound on Cost Difference Matrix

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

Lemma 4 (Bound on Cost Difference Matrix)

For any $h > 0$, define $\mathcal{K}_h := \{K \in \mathcal{K} \mid \text{Tr}(P_K^p - P^*) \leq h\}$. For any $K \in \mathcal{K}_h$, let $K' := R^{-1}B^\top P_K^p$, where P_K^p is the p 'th iterate's solution to (18), and $\Psi_{K_p} = (K_p - K'_p)^\top R(K_p - K'_p)$. Then, there exists $b(h) > 0$, such that

$$\|P_K^p - P^*\|_F \leq b(h) \|\Psi_{K_p}\|_F.$$

Bound on Cost Difference Matrix

- For $A^* = A - BR^{-1}B^T P^* + \gamma^{-2}DD^T P^*$, rewrite the closed-loop Riccati equation as

$$\begin{aligned} & A^{*\top} P_K^p + P_K^p A^* + Q_{K_p} + (K^* - K_p)^\top R K_p' \\ & + K_p'^\top R (K^* - K_p) - \gamma^{-2} P^* D D^\top P_K^p - \gamma^{-2} P_K^p D D^\top P^* \\ & + \gamma^{-2} P_K^p D D^\top P_K^p = 0. \end{aligned} \quad (19)$$

- Then do completion of squares so that

$$\begin{aligned} & A^{*\top} (P_K^p - P^*) + (P_K^p - P^*) A^* + \Psi_{K_p} \\ & + \gamma^{-2} (P_K^p - P^*) D D^\top (P_K^p - P^*) \\ & - (K_p' - K^*)^\top R (K_p' - K^*) = 0. \end{aligned} \quad (20)$$

Proof

- Implicit function theorem: $P_K^p = f(K_p \in \mathcal{K})$, $f(\cdot) \in \mathcal{C}^n$.
- There exists a ball $\mathcal{B}_\delta(K^*) := \{K \in \mathcal{K} \mid \|K - K^*\|_F \leq \delta\}$, such that $\mathcal{A}(K)$ is invertible for any $K \in \mathcal{K}_h \cap \mathcal{B}_\delta(K^*)$.
 - $\mathcal{A}(K_p) = I_n \otimes A^{*\top} + (A - BR^{-1}B^\top P_K^p + \gamma^{-2}DD^\top P_K^p)^\top \otimes I_n$.
- Therefore, for any $K \in \mathcal{K}_h \cap \mathcal{B}_\delta(K^*)$,
 - $\|\tilde{P}_K^p\|_F \leq \underline{\sigma}^{-1}(\mathcal{A}(K_p)) \|\Psi_{K_p}\|_F$.
- Similarly, for any $K \in \mathcal{K}_h \cap \mathcal{B}_\delta^c(K^*)$, where \mathcal{B}^c is a complement of \mathcal{B} , $\Psi_{K_p} \neq 0$ and there exists a constant $b_1 > 0$ such that $\|\Psi_{K_p}\| \geq b_1$.
- Set $b_2 = \max_{K \in \mathcal{K}_h \cap \mathcal{B}_\delta(K^*)} \underline{\sigma}^{-1}(\mathcal{A}(K))$ and $b(h) = \max\{b_2, \frac{h + \text{Tr}(P^*)}{b_1}\}$, then the proof follows immediately.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Outer Loop Convergence: Exponential Stability of P_K^p

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

Theorem 2

For any $h > 0$ and $K_0 \in \mathcal{K}_h$, there exists $\alpha(h) \in \mathbb{R}$ such that $\text{Tr}(P_K^{p+1} - P^) \leq \alpha(h) \text{Tr}(P_K^p - P^*)$. That is, P^* is an exponentially stable equilibrium.*

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Convergence Analysis: Inner Loop

- Now, we analyze the monotonic convergence rate of the inner loop.
- Given arbitrary gains $K_p \in \mathcal{K}$ and $L_q(K_p) \in \mathcal{L}$, and $P_{K,L}^{p,q} \succ 0$ solution of the inner-loop Lyapunov equation, the cost matrix $P_{K,L}^{p,q}$ monotonically converges to the solution of (15).

$$A_{K,L}^{(p,q)\top} P_{K,L}^{p,q} + P_{K,L}^{p,q} A_{K,L}^{p,q} + Q_K^p - \gamma^2 L_q^\top(K_p) L_q(K_p) = 0 \quad (21a)$$

$$K_{p+1} = R^{-1} B^\top P_{K,L}^{p,q}, \quad L_{q+1}(K_p) = \gamma^{-2} D^\top P_{K,L}^{p,q}. \quad (21b)$$

Convergence Analysis: Inner Loop I

Lemma 5

Suppose that $L_0(K_0)$ is stabilizing, then for any $q \in \mathbb{N}_+$ (with $P_{K,L}^{p,\bar{q}}$ as the solution to (15)), i.e.

$$A_{K,L}^{(p,q)\top} P_{K,L}^{p,q} + P_{K,L}^{p,q} A_{K,L}^{p,q} + Q_K^p - \gamma^2 L_q^\top(K_p) L_q(K_p) = 0 \quad (22a)$$

$$K_{p+1} = R^{-1} B^\top P_{K,L}^{p,q}, \quad L_{q+1}(K_p) = \gamma^{-2} D^\top P_{K,L}^{p,q}. \quad (22b)$$

Then, the following statements hold

- 1 $A_{K,L}^{p,q}$ is Hurwitz;
- 2 $P_{K,L}^{p,\bar{q}} \succeq \dots \succeq P_K^{(p,q+1)} \succeq P_K^{p,q} \succeq \dots \succeq P_{K,L}^{p,0}$; and
- 3 $\lim_{q \rightarrow \infty} \|P_{K,L}^{p,q} - P_{K,L}^{p,\bar{q}}\|_F = 0$.

Convergence Rate – Inner Loop

Lemma 6 (Monotonic Convergence of the Inner-Loop)

For any $K \in \mathcal{K}$, let $L(K)$ be the control gain for the player w such that $A_K + DL(K)$ is Hurwitz. Let P_K^L be the solution of

$$(A_K + DL(K))^{\top} P_K^L + P_K^L (A_K + DL(K)) + Q_K - \gamma^2 L(K)^{\top} L(K) = 0. \quad (23)$$

Let $L'(K) = \gamma^{-2} D^{\top} P_K^L$ and $\Psi_K^L = \gamma^{-2} (L'(K) - L(K))^{\top} (L'(K) - L(K))$. Then, for a $c(K) = \text{Tr} \left(\int_0^{\infty} e^{(A_K + DL(K^*))t} e^{(A_K + DL(K^*))^{\top} t} dt \right)$, the following inequality holds $\text{Tr}(P_K - P_K^L) \leq \|\Psi_K^L\| c(K)$.

Convergence of the Inner Loop Iteration

Theorem 3

For a $K \in \check{\mathcal{K}}$, and for any $(p, q) \in \mathbb{N}_+$, there exists $\beta(K) \in \mathbb{R}$ such that

$$\text{Tr}(P_K^p - P_{K,L}^{p,q+1}) \leq \beta(K) \text{Tr}(P_K^p - P_{K,L}^{p,q}). \quad (24)$$

Remark 2

As seen from Lemma 5, $P_K^p - P_{K,L}^{p,q} \succeq 0$. By the norm on a matrix trace (Cui and Molu, 2023, Lemma 13) and the result of Theorem 3, we have

$\|P_K - P_{K,L}^{p,q}\|_F \leq \text{Tr}(P_K - P_{K,L}^{p,q}) \leq \beta(K) \text{Tr}(P_K)$, i.e. $P_{K,L}^{p,q}$ exponentially converges to P_K in the Frobenius norm.

Algorithm as a Policy Iteration Scheme

- Choosing a stabilizing K_p we first evaluate u 's performance by solving (14).
 - This is the policy evaluation step in PI.
- The policy is then improved in a following iteration by solving for the cost matrix in (15b);
 - This is the policy improvement step.
- Essentially, a policy iteration algorithm whereupon
 - Performance of an initial control gain K_p is first evaluated against a cost function.
 - A newer evaluation of the cost matrix $P_{K,L}^{p,q}$ is then used to improve the controller gain K_{p+1} in the outer loop.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Sampling-based PO Scheme

- A, B, C, D, E are often unavailable so that the policy evaluation step will result in biased estimates.
- There is the possibility for a divergence from the stability-robustness feasibility set $\tilde{\mathcal{K}}$:
 - When errors are present from I/O or state data;
 - Residuals from early termination of numerically solving Riccati equations;
 - Using an approximate cost function owing to inexact values of Q and R ;
 - Since the inner loop is computed in a finite number of steps;
 - In a data sampling scheme, we must guarantee the stability and robustness of the closed-loop system.

Sampling-based PO: Statement of the Problem

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

Problem 4 (Sampling-based Policy Optimization)

If A, B, C, D, E, P are all replaced by approximate matrices $\hat{A}, \hat{B}, \hat{C}, \hat{D}, \hat{E}, \hat{P}$, under what conditions will the sequences $\{\hat{P}_{K,L}^{p,q}\}_{(p,q)=1}^{\infty}$, $\{\hat{K}_p\}_{p=0}^{\infty}$, $\{\hat{L}_q\}_{q=0}^{\infty}$ converge to a small neighborhood of the optimal values $\{P_{K,L}^\}_{(p,q)=0}^{\infty}$, $\{K_p^*\}_{p=0}^{\infty}$, and $\{L_q^*\}_{q=0}^{\infty}$?*

Discrete-Time Nonlinear System Interpretation

- From assumptions, a $P_K^0 \in \mathbb{S}^n$ exists such that when applied to find a K_0 such a K_0 will be stabilizing.
- Approximation errors between the nested iteration steps yield a hybrid of a continuous-time policy gain pair $(\hat{K}_p, \hat{L}_q(\hat{K}_p))$ and a learning scheme.
 - This learning scheme is essentially a discrete sampled data from a nonlinear system (owing to errors from various sources).
- Task: under inexact loop updates, lump iterates of gain errors into system inputs to the online PO scheme;

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Discrete-Time Nonlinear System Interpretation

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- How do we converge to the optimal solution and preserve closed-loop dynamic stability?
- What does input-to-state stability (ISS) Sontag (2008) have to do with it?

Online Model-free Reparameterization

- Suppose that $\hat{P}_K^0 \in \mathbb{S}^n$ is chosen following the controllability and stabilizability assumptions.
 - Then $\hat{K}_k^1 = R^{-1}B^\top \hat{P}_K^0$ will be stabilizing since $\tilde{K}_k^1 = \hat{K}_k^1 - K_k^1 \triangleq 0$.
- Ditto argument for L_1 .

Problem 5

For $(p, q) > 0$, show that for $\tilde{K}_k^p = \hat{K}_k^p - K_k^p \triangleq 0$ so that the sequence $\{P_{K,L}^{p,q}\}_{(p,q)=0}^\infty$ converges to the locally exponentially stable $\hat{P}_{K,L}^*$.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Hybrid System Reparameterization

- Lump estimate errors as an input into the gain terms to be computed in the PO algorithm.
- With inexact outer loop update, K_{p+1} becomes biased so that the inexact outer-loop GARE value iteration involves the recursions

$$\hat{A}_K^{p\top} \hat{P}_K^p + \hat{P}_K^p \hat{A}_K^p + \hat{Q}_K^p + \gamma^{-2} \hat{P}_K^p D D^\top \hat{P}_K^p = 0, \quad (25a)$$

$$\hat{K}_{p+1} = R^{-1} B^\top \hat{P}_K^p + \tilde{K}_{p+1} \triangleq \bar{K}_{p+1} + \tilde{K}_{p+1}, \quad (25b)$$

- NB: $\hat{A}_K^p = A - B \hat{K}_p$ and $\hat{Q}_K^p = Q + \hat{K}_p^\top R \hat{K}_p$.

Discrete-Time System Closed-loop System

- Same argument for the inner-loop inexact GARE value iteration updates:

$$\hat{A}_{K,L}^{p,q\top} \hat{P}_{K,L}^{p,q} + \hat{P}_{K,L}^{p,q} \hat{A}_{K,L}^{p,q} + \hat{Q}_K^p - \gamma^2 \hat{L}_q^\top \hat{L}_q(\hat{K}_p) = 0 \quad (26a)$$

$$\hat{K}_{p+1} = R^{-1} B^\top \hat{P}_{K,L}^{p,q} + \tilde{K}_p, \quad (26b)$$

$$\hat{L}_{q+1}(\hat{K}_p) = \gamma^{-2} D^\top \hat{P}_{K,L}^{p,q} + \tilde{L}_{q+1}(\tilde{K}_p) \quad (26c)$$

$$\triangleq \bar{L}_{q+1}(\bar{K}_p) + \tilde{L}_{q+1}(\tilde{K}_p). \quad (26d)$$

- Rewrite the infinite-dimensional stochastic differential equation as the discrete-time system (for iterates $(p, q) > 0$):

$$dx = [\hat{A}_{K,L}^{p,q} x + B(\hat{K}_p x - D\hat{L}_q(K_p) + u)]dt + Ddw. \quad (27)$$

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

System Trajectories from HJB Interpretation

- On a time interval $[s, s + \delta s]$, it follows from Itô's stochastic calculus and the Hamilton-Jacobi-Bellman equation that

$$\begin{aligned} d \left[x^\top(s + \delta s) \hat{P}_{K,L}^{p,q} x(s + \delta s) - x^\top(s) \hat{P}_{K,L}^{p,q} x(s) \right] = \\ (dx)^\top \hat{P}_{K,L}^{p,q} x + x^\top \hat{P}_{K,L}^{p,q} dx + (dx)^\top \hat{P}_{K,L}^{p,q} (dx). \end{aligned} \quad (28)$$

- Along the trajectories of equation (27) and using the gains in (15), *i.e.*

$$K_{p+1} = R^{-1} B^\top P_K^{p,q}, \quad L_{q+1}(K_p) = \gamma^{-2} D^\top P_{K,L}^{p,q}.$$

System Trajectories

- The r.h.s. in (28) becomes

$$x^\top \left[\hat{A}_{K,L}^{p,q\top} \hat{P}_{K,L}^{p,q} + \hat{P}_{K,L}^{p,q} \hat{A}_{K,L}^{p,q} \right] x dt + 2x^\top \hat{P}_{K,L}^{p,q} D dw \quad (29)$$

$$+ 2x^\top \hat{P}_{K,L}^{p,q} B (K_p x - D \hat{L}_q(K_p) + u) dt + \text{Tr}(D^\top P D),$$

$$= -x^\top \hat{Q}_K^p x dt - \gamma^{-2} x^\top \hat{P}_{K,L}^{p,q} D D^\top \hat{P}_{K,L}^{p,q} x dt + \text{Tr}(D^\top \hat{P}_{K,L}^{p,q}$$

$$D) + 2x^\top \hat{P}_{K,L}^{p,q} B \left[\hat{K}_p x - D \hat{L}_q(K_p) + u \right] dt + 2x^\top \hat{P}_{K,L}^{p,q} D dw \quad (30)$$

System Trajectories via HJB Expansions

- So that

$$\begin{aligned} & x^\top(s + \delta s) \hat{P}_{K,L}^{p,q}(s + \delta s) - x^\top(s) \hat{P}_{K,L}^{p,q}(s) \\ &= \int_s^{s+\delta s} \left[(-x^\top \hat{Q}_K^p x - \gamma^2 w^\top w) dt + 2\gamma^2 x^\top \hat{L}_{q+1}^\top(K_p) dw \right] \\ &+ \int_s^{s+\delta s} 2x^\top \hat{K}_{p+1}^\top R \left[\hat{K}_p x - D \hat{L}_q(\hat{K}_p) + u \right] dt \\ &+ \int_s^{s+\delta s} \text{Tr}(D^\top \hat{P}_{K,L}^{p,q} D) dt. \end{aligned} \tag{31}$$

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Input To State System Interpretation

- System matrices $\hat{A}_{K,L}^{p,q}$, B , C , D now embedded within input and state terms: \hat{Q}_K^p , \hat{K}_{p+1} , and \hat{L}_{q+1} ;
- Retrievable via online measurements.
- We essentially end up with an input-to-state system!
- The price that we pay is that the noise feedthrough matrix D must be known precisely.
 - No marvel: in many linear stochastic system with Brownian motion, D is identity (Duncan et al., 2011; Duncan and Pasik-Duncan, 2010).

Sampling-based Scheme

- Explore system model until we achieve exact equality in

$$\hat{A}_{K,L}^{p,q} \equiv A_{K,L}^{p,q}, \hat{P}_{K,L}^{p,q}, \hat{K}_{p+1} \equiv K_{p+1}, \text{ and}$$

$$\hat{L}_{q+1}(K_p) \equiv L_{q+1}(K_p).$$

- Choose $u = -K_0x + \eta_p$ and $w = -L_0x + \eta_q$ where (η_p, η_q) is drawn uniformly at random over matrices with a Frobenium norm r similar to (Gravell et al., 2021; Fazel et al., 2018).

Sampled System Parameterization

- Consider the identities

$$\begin{aligned}x^\top \hat{Q}_K^p x &= (x^\top \otimes x^\top) \text{vec}(\hat{Q}_K^p), \\ \gamma^2 w^\top w &= \gamma^2 (w^\top \otimes w^\top) \text{vec}(I_v), \\ 2\gamma^2 x^\top \hat{L}_{q+1}^\top(\hat{K}_p) dw &= 2\gamma^2 (I_n \otimes x^\top) dw \text{vec}(\hat{L}_{q+1}^\top(\hat{K}_p)), \\ 2x^\top \hat{K}_{p+1}^\top R \hat{K}_p x &= 2(x^\top \otimes x^\top) (I_n \otimes \hat{K}_p^\top) \text{vec}(\hat{K}_{p+1}^\top R), \\ 2x^\top \hat{K}_{p+1}^\top R D \hat{L}_q(\hat{K}_p) &= 2(\hat{L}_q^\top(\hat{K}_p) D^\top \otimes x^\top) \text{vec}(\hat{K}_{p+1}^\top R), \\ 2x^\top \hat{K}_{p+1}^\top R u &= 2(u^\top \otimes x^\top) \text{vec}(\hat{K}_{p+1}^\top R), \\ \text{Tr}(D^\top \hat{P}_{K,L}^{p,q} D) &= \text{vec}^\top(D) \text{vec}(\hat{P}_{K,L}^{p,q} D).\end{aligned}\tag{32}$$

Sampled System Parameterization I

- Let $\Delta_{xx} \in \mathbb{R}^{\frac{n(n+1)}{2}l}$, $\Delta_{ww} \in \mathbb{R}^{\frac{v(v+1)}{2}l}$, $l_{xx} \in \mathbb{R}^{l \times n^2}$, and $l_{ux} \in \mathbb{R}^{l \times mn}$ for $l \in \mathbb{N}_+$

- It follows that

$$\begin{aligned}\Delta_{xx} &= [\text{vecv}(x_1), \dots, \text{vecv}(x_l)]^\top, \quad x_l = x_{l+1} - x_l, \\ \Delta_{ww} &= [\text{vecv}(w_1), \dots, \text{vecv}(w_l)]^\top, \quad w_l = w_{l+1} - w_l, \\ l_{xx} &= \left[\int_{s_0}^{s_1} x \otimes x \, dt, \dots, \int_{s_{l-1}}^{s_l} x \otimes x \, dt \right]^\top,\end{aligned}$$

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

This page is left blank intentionally.

Sampled System Parameterization

$$I_{xw} = \left[\int_{s_0}^{s_1} (I_n \otimes x) dw, \dots, \int_{s_{l-1}}^{s_l} (I_n \otimes x) dw \right]^T,$$
$$I_{ux} = \left[\int_{s_0}^{s_1} u \otimes x dt, \dots, \int_{s_{l-1}}^{s_l} u \otimes x dt \right]^T. \quad (33)$$

Next, set

$$\Theta_{K,L}^{p,q} = \left[\Delta_{xx}, -2I_{xx}(I_n \otimes \hat{K}_p^\top) + 2(\hat{L}_q^\top(\hat{K}_p)D^\top \otimes x^\top) \right. \\ \left. -2I_{ux}, -2\gamma^2 I_{xw}, -\text{vec}^\top(D)\text{vec}(\hat{P}_{K,L}^{p,q}D) \right], \quad (34a)$$

$$\Upsilon_{K,L}^{p,q} = \left[-I_{xx}\text{vec}(\hat{Q}_K^p), -\gamma^2 I_{ww}\text{vec}(I_v) \right]. \quad (34b)$$

Sampled System Parameterization

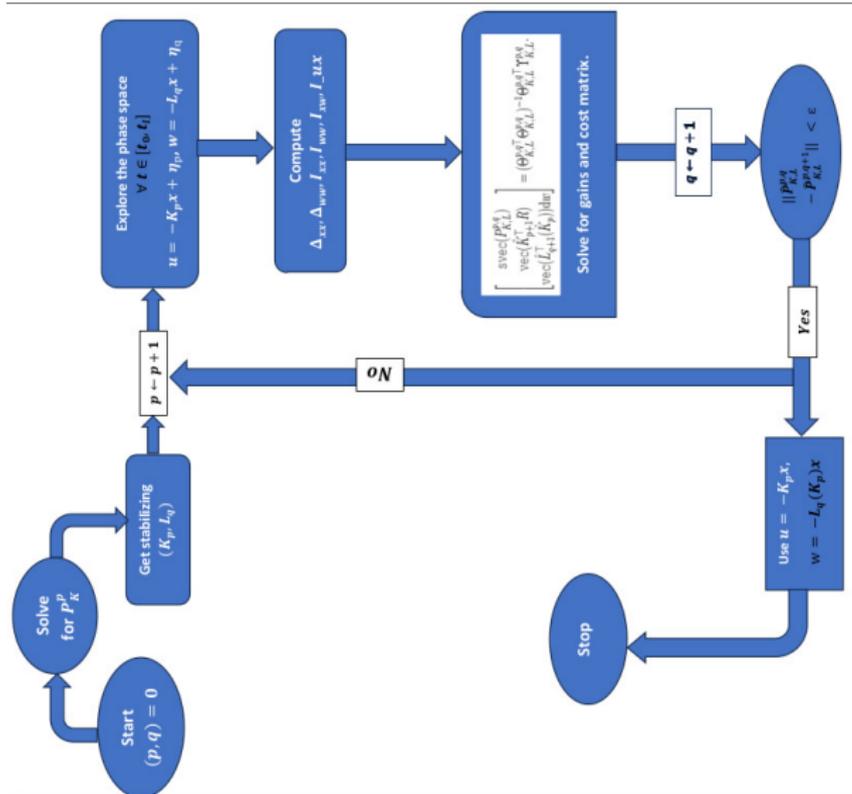
Define $\mathbf{1}_{q^2}$ as a one-vector with dimension q^2 . Thus,

$$\Theta_{K,L}^{p,q} \left[\text{svec}(P_{K,L}^{p,q}) \quad \text{vec}(\hat{K}_{p+1}^\top R) \quad \text{vec}(\hat{L}_{q+1}^\top (\hat{K}_p)) \quad \mathbf{1}_{q^2} \right]^\top = \Upsilon_{K,L}^{p,q}. \quad (35)$$

Suppose that $\Theta_{K,L}^{p,q}$ is of full rank, then we can retrieve the unknown matrices via least squares estimation *i.e.*

$$\begin{bmatrix} \text{svec}(P_{K,L}^{p,q}) \\ \text{vec}(\hat{K}_{p+1}^\top R) \\ \text{vec}(\hat{L}_{q+1}^\top (\hat{K}_p)) \\ \mathbf{1}_{q^2} \end{bmatrix} dw = (\Theta_{K,L}^{p,q \top} \Theta_{K,L}^{p,q})^{-1} \Theta_{K,L}^{p,q \top} \Upsilon_{K,L}^{p,q}. \quad (36)$$

Sampling-based Algorithm



Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

This page is left blank intentionally.

Robustness Analyses

- Define $\tilde{P} = P_K - \hat{P}_K$ and $\tilde{K} = K - \hat{K}$.
- Keep $\|\tilde{K}\| < \epsilon$, start with a $K \in \mathcal{K}$: iterates stay in \mathcal{K} .

Lemma 7 (Lemma 10, C&M, '23)

For any $K \in \mathcal{K}$, there exists an $e(K) > 0$ such that for a perturbation \tilde{K} , $K + \tilde{K} \in \mathcal{K}$, as long as $\|\tilde{K}\| < e(K)$.

Theorem 6

The inexact outer loop is small-disturbance ISS. That is, for any $h > 0$ and $\hat{K}_0 \in \mathcal{K}_h$, if $\|\tilde{K}\| < f(h)$, there exist a \mathcal{KL} -function $\beta_1(\cdot, \cdot)$ and a \mathcal{K}_∞ -function $\gamma_1(\cdot)$ such that

$$\|P_{\hat{K}}^p - P^*\| \leq \beta_1(\|P_{\hat{K}}^0 - P^*\|, p) + \gamma_1(\|\tilde{K}\|). \quad (37)$$

ISS Outer Loop Robustness Proof

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control

Contributions

Setup

Assumptions

Optimal Gain

Model-based PO

Outer loop

Stabilization and Convergence

Sampling-based PO

Discrete-time system

Sampling-based nonlinear system

Robustness Analysis

- Prelim result (Lemma 12, C&M, '23): For any $h > 0$ and $K \in \mathcal{K}_h$, let $K' = R^{-1}B^\top P_K$, where P_K is the solution of (18), and $\hat{K}' = K' + \tilde{K}$. Then, there exists $f(h) > 0$, such that $\hat{K}' \in \mathcal{K}_h$ as long as $\|\tilde{K}\| < f(h)$.
- Therefore, $\hat{K}'_K^p \in \mathcal{K}_h$ for any $p \in \mathbb{N}_+$.
- Let

$$f_1(\hat{K}') = \frac{\log(5/4)b(h)}{2n\|A_{\hat{K}'}^*\|}, f_2(\hat{K}') = \text{Tr} \left(\int_0^\infty e^{A_{\hat{K}'}^* \top t} e^{A_{\hat{K}'}^* t} dt \right).$$

ISS Outer Loop Robustness Proof



$$\underline{f}_1(h) = \inf_{\hat{K}' \in \mathcal{K}_h} f_1(\hat{K}') > 0, \bar{f}_2(h) = \sup_{\hat{K}' \in \mathcal{K}_h} f_2(\hat{K}') < \infty. \quad (38)$$

- This implies

$$\text{Tr}(P_{\hat{K}}^p - P^*) \leq [1 - \underline{f}_1(h)] \text{Tr}(P_{\hat{K}}^{p-1} - P^*) + \bar{f}_2(h) \|R\| \| \tilde{K}_{\hat{K}}^p \|^2. \quad (39)$$

- Repeating (39) for $p, p-1, \dots, 1$,

$$\text{Tr}[P_{\hat{K}}^p - P^*] \leq (1 - \underline{f}_1)^p \text{Tr}(P_{\hat{K}}^1 - P^*) + \frac{\bar{f}_2 \|R\| \| \tilde{K}_{\hat{K}} \|^2_{\infty}}{\underline{f}_1(h)}. \quad (40)$$

Outer Loop Robustness

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

It follows from (40) and (Mori, 1988, Theorem 2) that

$$\|P_{\hat{K}}^p - P^*\|_F \leq (1 - \underline{f}_1)^p \sqrt{n} \|P_{\hat{K}}^1 - P^*\|_F + \frac{\bar{f}_2 \|R\| \|\tilde{K}\|_\infty^2}{\underline{f}_1}. \quad (41)$$

As $p \rightarrow \infty$, $P_{\hat{K}}^p \rightarrow P^*$. Whence, a radius of P^* 's neighbor is proportional to $\|\tilde{K}\|_\infty^2$.

Inner Loop Robustness

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

The perturbed inner-loop iteration (26) has inexact matrix $\hat{A}_{K,L}^{p,q}$, and sequences $\{\hat{L}_{q+1}(K_p)\}_{q=0}^{\infty}$, and $\{\hat{P}_{K,L}^{p,q}\}_{q=0}^{\infty}$.

Lemma 8 (Stability of the Inner-Loop's System Matrix)

Given $K \in \check{\mathcal{K}}$, there exists a $g \in \mathbb{R}_+$, such that if $\|\tilde{L}_{q+1}(K_p)\|_F \leq g$, $\hat{A}_{K,L}^{p,q}$ is Hurwitz for all $q \in \mathbb{N}_+$.

Inner Loop Robustness

Theorem 7

Assume $\|\tilde{L}_q(K_p)\| < e$ for all $q \in \mathbb{N}_+$. There exists $\hat{\beta}(K) \in [0, 1)$, and $\lambda(\cdot) \in \check{\mathcal{K}}_\infty$, such that

$$\|\hat{P}_{K,L}^{p,q} - P_{K,L}^{p,q}\|_F \leq \hat{\beta}^{q-1}(K) \text{Tr}(P_{K,L}^{p,q}) + \lambda(\|\tilde{L}\|_\infty). \quad (42)$$

- From Theorem 7, as $q \rightarrow \infty$, $\hat{P}_{K,L}^{p,q}$ approaches the solution P_K and enters the ball centered at $P_{K,L}^{p,q}$ with radius proportional to $\|\tilde{L}\|_\infty$.
- The proposed inner-loop iterative algorithm well approximates $P_{K,L}^{p,q}$.

Transition Slide

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

This page is left blank intentionally.

Numerical Results – Car Cruise Control System

- (Åström and Murray, 2021, §3.1):

$$m \frac{dv}{dt} = \alpha_n u T(\alpha_n v) - mg C_r \operatorname{sgn}(u) - \frac{1}{2} \rho C_d A |v| v - mg \sin \theta \quad (43)$$

- $u(x(t)) = [u_1(t), u_2(t)]$ must maintain a constant velocity v (the state), whilst automatically adjusting the car's throttle, $u_1(t)$, $t \in [0, T]$
 - despite disturbances characterized by road slope changes ($u_3 = \theta$),
 - rolling friction (F_r), and
 - aerodynamic drag forces (F_d).

Numerical Results – Car Cruise Control System

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis

- Well-suited to our robust control formulation because
 - the disturbances and state variables are separable and can be lumped into the form of the stochastic differential equations;
 - it is a multiple-input (throttle, gear, vehicle speed) single-output (vehicle acceleration) system that introduces modeling challenges;
 - the entire operating range of the system is nonlinear though there is a reasonable linear bandwidth that characterize the input/output (I/O) system as we will see shortly.

Road (Disturbance) Profile

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

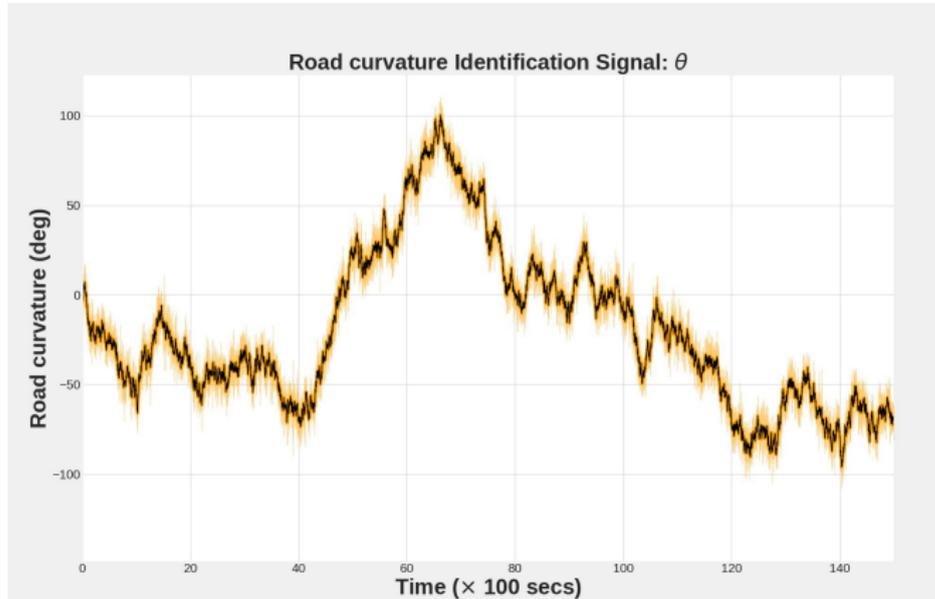
Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system
Robustness Analysis



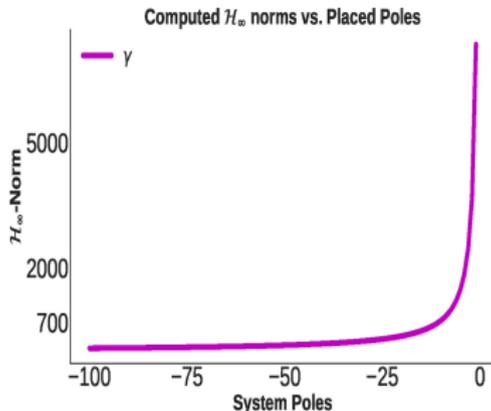
Search for initial stabilizing gain and \mathcal{H}_∞ -norm bound.

Proposition 1

(Bruinsma and Steinbuch, 1990) For all $\omega_p \in \mathbb{R}$, we have that $j\omega_p$ is an eigenvalue of the Hamiltonian $H(\gamma_1)$ if and only if γ_1 is a singular value of $T_{zw}(j\omega_p)$.

Algorithm 1 Search for the closed-loop \mathcal{H}_∞ -norm

```
1: Given a user-defined step size  $\eta > 0$ 
2: Set the initial upper bound on  $\gamma$  as  $\gamma_{ub} = \infty$ .
3: Initialize a buffer for possible  $\mathcal{H}_\infty$  norms for each  $K_1$ 
   to be found,  $\Gamma_{buf} = \{\}$ .
4: Initialize ordered poles  $\mathcal{P} = \{p_i \in \text{Re}(s) < 0 \mid i =$ 
    $1, 2, \dots\}$   $\triangleright p_1 < p_2 < \dots$ 
5: for  $p_i \in \mathcal{P}$  do
6:   Place  $p_i$  on (2);  $\triangleright$  (Tits and Yang, 1996)
7:   Compute stabilizing  $K_1^{p_i}$ 
8:   Find lower bound  $\gamma_{lb}$  for  $H(\gamma, K_1^{p_i})$ ;  $\triangleright$  using (22)
9:    $\Gamma_{buf}(i) = \text{get\_hinf\_norm}(T_{zw}, \gamma_{lb}, K_1^{p_i})$ .
10: end for
11: function  $\text{get\_hinf\_norm}(T_{zw}, \gamma_{lb}, K_1^{p_i})$ 
12:   while  $\gamma_{ub} = \infty$  do
13:      $\gamma := (1 + 2\eta)\gamma_{lb}$ ;
14:     Get  $\lambda_i(H(\gamma, K_1^{p_i}))$   $\triangleright$  c.f. (14)
15:     if  $\text{Re}(\Lambda) \neq \emptyset$  for  $\Lambda = \{\lambda_1, \dots, \lambda_n\}$  then
16:       Set  $\gamma_{ub} = \gamma$ ; exit
17:     else
18:       Set buffer  $\Gamma_{lb} = \{\}$ 
19:       for  $\lambda_k \in \{\text{Imag}(\Lambda)\}_{p-1}$  do  $\triangleright k = 1$  to  $K$ 
20:         Set  $m_k = \frac{1}{2}(\omega_k + \omega_{k+1})$ 
21:         Set  $\Gamma_{lb}(k) = \max\{\sigma[T_{zw}(jm_k)]\}$ ;
22:       end for
23:        $\gamma_{lb} = \max(\Gamma_{lb})$ 
```



Cost Matrix and Gains Convergence

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

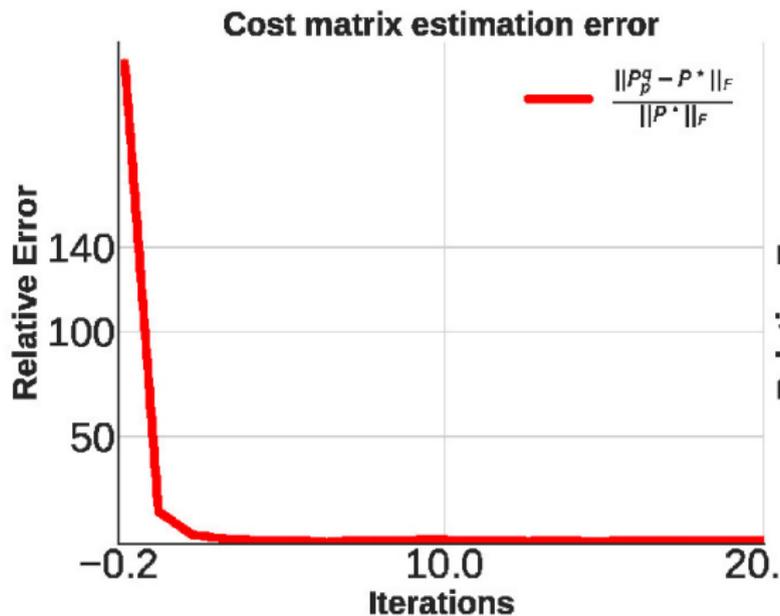
Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis



Pendulums Experiment – Comparison to NPG

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

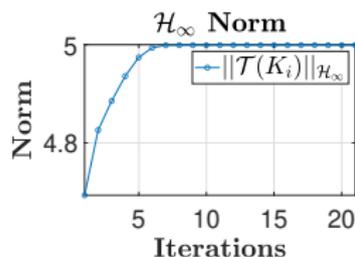
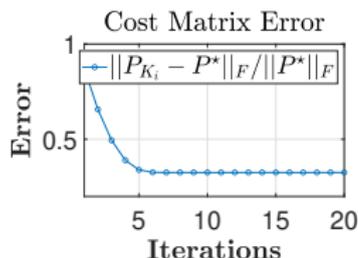
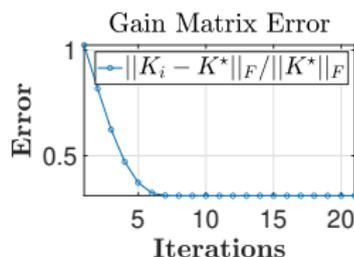
Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis



Model-free design: $\|\tilde{K}\|_\infty = 0.15$.

Pendulums Experiment – Comparison to NPG

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

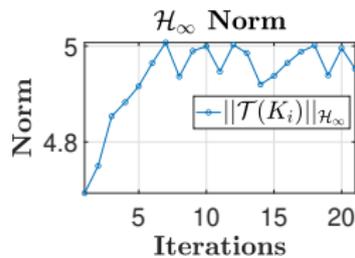
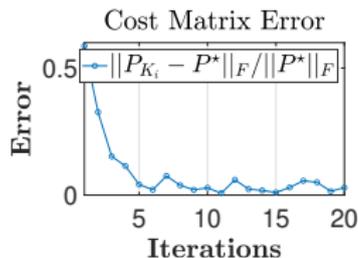
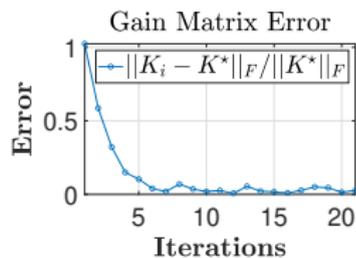
Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis



Model-based design: $\|\tilde{K}\|_\infty = 0.15$.

Double Pendulum and Acrobot Experiment – Comparison to NPG

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

Table: Computational Time: Model-based PO vs. Model-free PO vs. NPG.

Policy Optimization Computational time (secs)					
Double Inverted Pendulum			Triple Inverted Pendulum		
Model-based	Model-free	NPG	Model-based	Model-free	NPG
0.0901	0.3061	2.1649	0.1455	0.7829	2.3209

References I

- Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-End Training of Deep Visuomotor Policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- Draguna Vrabie and Frank Lewis. Adaptive dynamic programming for online solution of a zero-sum differential game. *J. Contr. Theory Appl.*, 9:353–360, 08 2011. doi: 10.1007/s11768-011-0166-4.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Bin Hu, Kaiqing Zhang, Na Li, Mehran Mesbahi, Maryam Fazel, and Tamer Başar. Toward a theoretical foundation of policy optimization for learning control policies. *Annual Review of Control, Robotics, and Autonomous Systems*, 6:123–158, 2023.
- K. Glover. Minimum entropy and risk-sensitive control: the continuous time case. In *Proceedings of the 28th IEEE Conference on Decision and Control*,, pages 388–391 vol.1, 1989.
- P.P. Khargonekar, I.R. Petersen, and M.A. Rotea. \mathcal{H}_∞ optimal control with state-feedback. *IEEE Transactions on Automatic Control*, 33(8):786–788, 1988. doi: 10.1109/9.1301.
- Tamer Basar. Minimax disturbance attenuation in ltv plants in discrete time. In *1990 American Control Conference*, pages 3112–3113. IEEE, 1990.
- Maryam Fazel, Rong Ge, Sham Kakade, and Mehran Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1467–1476. PMLR, 10–15 Jul 2018.

References III

Continuous-Time Stochastic Policy Optimization

Lekan Molu

Outline and Overview

Risk-sensitive control
Contributions

Setup

Assumptions
Optimal Gain

Model-based PO

Outer loop
Stabilization and Convergence

Sampling-based PO

Discrete-time system
Sampling-based nonlinear system

Robustness Analysis

- Eduardo D. Sontag. *Input to State Stability: Basic Concepts and Results*, pages 163–220. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.
- Tyrone E Duncan, B Maslowski, and Bozenna Pasik-Duncan. Control of some linear stochastic systems in a hilbert space with fractional brownian motions. In *2011 16th International Conference on Methods & Models in Automation & Robotics*, pages 107–110. IEEE, 2011.
- Tyrone E Duncan and Bozenna Pasik-Duncan. Stochastic linear-quadratic control for systems with a fractional brownian motion. In *49th IEEE Conference on Decision and Control (CDC)*, pages 6163–6168. IEEE, 2010.
- T. Mori. Comments on "a matrix inequality associated with bounds on solutions of algebraic Riccati and Lyapunov equation" by J. M. Saniuk and I.B. Rhodes. *IEEE Transactions on Automatic Control*, 33(11): 1088–, 1988. doi: 10.1109/9.14428.
- Karl Johan Åström and Richard M Murray. *Feedback Systems: An Introduction for Scientists and Engineers*. Princeton University Press, 2021.
- NA Bruinsma and M Steinbuch. A Fast Algorithm to Compute the H_∞ -norm of a Transfer Function Matrix. *Systems & Control Letters*, 14:287–293, 1990.