
Thompson Sampling Efficiently Learns to Control Diffusion Processes

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Diffusion processes that evolve according to linear stochastic differential equations
2 are an important family of continuous-time dynamic decision-making models.
3 Optimal policies are well-studied for them, under full certainty about the drift
4 matrices. However, little is known about data-driven control of diffusion processes
5 with uncertain drift matrices as conventional discrete-time analysis techniques are
6 not applicable. In addition, while the task can be viewed as a reinforcement learning
7 problem involving exploration and exploitation trade-off, ensuring system stability
8 is a fundamental component of designing optimal policies. We establish that
9 the popular Thompson sampling algorithm learns optimal actions fast, incurring
10 only a square-root of time regret, and also stabilizes the system in a short time
11 period. To the best of our knowledge, this is the first such result for Thompson
12 sampling in a diffusion process control problem. We validate our theoretical results
13 through empirical simulations with real parameter matrices from two settings
14 of airplane and blood glucose control. Moreover, we observe that Thompson
15 sampling significantly improves (worst-case) regret, compared to the state-of-the-
16 art algorithms, suggesting Thompson sampling explores in a more guarded fashion.
17 Our theoretical analysis involves characterization of a certain *optimality manifold*
18 that ties the local geometry of the drift parameters to the optimal control of the
19 diffusion process. We expect this technique to be of broader interest.

20 1 Introduction

21 One of the most natural reinforcement learning (RL) algorithms for controlling a diffusion process
22 with unknown parameters is based on Thompson sampling (TS) [1]: a Bayesian posterior for the
23 model is calculated based on its time evolution, and a control policy is then designed by treating
24 a sampled model from the posterior as the truth. Despite its simplicity, guaranteeing efficiency
25 and whether sampling the actions from the posterior could lead to unbounded future trajectories is
26 unknown. In fact, the only known such theoretical result for control of a diffusion process is for an
27 epsilon-greedy type policy that requires selecting purely random actions at a certain rate [2].

28 In this work, we consider a p dimensional state signal $\{\mathbf{x}_t\}_{t \geq 0}$ that obeys the (Ito) stochastic
29 differential equation (SDE)

$$d\mathbf{x}_t = (A_0\mathbf{x}_t + B_0\mathbf{u}_t) dt + d\mathbb{W}_t, \quad (1)$$

30 where the *drift matrices* A_0 and B_0 are unknown, $\mathbf{u}_t \in \mathbb{R}^q$ is the control action at any time $t \geq 0$,
31 and it is designed based on values of \mathbf{x}_s for $s \in [0, t]$. The matrix $B_0 \in \mathbb{R}^{p \times q}$ models the influence
32 of the control action on the state evolution over time, while $A_0 \in \mathbb{R}^{p \times p}$ is the (open-loop) transition
33 matrix reflecting interactions between the coordinates of the state vector \mathbf{x}_t . The diffusion term in (1)
34 consists of a non-standard Wiener process \mathbb{W}_t that will be defined in the next section. The goal is to

35 study efficient RL policies that can design u_t to minimize a quadratic cost function, defined in the
36 next section, subject to uncertainties around A_0 and B_0 .

37 At a first glance, this problem is similar to most RL problems since the optimal policy must balance
38 between the two objectives of learning the unknown matrices A_0 and B_0 (exploration) and optimally
39 selecting the control signals u_t to minimize the cost (exploitation). However, unlike most RL
40 problems that have finite or bounded-support state space, ensuring *stability*, that x_t stays bounded, is
41 a crucial part of designing optimal policies. For example, in the discrete-time version of the problem,
42 robust exploration is used to protect against unpredictably unstable trajectories [3–6].

43 **Related literature.** The existing literature studies efficiency of TS for learning optimal decisions
44 in finite action spaces [7–12]. In this stream of research, it is shown that, over time, the posterior
45 distribution concentrates around low-cost actions [13–15]. TS is also studied in further discrete-time
46 settings with the environment represented by parameters that belong to a continuum, and Bayesian and
47 frequentist regret bounds are shown for linear-quadratic regulators [16–19]. However, effectiveness
48 of TS in highly noisy environments that are modeled by diffusion processes remains unexplored to
49 date, due to technical challenges that will be described below.

50 For continuous-time linear time invariant dynamical systems, infinite-time consistency results are
51 shown under a variety of technical assumptions, followed by alternating policies that cause (small)
52 linear regrets [20–24]. From a computational viewpoint, pure exploration algorithms for computing
53 optimal policies based on multiple trajectories of the state and action data are studied as well [25–27],
54 for which a useful survey is available [28]. However, papers that study exploration versus exploitation,
55 and provide non-asymptotic estimation rates or regret bounds are limited to a few recent work about
56 offline RL or stabilized processes [29, 2, 30].

57 **Contributions.** This work, first establishes that TS learns to stabilize the diffusion process (1).
58 Specifically, in Theorem 1 of Section 3, we provide the first theoretical stabilization guarantee
59 for diffusion processes, showing that the probability of preventing the state process from growing
60 unbounded grows to 1, at an exponential rate that depends on square-root of the time length devoted
61 to stabilization. As mentioned above, for RL problems with finite state spaces, the process is by
62 definition stabilized, regardless of the policy. However, for the Euclidean state space of x_t in (1),
63 stabilization is necessary to ensure that the state and the cost do not grow unbounded.

64 Then, efficiency of TS in balancing exploration versus exploitation for minimizing a cost function
65 that has a quadratic form of both the state and the control action is shown. Indeed, we establish
66 in Theorem 2 of Section 4 that the regret TS incurs, grows as the *square-root of time*, while the
67 squared estimation error decays with the same rate. It is also shown that both the above quantities
68 grow quadratically with the dimension. To the authors’ knowledge, the presented results are the first
69 theoretical analyses of TS for learning to control diffusion processes.

70 Additionally, through extensive simulations we illustrate that TS enjoys smaller average regret and
71 substantially lower worst-case regret than the existing RL policies, thanks to its informed exploration.

72 It is important to highlight that theoretical analysis of RL policies for diffusion processes is highly
73 non-trivial. Specifically, the conventional discrete-time RL technical tools are not applicable, due
74 to uncountable cardinality of the random variables involved in a diffusion process, the unavoidable
75 dependence between them, and the high level of processing and estimation noise. To address these, we
76 make four main contributions. First, non-asymptotic and uniform upper bounds for continuous-time
77 martingales and for Ito integrals are required to quantify the estimation accuracy. For that purpose,
78 we establish concentration inequalities and show sub-exponential tail bounds for *double stochastic*
79 *integrals*. Second, one needs sharp bounds for the impact of estimation errors on eigenvalues of
80 certain non-linear matrices of the drift parameters that determine actions taken by TS policy. To tackle
81 that, we perform a novel and tight *eigenvalue perturbation-analysis* based on the approximation error,
82 dimension, and spectrum of the matrices. We also establish *Lipschitz continuity* of the control policy
83 with respect to the drift matrices, by developing new techniques based on matrix-valued curves. Third,
84 to capture evaluation of both immediate and long-term effects of sub-optimal actions, we employ
85 *Ito calculus* to bound the stochastic regret and specify effects of all problem parameters. Finally, to
86 study learning from data trajectories that the condition number of their information matrix grows
87 unbounded, we develop stochastic inequalities for *self-normalized continuous-time martingales*, and
88 *spectral analysis* of non-linear functions of random matrices.

89 **Organization.** The organization of the subsequent sections is as follows. We formulate the problem
90 in Section 2, while Algorithm 1 that utilizes TS for learning to stabilize the process and its high-
91 probability performance guarantee are presented in Section 3. Then, in Section 4, TS is considered for
92 learning to minimize a quadratic cost function, and the rates of estimation and regret are established.
93 Next, theoretical analysis are provided in Section 5, followed by real-world numerical results of
94 Section 6. Detailed proofs and auxiliary lemmas are delegated to the appendices.

95 **Notation.** The smallest (the largest) eigenvalue of matrix M , in magnitude, is denoted by $\underline{\lambda}(M)$
96 ($\bar{\lambda}(M)$). For a vector a , $\|a\|$ is the ℓ_2 norm, and for a matrix M , $\|M\|$ is the operator norm that is
97 the supremum of $\|Ma\|$ for a on the unit sphere. $\mathcal{N}(\mu, \Sigma)$ is Gaussian distribution with mean μ and
98 covariance Σ . If μ is a matrix (instead of vector), then $\mathcal{N}(\mu, \Sigma)$ denotes a distribution on matrices of
99 the same dimension as μ , such that all columns are independent and share the covariance matrix Σ . In
100 this paper, transition matrices $A \in \mathbb{R}^{p \times p}$ together with input matrices $B \in \mathbb{R}^{p \times q}$ are jointly denoted
101 by the $(p+q) \times p$ parameter matrix $\theta = [A, B]^\top$. We employ \vee (\wedge) for maximum (minimum).
102 Finally, $a \lesssim b$ expresses that $a \leq \alpha_0 b$, for some fixed constant α_0 .

103 2 Problem Statement

104 We study the problem of designing provably efficient reinforcement learning policies for minimizing
105 a quadratic cost function in an uncertain linear diffusion process. To proceed, fix the complete
106 probability space $(\Omega, \{\mathcal{F}_t\}_{t \geq 0}, \mathbb{P})$, where Ω is the sample space, $\{\mathcal{F}_t\}_{t \geq 0}$ is a continuous-time
107 filtration (i.e., increasing sigma-fields), and \mathbb{P} is the probability measure defined on \mathcal{F}_∞ .

108 The state comprises the diffusion process \mathbf{x}_t in (1), where $\theta_0 = [A_0, B_0]^\top \in \mathbb{R}^{(p+q) \times p}$ is the un-
109 known drift parameter. The diffusion term in (1) follows infinitesimal variations of the p dimensional
110 Wiener process $\{\mathbb{W}_t\}_{t \geq 0}$. That is, $\{\mathbb{W}_t\}_{t \geq 0}$ is a multivariate Gaussian process with independent
111 increments and with the stationary covariance matrix $\Sigma_{\mathbb{W}}$, such that for all $0 \leq s_1 \leq s_2 \leq t_1 \leq t_2$,

$$\begin{bmatrix} \mathbb{W}_{t_2} - \mathbb{W}_{t_1} \\ \mathbb{W}_{s_2} - \mathbb{W}_{s_1} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0_p \\ 0_p \end{bmatrix}, \begin{bmatrix} (t_2 - t_1)\Sigma_{\mathbb{W}} & 0_{p \times p} \\ 0_{p \times p} & (s_2 - s_1)\Sigma_{\mathbb{W}} \end{bmatrix} \right). \quad (2)$$

112 Existence, construction, continuity, and non-differentiability of Wiener processes are well-known [31].
113 It is standard to assume that $\Sigma_{\mathbb{W}}$ is positive definite, which is a common condition in learning-based
114 control [28, 29, 2, 30] to ensure accurate estimation over time.

115 The RL policy designs the action $\{\mathbf{u}_t\}_{t \geq 0}$, based on the observed system state by the time, as well as
116 the previously applied actions, to minimize the long-run average cost

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \int_0^T [\mathbf{x}_t^\top, \mathbf{u}_t^\top] Q \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} dt, \quad \text{for } Q = \begin{bmatrix} Q_x & Q_{xu} \\ Q_{xu}^\top & Q_u \end{bmatrix}. \quad (3)$$

117 Above, the cost is determined by the positive definite matrix Q , where $Q_x \in \mathbb{R}^{p \times p}$, $Q_u \in \mathbb{R}^{q \times q}$,
118 $Q_{xu} \in \mathbb{R}^{p \times q}$. In fact, Q determines the weights of different coordinates of $\mathbf{x}_t, \mathbf{u}_t$ in the cost function,
119 so that the policy aims to make the states small, by deploying small actions. The cost matrix Q is
120 assumed known to the policy. Formally, the problem is to minimize (3) by the policy

$$\mathbf{u}_t = \hat{\pi} \left(Q, \{\mathbf{x}_s\}_{0 \leq s \leq t}, \{\mathbf{u}_s\}_{0 \leq s < t} \right). \quad (4)$$

121 Without loss of generality, and for the ease of presentation, we follow the canonical formulation
122 that sets $Q_{xu} = 0$; one can simply convert the case $Q_{xu} \neq 0$ to the canonical form, by employing a
123 rotation to $\mathbf{x}_t, \mathbf{u}_t$ [32–35]. It is well-known that if, hypothetically, the truth θ_0 was known, an optimal
124 policy π_{opt} could be explicitly found by solving the continuous-time algebraic Riccati equation. That
125 is, for a generic drift matrix $\theta = [A, B]^\top$, finding the symmetric $p \times p$ matrix $P(\theta)$ that satisfies

$$A^\top P(\theta) + P(\theta)A - P(\theta)BQ_u^{-1}B^\top P(\theta) + Q_x = 0. \quad (5)$$

126 This means, for the true parameter $\theta_0 = [A_0, B_0]^\top$, we can let $P(\theta_0)$ solve the above equation, and
127 define the policy

$$\pi_{\text{opt}} : \quad \mathbf{u}_t = -Q_u^{-1}B_0^\top P(\theta_0) \mathbf{x}_t, \quad \forall t \geq 0. \quad (6)$$

128 It is known that the linear time-invariant policy π_{opt} minimizes the average cost in (3) [32–35].

129 **Definition 1** *The process in (1) is stabilizable, if all eigenvalues of $\bar{A} = A_0 + B_0K$ have negative*
 130 *real-parts, for a matrix K . Such K, \bar{A} are called a stabilizer and the stable closed-loop matrix.*

131 We assume that the process (1) with the drift parameter θ_0 is stabilizable. Therefore, $P(\theta_0)$ exists,
 132 is unique, and can be computed using continuous-time Riccati differential equations similar to (5),
 133 except that the zero matrix on the right-hand side will be replaced by the derivative of $P(\theta)$ [32–35].
 134 Furthermore, it is known that real-parts of all eigenvalues of $\bar{A}_0 = A_0 - B_0Q_u^{-1}B_0^\top P(\theta_0)$ are
 135 negative, i.e., $|\bar{\lambda}(\exp(\bar{A}_0t))| < 1$, which means the matrix $\exp(\bar{A}_0t)$ decays exponentially fast as
 136 t grows [32–35]. In the sequel, we use (5) and refer to the solution $P(\theta)$ for different stabilizable θ .
 137 More details about the above optimal feedback policy can be found in the aforementioned references.

138 In absence of exact knowledge of θ_0 , a policy $\hat{\pi}$ collects data and leverages it to approximate π_{opt}
 139 in (6). Therefore, at all (finite) times, there is a gap between the cost of $\hat{\pi}$, compared to that of π_{opt} .
 140 The cumulative performance degradation due to this gap is the *regret* of the policy $\hat{\pi}$, that we aim to
 141 minimize. Technically, whenever the control action u_t is designed by the policy $\hat{\pi}$ according to (4),
 142 concatenate the resulting state and input signals to get the observation $z_t(\hat{\pi}) = [x_t^\top, u_t^\top]^\top$. If it is
 143 clear from the context, we drop $\hat{\pi}$. Similarly, $z_t(\pi_{\text{opt}})$ denotes the observation signal of π_{opt} . Now,
 144 the regret at time T is defined by:

$$\text{Reg}_{\hat{\pi}}(T) = \int_0^T \left(\|Q^{1/2}z_t(\hat{\pi})\|^2 - \|Q^{1/2}z_t(\pi_{\text{opt}})\|^2 \right) dt.$$

145 A secondary objective is the learning accuracy of θ_0 from the single trajectory of the data generated
 146 by $\hat{\pi}$. Letting $\hat{\theta}_t$ be the parameter estimate at time t , we are interested in scaling of $\|\hat{\theta}_t - \theta_0\|$ with
 147 respect to t, p , and q .

148 3 Stabilizing the Diffusion Process

149 This section focuses on establishing that Thompson sampling (TS) learns to stabilize the diffusion
 150 process (1). First, let us intuitively discuss the problem of stabilizing unknown diffusion processes.
 151 Given that the optimal policy in (6) stabilizes the process in (1), a natural candidate to obtain
 152 a stable process under uncertainty of the drift matrices A_0, B_0 , is a linear feedback of the form
 153 $u_t = Kx_t$. So, letting $\bar{A} = A_0 + B_0K$, the solution of (1) is the Ornstein–Uhlenbeck process
 154 $x_t = e^{\bar{A}t}x_0 + \int_0^t e^{\bar{A}(t-s)}d\mathbb{W}_s$ [31]. Thus, if real-part of an eigenvalue of \bar{A} is non-negative, then
 155 the magnitude of x_t grows unbounded with t [31]. Therefore, addressing instabilities of this form is
 156 important, *prior* to minimizing the cost. Otherwise, the regret grows (super) linearly with time. In
 157 particular, if A_0 has some eigenvalue(s) with non-negative real-part(s), then it is necessary to employ
 158 feedback to preclude instabilities.

159 In addition to minimizing the cost, the algebraic Riccati equation in (5) provides a reliable and
 160 widely-used framework for stabilization, as discussed after (6). Accordingly, due to uncertainty
 161 about θ_0 , one can solve (5) and find $P(\hat{\theta})$, only for an approximation $\hat{\theta}$ of θ_0 . Then, we expect to
 162 stabilize the system in (1) by applying a linear feedback that is designed for the approximate drift
 163 matrix $\hat{\theta}$. Technically, we need to ensure that all eigenvalues of $A_0 - B_0Q_u^{-1}\hat{B}^\top P(\hat{\theta})$ lie in the
 164 open left half-plane. To ensure that these requirements are met in a sustainable manner, the main
 165 challenges are

- 166 (i) fast and accurate learning of θ_0 so that after a short time period, a small error $\hat{\theta} - \theta_0$ is guaranteed,
- 167 (ii) specifying the effect of the error $\hat{\theta} - \theta_0$, on stability of $A_0 - B_0Q_u^{-1}\hat{B}^\top P(\hat{\theta})$, and
- 168 (iii) devising a remedy for the case that the stabilization procedure fails.

169 Note that the last challenge is unavoidable, since learning from finite data can never be perfectly accu-
 170 rate, and so any finite-time stabilization procedure has a (possibly small) positive failure probability.

171 Algorithm 1 addresses the above challenges by applying additionally randomized control actions, and
 172 using them to provide a posterior belief \mathcal{D} about θ_0 . Note that the posterior is *not* concentrated at
 173 θ_0 , and a sample $\hat{\theta}$ from \mathcal{D} approximates θ_0 , crudely. Still, the theoretical analysis of Theorem 1

174 indicates that the failure probability of Algorithm 1 decays exponentially fast with the length of the
 175 time interval it is executed. Importantly, this small failure probability can shrink further by repeating
 176 the procedure of sampling from \mathcal{D} . So, stabilization under uncertainty is guaranteed, after a limited
 177 time of interacting with the environment.

178 To proceed, let $\{w_n\}_{n=0}^{\kappa}$ be a sequence of independent Gaussian vectors with the distribution
 179 $w_n \sim \mathcal{N}(0, \sigma_w^2 I_q)$, for some fixed constant σ_w . Suppose that we aim to devote the time length τ
 180 to collect observations for learning to stabilize. Note that since stabilization is performed before
 181 moving forward to the main objective of minimizing the cost functions, the stabilization time length
 182 τ is desired to be as short as possible. We divide this time interval of length τ to κ sub-intervals
 183 of equal length, and randomize an initial linear feedback policy by adding $\{w_n\}_{n=0}^{\kappa}$. That is, for
 184 $n = 0, 1, \dots, \kappa - 1$, Algorithm 1 employs the control action

$$\mathbf{u}_t = K\mathbf{x}_t + w_n, \quad \text{for } \frac{n\tau}{\kappa} \leq t < \frac{(n+1)\tau}{\kappa}, \quad (7)$$

185 where K is an initial stabilizing feedback so that all eigenvalues of $A_0 + B_0K$ lie in the open
 186 left half-plane. In practice, such K is easily found using physical knowledge of the model, e.g.,
 187 via conservative control sequence for an airplane [36, 37]. However, note that such actions are
 188 sub-optimal involving large regrets. Therefore, they are only temporarily applied, for the sake of
 189 data collection. Then, the data collected during the time interval $0 \leq t \leq \tau$ will be utilized by the
 190 algorithm to determine the posterior belief \mathcal{D}_τ , as follows. Recalling the notation $\mathbf{z}_t^\top = [\mathbf{x}_t^\top, \mathbf{u}_t^\top]$,
 191 let $\hat{\boldsymbol{\mu}}_0, \hat{\boldsymbol{\Sigma}}_0$ be the mean and the precision matrix of a prior normal distribution on $\boldsymbol{\theta}_0$ (using the
 192 notation defined in Section 1 for random matrices). Nonetheless, if there is no such prior, we simply
 193 let $\hat{\boldsymbol{\mu}}_0 = \mathbf{0}_{(p+q) \times p}$ and $\hat{\boldsymbol{\Sigma}}_0 = I_{p+q}$. Then, define

$$\hat{\boldsymbol{\Sigma}}_\tau = \hat{\boldsymbol{\Sigma}}_0 + \int_0^\tau \mathbf{z}_s \mathbf{z}_s^\top ds, \quad \hat{\boldsymbol{\mu}}_\tau = \hat{\boldsymbol{\Sigma}}_\tau^{-1} \left(\hat{\boldsymbol{\Sigma}}_0 \hat{\boldsymbol{\mu}}_0 + \int_0^\tau \mathbf{z}_s d\mathbf{x}_s^\top \right). \quad (8)$$

194 Using $\hat{\boldsymbol{\Sigma}}_\tau \in \mathbb{R}^{(p+q) \times (p+q)}$ together with the mean matrix $\hat{\boldsymbol{\mu}}_\tau$, Algorithm 1 forms the posterior belief

$$\mathcal{D}_\tau = \mathcal{N}(\hat{\boldsymbol{\mu}}_\tau, \hat{\boldsymbol{\Sigma}}_\tau^{-1}), \quad (9)$$

195 about the drift parameter $\boldsymbol{\theta}_0$. So, as defined in the notation, the posterior distribution of every column
 196 $i = 1, \dots, p$ of $\boldsymbol{\theta}_0$, is an independent multivariate normal with the covariance matrix $\hat{\boldsymbol{\Sigma}}_\tau^{-1}$, while the
 mean is the column i of $\hat{\boldsymbol{\mu}}_\tau$. The final step of Algorithm 1 is to output a sample $\hat{\boldsymbol{\theta}}$ from \mathcal{D}_τ .

Algorithm 1 : Stabilization under Uncertainty

Inputs: initial feedback K , stabilization time length τ
for $n = 0, 1, \dots, \kappa - 1$ **do**
 while $n\tau\kappa^{-1} \leq t < (n+1)\tau\kappa^{-1}$ **do**
 Apply control action \mathbf{u}_t in (7)
 end while
end for
Calculate $\hat{\boldsymbol{\Sigma}}_\tau, \hat{\boldsymbol{\mu}}_\tau$ according to (8)
Return sample $\hat{\boldsymbol{\theta}}$ from the distribution \mathcal{D}_τ in (9)

197

198 Next, to establish performance guarantees for Algorithm 1, let us quantify the *ideal* stability by

$$\zeta_0 = -\log \bar{\lambda}(\exp[A_0 - B_0 Q_u^{-1} B_0^\top P(\boldsymbol{\theta}_0)]). \quad (10)$$

199 By definition, ζ_0 is positive. In fact, it is the smallest distance between the imaginary axis in the
 200 complex-plane, and the eigenvalues of the transition matrix $\bar{A}_0 = A_0 - B_0 Q_u^{-1} B_0^\top P(\boldsymbol{\theta}_0)$, under
 201 the optimal policy in (6). Since $\boldsymbol{\theta}_0$ is unavailable, it is *not* realistic to expect that after applying
 202 a policy based on $\hat{\boldsymbol{\theta}}$ given by Algorithm 1, real-parts of all eigenvalues of the resulting matrix
 203 $A_0 - B_0 Q_u^{-1} \hat{B}^\top P(\hat{\boldsymbol{\theta}})$ are at most $-\zeta_0$. However, ζ_0 is crucial in studying stabilization, such
 204 that stabilizing controllers for systems with larger ζ_0 can be learned faster. The exact effect of this
 205 quantity, as well as those of other properties of the diffusion process, are formally established in the
 206 following result. Informally, the failure probability of Algorithm 1 decays exponentially with $\tau^{1/2}$.

207 **Theorem 1 (Stabilization Guarantee)** For the sample $\hat{\theta}$ given by Algorithm 1, let \mathcal{E}_τ be the failure
 208 event that $A_0 - B_0 Q_u^{-1} \hat{B}^\top P(\hat{\theta})$ has an eigenvalue in the closed right half-plane. Then, if $\kappa \gtrsim \tau^2$,
 209 we have

$$\log \mathbb{P}(\mathcal{E}_\tau) \lesssim - \frac{\lambda(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2}{\bar{\lambda}(\Sigma_{\mathbb{W}}) \vee \sigma_w^2} \frac{1 \wedge \zeta_0^p}{1 \vee \|K\|^3} \sqrt{\frac{\tau}{p^3 q}}. \quad (11)$$

210 The above result indicates that more heterogeneity in coordinates of the Wiener noise renders
 211 stabilization harder. Moreover, using (10), the term $1 \wedge \zeta_0^p$ reflects that less stable diffusion processes
 212 with smaller ζ_0 , are significantly harder to stabilize under uncertainty. Also as one can expect, larger
 213 dimensions make learning to stabilize harder. This is contributed by higher number of parameters
 214 to learn, as well as higher sensitivity of eigenvalues for processes of larger dimensions. Finally, the
 215 failure probability decays as $\tau^{1/2}$, mainly because continuous-time martingales have sub-exponential
 216 distributions, unlike sub-Gaussianity of discrete-time counterparts [38–40].

217 4 Thompson Sampling for Efficient Control: Algorithm and Theory

218 In this section, we proceed towards analysis of Thompson sampling (TS) for minimizing the quadratic
 219 cost in (3), and show that it efficiently learns the optimal control actions. That is, TS balances the
 220 exploration versus exploitation, such that its regret grows with (nearly) the square-root rate, as time
 221 grows. In the sequel, we introduce Algorithm 2 and discuss the conceptual and technical frameworks
 222 it relies on. Then, we establish efficiency by showing regret bounds in terms of different problem
 223 parameters and provide the rates of estimating the unknown drift matrices.

224 In Algorithm 2, first the learning-based stabilization Algorithm 1 is run during the time period
 225 $0 \leq t < \tau_0$. So, according to Theorem 1, the optimal feedback of $\hat{\theta}_0$ stabilizes the system with a
 226 high probability, as long as τ_0 is sufficiently large. Note that if growth of the state indicates that
 227 Algorithm 1 failed to stabilize, one can repeat sampling from \mathcal{D}_{τ_0} . So, we can assume that the
 228 evolution of the controlled diffusion process remains stable when Algorithm 2 is being executed. On
 229 the other hand, the other benefit of running Algorithm 1 at the beginning is that it performs an initial
 230 exploration phase that will be utilized by Algorithm 2 to minimize the regret.

231 Then, in order to learn the optimal policy π_{opt} with minimal sub-optimality, RL algorithms need
 232 to cope with a fundamental challenge, commonly known as the exploration-exploitation dilemma.
 233 To see that, first note that an acceptable policy that aims to have sub-linear regret, needs to take
 234 near-optimal control actions in a long run; $\mathbf{u}_t \approx -Q_u^{-1} B_0^\top P(\theta_0) \mathbf{x}_t$. Although such policies exploit
 235 well and their control actions are close to that of π_{opt} , their regret grows large since they fail to
 236 explore. Technically, the trajectory of observations $\{\mathbf{z}_t\}_{t \geq 0}$ is not rich enough to provide accurate
 237 estimations, since in $\mathbf{z}_t^\top = [\mathbf{x}_t^\top, \mathbf{u}_t^\top]$, the signal \mathbf{u}_t is (almost) a linear function of the state signal
 238 \mathbf{x}_t , and so does not contribute towards gathering information about the unknown parameter θ_0 .
 239 Conversely, for sufficient explorations, RL policies need to take actions that deviate from those of
 240 π_{opt} , which imposes large regret (as quantified in Lemma 7). Accordingly, the above trade-off needs
 241 to be delicately balanced; what we show that TS does.

242 Algorithm 2 is episodic; the parameter estimates $\hat{\theta}_n$ are updated only at the end of the episodes at
 243 times $\{\tau_n\}_{n=0}^\infty$, while during every episode, actions are taken as if $\hat{\theta}_n = [\hat{A}_n, \hat{B}_n]^\top$ is the unknown
 244 truth θ_0 . That is, for $\tau_{n-1} \leq t < \tau_n$, using $P(\hat{\theta}_n)$ in (5), we let $\mathbf{u}_t = -Q_u^{-1} \hat{B}_n^\top P(\hat{\theta}_n) \mathbf{x}_t$.
 245 Then, for each $n = 1, 2, \dots$, at time τ_n , we use all the observations collected so far, to find $\hat{\Sigma}_{\tau_n}, \hat{\mu}_{\tau_n}$
 246 according to (8). Next, we use them to sample $\hat{\theta}_n$ from the posterior \mathcal{D}_{τ_n} in (9).

247 The episodes in Algorithm 2 are chosen such that their end points satisfy

$$0 < \underline{\alpha} \leq \inf_{n \geq 0} \frac{\tau_{n+1} - \tau_n}{\tau_n} \leq \sup_{n \geq 0} \frac{\tau_{n+1} - \tau_n}{\tau_n} \leq \bar{\alpha} < \infty, \quad (12)$$

248 for some fixed constants $\underline{\alpha}, \bar{\alpha}$. Broadly speaking, (12) lets the episode lengths of Algorithm 2 scale
 249 properly to avoid unnecessary updates of parameter estimates, while at the same time performing
 250 sufficient exploration. To see that, first note that since $\hat{\Sigma}_\tau$ grows with τ , the estimation error $\hat{\theta}_n - \theta_0$
 251 decays (at best polynomially fast) with τ_n . So, until ensuring that updating the posterior yields to

252 significantly better approximations, it will not be beneficial to update it, sample from it, and solve
 253 (5). So, the period $\tau_{n+1} - \tau_n$ that the data up to time τ_n is utilized, is set to be as long as $\underline{\alpha}\tau_n$.
 254 On the other hand, the above period cannot be too long, since we aim to improve the parameter
 255 estimates after collecting enough new observations; $\tau_{n+1} \leq (1 + \bar{\alpha})\tau_n$. A simple setting is to let
 256 $\underline{\alpha} = \bar{\alpha}$, which yields to exponential episodes $\tau_n = \tau_0 (1 + \bar{\alpha})^n$. Note that for TS in continuous time,
 257 posterior updates should be limited to sufficiently-apart time points. Otherwise, repetitive updates are
 258 computationally impractical, and also can degrade the performance by preventing control actions
 259 from having enough time to effectively influence.

Algorithm 2 : Thompson Sampling for Efficient Control of Diffusion Processes

Inputs: stabilization time τ_0
 Calculate sample $\hat{\theta}_0$ by running Algorithm 1 for time τ_0
for $n = 1, 2, \dots$ **do**
 while $\tau_{n-1} \leq t < \tau_n$ **do**
 Apply control action $u_t = -Q_u^{-1} \hat{B}_{n-1}^\top P(\hat{\theta}_{n-1}) x_t$
 end while
 Letting $\hat{\Sigma}_{\tau_n}, \hat{\mu}_{\tau_n}$ be as (8), sample $\hat{\theta}_n$ from \mathcal{D}_{τ_n} given in (9)
end for

260 We show next that Algorithm 2 addresses the exploration-exploitation trade-off efficiently. To see
 261 the intuition, consider the sequence of posteriors \mathcal{D}_{τ_n} . The explorations Algorithm 2 performs by
 262 sampling $\hat{\theta}_n$ from \mathcal{D}_{τ_n} , depends on $\hat{\Sigma}_{\tau_n}$. Now, if hypothetically $\underline{\lambda}(\hat{\Sigma}_{\tau_n})$ is not large enough, then
 263 \mathcal{D}_{τ_n} does not sufficiently concentrate around $\hat{\mu}_{\tau_n}$ and so $\hat{\theta}_n$ will probably deviate from the previous
 264 samples $\{\hat{\theta}_i\}_{i=1}^{n-1}$. So, the algorithm explores more and obtains richer data z_t by diversifying the
 265 control signal u_t . This renders the next mean $\hat{\mu}_{\tau_{n+1}}$ a more accurate approximation of θ_0 , and also
 266 makes $\underline{\lambda}(\hat{\Sigma}_{\tau_{n+1}})$ grow faster than before. Thus, the next posterior $\mathcal{D}_{\tau_{n+1}}$ provides a better sample
 267 with smaller estimation error $\hat{\theta}_{n+1} - \theta_0$. Similarly, if a posterior is excessively concentrated, in a few
 268 episodes the posteriors adjust accordingly to the proper level of exploration. Hence, TS eventually
 269 balances the exploration versus the exploitation. This is formalized below.

270 **Theorem 2 (Regret and Estimation Rates)** *Parameter estimates and regret of Algorithm 2, satisfy*

$$\begin{aligned} \|\hat{\theta}_n - \theta_0\|^2 &\lesssim \frac{\bar{\lambda}(\Sigma_{\mathbb{W}})}{\underline{\lambda}(\Sigma_{\mathbb{W}})} \log(1 + \bar{\alpha}) \quad (p+q)p \quad \tau_n^{-1/2} \log \tau_n, \\ \text{Reg}(T) &\lesssim (\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2) \tau_0 + \frac{\bar{\lambda}(\Sigma_{\mathbb{W}})^2}{\underline{\lambda}(\Sigma_{\mathbb{W}})} \frac{\bar{\alpha} \|P(\theta_0)\|^6}{\log(\underline{\alpha} + 1) \underline{\lambda}(Q)^6} \quad (p+q)p \quad T^{1/2} \log T. \end{aligned}$$

271 In the above regret and estimation rates, and similar to Theorem 1, $\bar{\lambda}(\Sigma_{\mathbb{W}})/\underline{\lambda}(\Sigma_{\mathbb{W}})$ reflects the
 272 impact of heterogeneity in coordinates of \mathbb{W}_t on the quality of learning. Also, larger $\log(1 + \bar{\alpha})$
 273 corresponds to longer episodes which compromises the estimation. Further, $p(p+q)$ shows that larger
 274 number of parameters linearly worsens the learning accuracy. In the regret bound, $\|P(\theta_0)\|/\underline{\lambda}(Q)$
 275 indicates effect of the true problem parameters θ_0, Q . Finally, $(\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2) \tau_0$ captures the initial
 276 phase that Algorithm 1 is run for stabilization, which takes sub-optimal control actions as in (7).

277 5 Intuition and Summary of the Analysis

278 The goal of this section is to provide a high-level roadmap of the proofs of Theorems 1 and 2, and
 279 convey the main intuition behind the analysis. Complete proofs and the technical lemmas are provided
 280 in Appendices A and B, respectively.

281 **Summary of the Proof of Theorem 1.** The main steps involve analyzing the estimation (Lemma 4),
 282 studying its effect on the solutions of (5) (Lemma 12), and characterizing impact of errors in entries
 283 of parameter matrices on their eigenvalues (Lemma 5). Next, we elaborate on these steps.

284 We show that the error satisfies $\|\hat{\theta} - \theta_0\| \lesssim p(p+q)^{1/2}\tau^{-1/2}$ (Lemma 4). More precisely, the error
 285 depends mainly on total strength of the observation signals z_t , which are captured in the precision
 286 matrix $\hat{\Sigma}_\tau$, as well as total interactions between the signal z_t and the noise \mathbb{W}_t in the form of the
 287 stochastic integral matrix $\int_0^\tau z_t d\mathbb{W}_t^\top$. However, we establish an upper bound $\bar{\lambda}(\hat{\Sigma}_\tau^{-1}) \lesssim \tau^{-1}$, that
 288 indicates the concentration rate of the posterior \mathcal{D}_τ (Lemma 3). Similarly, thanks to the randomization
 289 signal w_n , the signals z_t are diverse enough to effectively explore the set of matrices $\theta = [A, B]^\top$,
 290 leading to accurate approximation of θ_0 by the posterior mean matrix $\hat{\mu}_\tau$. Then, to bound the error
 291 terms caused by the Wiener noise \mathbb{W}_t , we establish the rate $p(p+q)^{1/2}\tau^{1/2}$ (Lemma 2). Indeed,
 292 we show that the entries of this error matrix are continuous-time martingales, and use exponential
 293 inequalities for quadratic forms and double stochastic integrals [39, 38] to establish that they have a
 294 sub-exponential distribution.

295 Moreover, the error rate of the feedback satisfies a similar property; $\|\hat{B}^\top P(\hat{\theta}) - B_0^\top P(\theta_0)\| \lesssim$
 296 $p(p+q)^{1/2}\tau^{-1/2}$ (Lemma 12). So, letting $\bar{A} = A_0 - B_0 Q_u^{-1} \hat{B}^\top P(\hat{\theta})$ and $\bar{A}_0 = A_0 -$
 297 $B_0 Q_u^{-1} B_0^\top P(\theta_0)$, it holds that $\|\bar{A} - \bar{A}_0\| \lesssim p(p+q)^{1/2}\tau^{-1/2}$. Next, to consider the effect
 298 of the errors on the eigenvalues of \bar{A} , we compare them to the eigenvalues of \bar{A}_0 , which are bounded
 299 by $-\zeta_0$ in (10). To that end, we establish a novel and tight perturbation analysis for eigenvalues of
 300 matrices, with respect to their entries and spectral properties (Lemma 5). Using that, we show that
 301 the difference between the eigenvalues of \bar{A} and \bar{A}_0 scales as $(1 \vee r^{1/2} \|\bar{A} - \bar{A}_0\|)^{1/r}$, where r is
 302 the size of the largest block in the Jordan block-diagonalization of \bar{A}_0 . Therefore, for stability of
 303 \bar{A} , we need $\|\bar{A} - \bar{A}_0\| \lesssim p^{-1/2} (1 \wedge \zeta_0^p)$, since $r \leq p$. Note that if \bar{A}_0 is diagonalizable, $r = 1$
 304 implies that we can replace the above upper bound by $1 \wedge \zeta_0$. Putting this stability result together
 305 with the estimation error in the previous paragraph, we obtain (11).

306 **Summary of the Proof of Theorem 2.** To establish the estimation rates, we develop multiple
 307 intermediate lemmas quantifying the exact amount of exploration Algorithm 2 performs. First, we
 308 utilize the fact that the bias of the posterior distribution \mathcal{D}_{τ_n} depends on its covariance matrix $\hat{\Sigma}_{\tau_n}$,
 309 as well as a self-normalized continuous-time matrix-valued martingale. For the effect of the former,
 310 i.e., $\bar{\lambda}(\hat{\Sigma}_{\tau_n}^{-1/2})$, we show an upper-bound of the order $\tau_n^{-1/4}$ (Lemma 9). To that end, the local
 311 geometry of the optimality manifolds that contain drift parameters θ that has the same optimal
 312 feedback as that of the unknown truth θ_0 in (6) are fully specified (Lemma 6), and spectral properties
 313 of non-linear functions of random matrices are studied. Then, we establish a stochastic inequality for
 314 the self-normalized martingale, indicating that its scaling is of the order $p(p+q) \log \tau_n$ (Lemma 8).
 315 Therefore, utilizing the fact that $\hat{\theta}_n - \hat{\mu}_{\tau_n}$ has the same scaling as the bias matrix $\hat{\mu}_{\tau_n} - \theta_0$, we
 316 obtain the estimation rates of Theorem 2.

317 Next, to prove the presented regret bound, we establish a delicate and tight analysis for the dominant
 318 effect of the control signal u_t on the regret Algorithm 2 incurs. Technically, by carefully examining
 319 the infinitesimal influences of the control actions at every time on the cost, we show that it suffices
 320 to integrate the squared deviations $\|u_t + Q_u^{-1} \hat{B}_n^\top P(\hat{\theta}_n) x_t\|^2$ to obtain $\text{Reg}(T)$ (Lemma 7). We
 321 proceed toward specifying the effect of the exploration Algorithm 2 performs on its exploitation
 322 performance by proving the Lipschitz continuity of the solutions of the Riccati equation (5) with
 323 respect to the drift parameters: $\|P(\hat{\theta}_n) - P(\theta_0)\| \lesssim \|\hat{\theta}_n - \theta_0\|$ (Lemma 12). This result is a very
 324 important property of (5) that lets the rates of deviations from the optimal action scale the same as the
 325 estimation error, and is proven by careful analysis of integration along matrix-valued curves in the
 326 space of drift matrices, as well as spectral analysis for approximate solutions of a Lyapunov equation
 327 (Lemma 10). Thus, the regret bound is achieved, using the estimation error result in Theorem 2.

328 6 Numerical Analysis

329 We empirically evaluate the theoretical results of Theorems 1 and 2 under three control problems. The
 330 first two are for the flight control of X-29A airplane at 2000 ft [36] and for Boeing 747 [37]. The third

331 simulation is for blood glucose control [41]. We present the results for X-29A airplane in this section,
 332 and defer the other two examples to the appendix. The true drift matrices of the X-29A airplane
 333 are $A_0 = \begin{bmatrix} -0.16 & 0.07 & -1.00 & 0.04 \\ -15.20 & -2.60 & 1.11 & 0.00 \\ 6.84 & -0.10 & -0.06 & 0.00 \\ 0.00 & 1.00 & 0.07 & 0.00 \end{bmatrix}$, $B_0 = \begin{bmatrix} -0.0006 & 0.0007 \\ 1.3430 & 0.2345 \\ 0.0897 & -0.0710 \\ 0.0000 & 0.0000 \end{bmatrix}$. Further, we let $\Sigma_{\mathbb{W}} = 0.5 I_p$,
 334 $Q_x = I_p$, and $Q_u = 0.1 I_q$ where I_n is the n by n identity matrix. To update the diffusion process \mathbf{x}_t
 335 in (1), time-steps of length 10^{-3} are employed. Then, in Algorithm 1, we let $\sigma_w = 5$, $\kappa = \lfloor \tau^{3/2} \rfloor$,
 336 while τ varies from 4 to 20 seconds. The initial feedback K is generated randomly. The results
 337 for 1000 repetitions are depicted on the left plot of Figure 1, confirming Theorem 1 that the failure
 338 probability of stabilization, decreases exponentially in τ .

339 On the right hand side of Figure 1, Algorithm 2 is executed for 600 second, for $\tau_n = 20 \times 1.1^n$. We
 340 compare TS with the *Randomized Estimate* algorithm [2] for 100 different repetitions. Average- and
 341 worst-case values of the estimation error and the regret are reported, both normalized by their scaling
 342 with time and dimension, as in Theorem 2. The graphs show that (especially the worst-case) regret of
 343 TS substantially outperforms, suggesting that TS explores in a more robust fashion. Simulations for
 Boeing 747 and for the blood glucose control, in the appendix, corroborate the above findings.

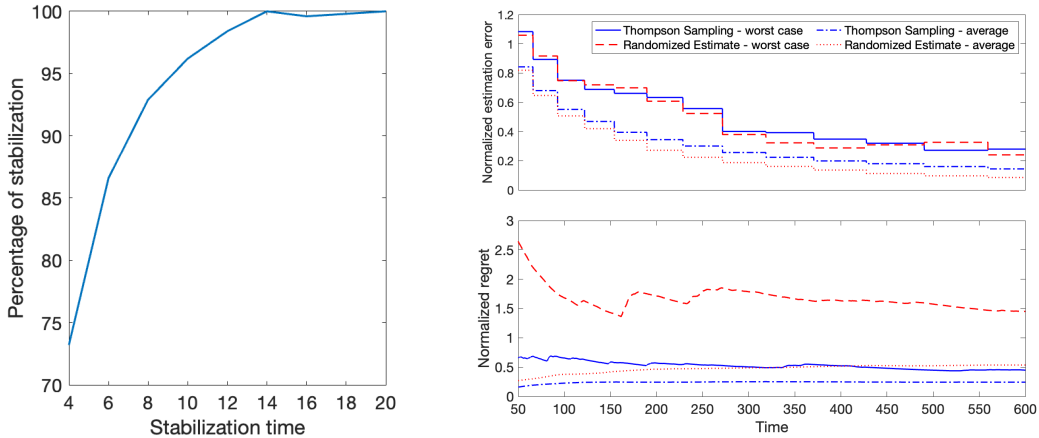


Figure 1: For the X-29A flight control problem, percentage of stabilization for 1000 runs of Algorithm 1 is plotted on the left. The graphs on the right depict the performance of Algorithm 2 (blue) compared to Randomized Estimate policy (red) [2]. The top graph plots the normalized squared estimation error, $\|\hat{\theta}_n - \theta_0\|^2$ divided by $p(p+q)\tau_n^{-1/2} \log \tau_n$, versus time, while the lower one showcases the regret $\text{Reg}(T)$, normalized by $p(p+q)T^{1/2} \log T$. Curves for the worst-case among 100 replications are provided for both quantities, as well as for the averages over all replicates.

344

345 7 Concluding Remarks and Future Work

346 We studied Thompson sampling (TS) RL policies to control a diffusion process with unknown drift
 347 matrices. First, we proposed a stabilization algorithm for linear diffusion processes, and established
 348 that its failure probability decays exponentially with time. Further, efficiency of TS in balancing
 349 exploration versus exploitation for minimizing a quadratic cost function is shown. More precisely,
 350 regret bounds growing as square-root of time and square of dimensions are established for Algorithm 2.
 351 Empirical studies showcasing superiority of TS over state-of-the-art are provided as well.

352 As the first theoretical analysis of TS for control of a continuous-time model, this work implies
 353 multiple important future directions. Establishing minimax regret lower-bounds for diffusion process
 354 control problem is yet unanswered. Moreover, studying the performance of TS for robust control
 355 of the diffusion processes aiming to simultaneously minimize the cost function for a family of drift
 356 matrices, is also an interesting direction for further investigation. Another problem of interest is
 357 efficiency of TS for learning to control under partial observation where the state is not observed and
 358 instead a noisy linear function of the state is available as the output signal.

References

- 359
- 360 [1] W. R. Thompson, “On the likelihood that one unknown probability exceeds another in view of
361 the evidence of two samples,” *Biometrika*, vol. 25, no. 3/4, pp. 285–294, 1933. (Cited on page
362 1)
- 363 [2] M. K. S. Faradonbeh and M. S. S. Faradonbeh, “Efficient estimation and control of unknown
364 stochastic differential equations,” *arXiv preprint arXiv:2109.07630*, 2021. (Cited on page 1, 2,
365 3, 9, 42, 43, 44)
- 366 [3] P. A. Ioannou and J. Sun, *Robust adaptive control*. PTR Prentice-Hall Upper Saddle River, NJ,
367 1996, vol. 1. (Cited on page 2)
- 368 [4] F. L. Lewis, L. Xie, and D. Popa, *Optimal and robust estimation: with an introduction to
369 stochastic control theory*. CRC press, 2017.
- 370 [5] A. Subrahmanyam and G. P. Rao, *Identification of Continuous-time Systems: Linear and Robust
371 Parameter Estimation*. CRC Press, 2019.
- 372 [6] J. Umenberger, M. Ferizbegovic, T. B. Schön, and H. Hjalmarsson, “Robust exploration in
373 linear quadratic reinforcement learning,” *Advances in Neural Information Processing Systems*,
374 vol. 32, 2019. (Cited on page 2)
- 375 [7] S. Agrawal and N. Goyal, “Analysis of thompson sampling for the multi-armed bandit problem,”
376 in *Conference on learning theory*. JMLR Workshop and Conference Proceedings, 2012, pp.
377 39–1. (Cited on page 2)
- 378 [8] —, “Further optimal regret bounds for thompson sampling,” in *Artificial intelligence and
379 statistics*. PMLR, 2013, pp. 99–107.
- 380 [9] A. Gopalan and S. Mannor, “Thompson sampling for learning parameterized markov decision
381 processes,” in *Conference on Learning Theory*. PMLR, 2015, pp. 861–898.
- 382 [10] M. J. Kim, “Thompson sampling for stochastic control: The finite parameter case,” *IEEE
383 Transactions on Automatic Control*, vol. 62, no. 12, pp. 6415–6422, 2017.
- 384 [11] M. Abeille and A. Lazaric, “Linear thompson sampling revisited,” in *Artificial Intelligence and
385 Statistics*. PMLR, 2017, pp. 176–184.
- 386 [12] N. Hamidi and M. Bayati, “On worst-case regret of linear thompson sampling,” *arXiv preprint
387 arXiv:2006.06790*, 2020. (Cited on page 2)
- 388 [13] D. Russo and B. Van Roy, “Learning to optimize via posterior sampling,” *Mathematics of
389 Operations Research*, vol. 39, no. 4, pp. 1221–1243, 2014. (Cited on page 2)
- 390 [14] —, “An information-theoretic analysis of thompson sampling,” *The Journal of Machine
391 Learning Research*, vol. 17, no. 1, pp. 2442–2471, 2016.
- 392 [15] D. Russo, B. Van Roy, A. Kazerouni, I. Osband, and Z. Wen, “A tutorial on thompson sampling,”
393 *arXiv preprint arXiv:1707.02038*, 2017. (Cited on page 2)
- 394 [16] M. Abeille and A. Lazaric, “Improved regret bounds for thompson sampling in linear quadratic
395 control problems,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 1–9.
396 (Cited on page 2)
- 397 [17] Y. Ouyang, M. Gagrani, and R. Jain, “Posterior sampling-based reinforcement learning for
398 control of unknown linear systems,” *IEEE Transactions on Automatic Control*, vol. 65, no. 8,
399 pp. 3600–3607, 2019.
- 400 [18] M. K. S. Faradonbeh, A. Tewari, and G. Michailidis, “On adaptive linear–quadratic regulators,”
401 *Automatica*, vol. 117, p. 108982, 2020.
- 402 [19] S. Sudhakara, A. Mahajan, A. Nayyar, and Y. Ouyang, “Scalable regret for learning to control
403 network-coupled subsystems with unknown dynamics,” *arXiv preprint arXiv:2108.07970*, 2021.
404 (Cited on page 2)

- 405 [20] P. Mandl, “Consistency of estimators in controlled systems,” in *Stochastic Differential Systems*.
406 Springer, 1989, pp. 227–234. (Cited on page 2)
- 407 [21] T. E. Duncan and B. Pasik-Duncan, “Adaptive control of continuous-time linear stochastic
408 systems,” *Mathematics of Control, signals and systems*, vol. 3, no. 1, pp. 45–60, 1990.
- 409 [22] P. Caines, “Continuous time stochastic adaptive control: non-explosion, ε -consistency and
410 stability,” *Systems & control letters*, vol. 19, no. 3, pp. 169–176, 1992.
- 411 [23] T. E. Duncan, L. Guo, and B. Pasik-Duncan, “Adaptive continuous-time linear quadratic
412 gaussian control,” *IEEE Transactions on automatic control*, vol. 44, no. 9, pp. 1653–1662,
413 1999.
- 414 [24] P. E. Caines and D. Levanony, “Stochastic ε -optimal linear quadratic adaptation: An alternating
415 controls policy,” *SIAM Journal on Control and Optimization*, vol. 57, no. 2, pp. 1094–1126,
416 2019. (Cited on page 2)
- 417 [25] S. A. A. Rizvi and Z. Lin, “Output feedback reinforcement learning control for the continuous-
418 time linear quadratic regulator problem,” in *2018 Annual American Control Conference (ACC)*.
419 IEEE, 2018, pp. 3417–3422. (Cited on page 2)
- 420 [26] K. Doya, “Reinforcement learning in continuous time and space,” *Neural computation*, vol. 12,
421 no. 1, pp. 219–245, 2000.
- 422 [27] H. Wang, T. Zariphopoulou, and X. Y. Zhou, “Reinforcement learning in continuous time and
423 space: A stochastic control approach.” *J. Mach. Learn. Res.*, vol. 21, pp. 198–1, 2020. (Cited
424 on page 2)
- 425 [28] Z.-P. Jiang, T. Bian, and W. Gao, “Learning-based control: A tutorial and some recent results,”
426 *Foundations and Trends® in Systems and Control*, vol. 8, no. 3, 2020. (Cited on page 2, 3)
- 427 [29] M. Basei, X. Guo, A. Hu, and Y. Zhang, “Logarithmic regret for episodic continuous-time
428 linear-quadratic reinforcement learning over a finite-time horizon,” *Available at SSRN 3848428*,
429 2021. (Cited on page 2, 3)
- 430 [30] L. Szpruch, T. Treetanhiplot, and Y. Zhang, “Exploration-exploitation trade-off for
431 continuous-time episodic reinforcement learning with linear-convex models,” *arXiv preprint*
432 *arXiv:2112.10264*, 2021. (Cited on page 2, 3)
- 433 [31] I. Karatzas and S. Shreve, *Brownian motion and stochastic calculus*. Springer Science &
434 Business Media, 2012, vol. 113. (Cited on page 3, 4, 16, 17, 19, 20, 21, 30, 32)
- 435 [32] G. Chen, G. Chen, and S.-H. Hsu, *Linear stochastic control systems*. CRC press, 1995, vol. 3.
436 (Cited on page 3, 4)
- 437 [33] J. Yong and X. Y. Zhou, *Stochastic controls: Hamiltonian systems and HJB equations*. Springer
438 Science & Business Media, 1999, vol. 43.
- 439 [34] H. Pham, *Continuous-time stochastic control and optimization with financial applications*.
440 Springer Science & Business Media, 2009, vol. 61.
- 441 [35] S. P. Bhattacharyya and L. H. Keel, *Linear Multivariable Control Systems*. Cambridge
442 University Press, 2022. (Cited on page 3, 4)
- 443 [36] J. T. Bosworth, *Linearized aerodynamic and control law models of the X-29A airplane and*
444 *comparison with flight data*. National Aeronautics and Space Administration, Office of
445 Management . . . , 1992, vol. 4356. (Cited on page 5, 8)
- 446 [37] T. Ishihara, H.-J. Guo, and H. Takeda, “A design of discrete-time integral controllers with
447 computation delays via loop transfer recovery,” *Automatica*, vol. 28, no. 3, pp. 599–603, 1992.
448 (Cited on page 5, 8, 42)
- 449 [38] P. Cheridito, H. M. Soner, and N. Touzi, “Small time path behavior of double stochastic integrals
450 and applications to stochastic control,” *The Annals of Applied Probability*, vol. 15, no. 4, pp.
451 2472–2495, 2005. (Cited on page 6, 8)

- 452 [39] B. Laurent and P. Massart, “Adaptive estimation of a quadratic functional by model selection,”
453 *Annals of Statistics*, pp. 1302–1338, 2000. (Cited on page 8, 17, 20, 21)
- 454 [40] J. A. Tropp, “User-friendly tail bounds for sums of random matrices,” *Foundations of computa-*
455 *tional mathematics*, vol. 12, no. 4, pp. 389–434, 2012. (Cited on page 6, 18)
- 456 [41] T. Zhou, J. L. Dickson, and J. Geoffrey Chase, “Autoregressive modeling of drift and random
457 error to characterize a continuous intravascular glucose monitoring sensor,” *Journal of Diabetes*
458 *Science and Technology*, vol. 12, no. 1, pp. 90–104, 2018. (Cited on page 9, 42)
- 459 [42] P. Hartman and A. Wintner, “The spectra of toeplitz’s matrices,” *American Journal of Mathe-*
460 *matics*, vol. 76, no. 4, pp. 867–882, 1954. (Cited on page 17)
- 461 [43] L. Reichel and L. N. Trefethen, “Eigenvalues and pseudo-eigenvalues of toeplitz matrices,”
462 *Linear algebra and its applications*, vol. 162, pp. 153–185, 1992. (Cited on page 17)
- 463 [44] S. I. Resnick, *Adventures in stochastic processes*. Springer Science & Business Media, 1992.
- 464 [45] P. Billingsley, “Convergence of probability measures,” *INC, New York*, vol. 2, no. 2.4, 1999.
- 465 [46] R. Durrett, *Probability: theory and examples*. Cambridge university press, 2019, vol. 49.
- 466 [47] R. Gondhalekar, E. Dassau, and F. J. Doyle III, “Periodic zone-mpc with asymmetric costs for
467 outpatient-ready safety of an artificial pancreas to treat type 1 diabetes,” *Automatica*, vol. 71, pp.
468 237–246, 2016. (Cited on page 42)

469 **Checklist**

- 470 1. For all authors...
- 471 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
472 contributions and scope? [Yes]
- 473 (b) Did you describe the limitations of your work? [Yes] ; See Sections 3 and 4.
- 474 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 475 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
476 them? [Yes]
- 477 2. If you are including theoretical results...
- 478 (a) Did you state the full set of assumptions of all theoretical results? [Yes] ; See Theo-
479 rems 1 and 2, as well as their discussions.
- 480 (b) Did you include complete proofs of all theoretical results? [Yes] ; Proofs are provided
481 as appendices.
- 482 3. If you ran experiments...
- 483 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
484 perimental results (either in the supplemental material or as a URL)? [Yes] ; See
485 Section 6
- 486 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
487 were chosen)? [N/A]
- 488 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
489 ments multiple times)? [N/A] ; The reported curves reflect the worst-case analysis and
490 so no error bars are needed.
- 491 (d) Did you include the total amount of compute and the type of resources used (e.g., type
492 of GPUs, internal cluster, or cloud provider)? [No] .
- 493 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 494 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 495 (b) Did you mention the license of the assets? [N/A]
- 496 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 497
- 498 (d) Did you discuss whether and how consent was obtained from people whose data you're
499 using/curating? [N/A]
- 500 (e) Did you discuss whether the data you are using/curating contains personally identifiable
501 information or offensive content? [N/A]
- 502 5. If you used crowdsourcing or conducted research with human subjects...
- 503 (a) Did you include the full text of instructions given to participants and screenshots, if
504 applicable? [N/A]
- 505 (b) Did you describe any potential participant risks, with links to Institutional Review
506 Board (IRB) approvals, if applicable? [N/A]
- 507 (c) Did you include the estimated hourly wage paid to participants and the total amount
508 spent on participant compensation? [N/A]

509 **Organization of Appendices**

510 This paper has four appendices. First, we prove Theorem 1 in Appendix A, together with multiple
 511 lemmas that the proof of the theorem relies on, and their statements and proofs are provided in
 512 Appendix A as well. Similarly, the proof for Theorem 2 together with intermediate steps, all are
 513 presented in Appendix B. Then, Appendix C consists of statements and proofs of other results that
 514 are used for establishing both theorems. Finally, empirical simulations beyond those presented in
 515 Section 6 are presented in Appendix D.

516 **A Proof of Theorem 1**

517 For analyzing the estimation error, we establish Lemma 4, which under the condition $\kappa \gtrsim \tau^2$
 518 provides that with probability at least $1 - \delta$,

$$\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| \lesssim \frac{p(p+q)^{1/2} \bar{\lambda}(\Sigma_{\mathbb{W}}) \vee \sigma_w^2}{\tau^{1/2} \underline{\lambda}(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2} (1 + \|K\|)^3 \log\left(\frac{pq\kappa}{\delta}\right).$$

519 Note that in the proof of Lemma 4, results of Lemmas 1, 2, and 3 are used.

520 Therefore, Lemma 12 implies that for solutions of (5), with probability at least $1 - \delta$, it holds that

$$\|P(\hat{\boldsymbol{\theta}}) - P(\boldsymbol{\theta}_0)\| \lesssim \frac{p(p+q)^{1/2} \bar{\lambda}(\Sigma_{\mathbb{W}}) \vee \sigma_w^2}{\tau^{1/2} \underline{\lambda}(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2} (1 + \|K\|)^3 \log\left(\frac{pq\kappa}{\delta}\right).$$

Note that we get the same expression as the right-hand-side above, as an upper bound for
 $\|\hat{B}^\top P(\hat{\boldsymbol{\theta}}) - B_0^\top P(\boldsymbol{\theta}_0)\|$. So, letting

$$\bar{A} = A_0 - B_0 Q_u^{-1} \hat{B}^\top P(\hat{\boldsymbol{\theta}}), \quad \bar{A}_0 = A_0 - B_0 Q_u^{-1} B_0^\top P(\boldsymbol{\theta}_0),$$

521 we obtain

$$\|\bar{A} - \bar{A}_0\| \lesssim \sqrt{\frac{p^2 q \bar{\lambda}(\Sigma_{\mathbb{W}}) \vee \sigma_w^2}{\tau \underline{\lambda}(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2}} (1 + \|K\|)^3 \log\left(\frac{pq\kappa}{\delta}\right), \quad (13)$$

522 with probability at least $1 - \delta$.

523 Next, to consider the effect of the above errors on the eigenvalues of \bar{A} , we compare them to that
 524 of \bar{A}_0 . Note that real-parts of all eigenvalues of \bar{A}_0 are at most $-\zeta_0$, as defined in (10). So, using
 525 the result and the notation of Lemma 5, for all eigenvalues of \bar{A} being on the open left half-plane it
 526 suffices to have $\Delta_{\bar{A}_0}(\bar{A} - \bar{A}_0) \leq \zeta_0$. Also, in lights of Lemma 5, suppose that r is the size of the
 527 largest block in the Jordan block-diagonalization of \bar{A}_0 . So, (36) implies that if

$$\|\bar{A} - \bar{A}_0\| \lesssim r^{-1/2} (1 \wedge \zeta_0^r),$$

528 then Algorithm 1 successfully stabilizes the diffusion process in (1). Thus, (13) shows that the failure
 529 probability of Algorithm 1; $\mathbb{P}(\mathcal{E}_\tau)$, satisfies

$$\sqrt{\frac{p^2 q \bar{\lambda}(\Sigma_{\mathbb{W}}) \vee \sigma_w^2}{\tau \underline{\lambda}(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2}} (1 + \|K\|)^3 \log\left(\frac{pq\kappa}{\mathbb{P}(\mathcal{E}_\tau)}\right) \gtrsim r^{-1/2} (1 \wedge \zeta_0^r).$$

530 Finally, $r \leq p$ together with $\log(pq\kappa) \lesssim \tau^{1/2}$, lead to the desired result.

531 In the remainder of this section, technical lemmas above are stated and their proofs will be provided.

532 **A.1 Bounding cross products of state and randomization**

533 **Definition 2** For a set \mathcal{S} , let $\mathbb{1}\{\mathcal{S}\}$ be the indicator function that is 1 on \mathcal{S} , and vanishes outside of \mathcal{S} .

534 **Lemma 1** In Algorithm 1, for $t \geq 0$, define the piecewise-constant signal $v(t)$ below according to
 535 the randomization sequence w_n :

$$v(t) = \sum_{n=0}^{\kappa-1} \mathbb{1}\left\{\frac{n\tau}{\kappa} \leq t < \frac{(n+1)\tau}{\kappa}\right\} w_n. \quad (14)$$

536 Then, with probability at least $1 - \delta$, we have

$$\begin{aligned} & \left\| \int_0^\tau \mathbf{x}_s v(s)^\top \mathbf{d}s - \frac{\tau^2}{2\kappa^2} B_0 \sum_{n=0}^{\kappa-1} w_n w_n^\top \right\| \\ & \lesssim (\sigma_w^2 + \bar{\lambda}(\Sigma_{\mathbb{W}})) \left(1 + \int_0^\tau \|e^{\bar{A}t}\| \mathbf{d}t \right) \left(pq^{1/2} \tau^{1/2} \log \frac{pq}{\delta} + q \left[1 + \frac{\tau^2}{\kappa^2} \right] \frac{\tau}{\kappa^{1/2}} \log^{3/2} \frac{\kappa q}{\delta} \right). \end{aligned}$$

537 Proof. First, after plugging the control signal \mathbf{u}_t in (1) and solving the resulting stochastic differential
538 equation, we obtain

$$\mathbf{x}_t = e^{\bar{A}t} \mathbf{x}_0 + \int_0^t e^{\bar{A}(t-s)} \mathbf{d}\mathbb{W}_s + \int_0^t e^{\bar{A}(t-s)} B_0 v(s) \mathbf{d}s. \quad (15)$$

539 This implies that

$$\int_0^\tau \mathbf{x}_t v(t)^\top \mathbf{d}t = \Phi_1 + \Phi_2 + \Phi_3,$$

540 where

$$\begin{aligned} \Phi_1 &= \int_0^\tau e^{\bar{A}t} \mathbf{x}_0 v(t)^\top \mathbf{d}t = \sum_{n=0}^{\kappa-1} \left(\int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} e^{\bar{A}t} \mathbf{d}t \right) \mathbf{x}_0 w_n^\top, \\ \Phi_2 &= \int_0^\tau \int_0^t e^{\bar{A}(t-s)} \mathbf{d}\mathbb{W}_s v(t)^\top \mathbf{d}t = \sum_{n=0}^{\kappa-1} \left(\int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} \int_0^t e^{\bar{A}(t-s)} \mathbf{d}\mathbb{W}_s \mathbf{d}t \right) w_n^\top, \\ \Phi_3 &= \int_0^\tau \int_0^t e^{\bar{A}(t-s)} B_0 v(s) \mathbf{d}s v(t)^\top \mathbf{d}t. \end{aligned}$$

541 To analyze Φ_1 , we use the fact that every entry of Φ_1 is a normal random variable with mean zero
542 and variance at most

$$\sigma_w^2 \sum_{n=0}^{\kappa-1} \left\| \int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} e^{\bar{A}t} \mathbf{d}t \right\|^2 \|\mathbf{x}_0\|^2 \leq \sigma_w^2 \left(\int_0^\tau \|e^{\bar{A}t}\| \mathbf{d}t \right)^2 \|\mathbf{x}_0\|^2.$$

543 Therefore, with probability at least $1 - \delta$, it holds that

$$\|\Phi_1\| \lesssim \sigma_w \left(\int_0^\tau \|e^{\bar{A}t}\| \mathbf{d}t \right) \|\mathbf{x}_0\| \sqrt{pq \log \left(\frac{pq}{\delta} \right)}. \quad (16)$$

544 Furthermore, to study Φ_2 , Fubini Theorem [31] gives

$$\begin{aligned}
\int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} \int_0^t e^{\bar{A}(t-s)} d\mathbb{W}_s dt &= \int_0^{n\tau\kappa^{-1}} \left(\int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} e^{\bar{A}(t-n\tau\kappa^{-1})} dt \right) e^{\bar{A}(n\tau\kappa^{-1}-s)} d\mathbb{W}_s \\
&+ \int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} \left(\int_s^{(n+1)\tau\kappa^{-1}} e^{\bar{A}(t-s)} dt \right) d\mathbb{W}_s \\
&= F \sum_{m=1}^n e^{\bar{A}(n-m)\tau\kappa^{-1}} \int_{(m-1)\tau\kappa^{-1}}^{m\tau\kappa^{-1}} e^{\bar{A}(m\tau\kappa^{-1}-s)} d\mathbb{W}_s \\
&+ \int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} G_s d\mathbb{W}_s
\end{aligned}$$

545 where the matrix $F = \int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} e^{\bar{A}(t-n\tau\kappa^{-1})} dt$, does not depend on s or n , and the matrix

546 $G_s = \int_s^{(n+1)\tau\kappa^{-1}} e^{\bar{A}(t-s)} dt$, does not depend on n , since $n\tau\kappa^{-1} \leq s \leq (n+1)\tau\kappa^{-1}$.

547 So, letting $e_i, i = 1, \dots, p$, be the standard basis of the Euclidean space, conditioned on the Wiener
548 process $\{\mathbb{W}_s\}_{s \geq 0}$, for every $j = 1, \dots, q$, the coordinate j of $e_i^\top \Phi_2$ is a mean zero normal random
549 variable. Thus, given $\{\mathbb{W}_s\}_{s \geq 0}$, with probability at least $1 - \delta$, it holds that

$$(e_i^\top \Phi_2 e_j)^2 \lesssim \text{var} \left(e_i^\top \Phi_2 e_j \middle| \mathcal{F}(\mathbb{W}_{0:\tau}) \right) \log \frac{1}{\delta}.$$

550 Now, to calculate the conditional variance, we can write

$$\frac{\text{var} \left(e_i^\top \Phi_2 e_j \middle| \mathcal{F}(\mathbb{W}_{0:\tau}) \right)}{\sigma_w^2} = \sum_{n=0}^{\kappa-1} \left[e_i^\top \int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} \int_0^t e^{\bar{A}(t-s)} d\mathbb{W}_s dt \right]^2 \lesssim \sum_{n=1}^{\kappa-1} \left[\left(\sum_{m=1}^n \beta_{m,n} \right)^2 + \alpha_n^2 \right],$$

551 where

$$\begin{aligned}
\beta_{m,n} &= e_i^\top F e^{\bar{A}(n-m)\tau\kappa^{-1}} \int_{(m-1)\tau\kappa^{-1}}^{m\tau\kappa^{-1}} e^{\bar{A}(m\tau\kappa^{-1}-s)} d\mathbb{W}_s, \\
\alpha_n &= e_i^\top \int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} G_s d\mathbb{W}_s.
\end{aligned}$$

552 To proceed, define the matrix $H = [H_{n,m}]$, where for $1 \leq m, n \leq \kappa - 1$, every block $H_{n,m} \in \mathbb{R}^{1 \times p}$
553 is

$$H_{n,m} = e_i^\top F e^{\bar{A}(n-m)\tau\kappa^{-1}},$$

554 for $m \leq n$, and is 0 for $m > n$. Then, denote

$$\Gamma = \begin{bmatrix} \int_{\tau\kappa^{-1}}^{\tau\kappa^{-1}} e^{\bar{A}(\tau\kappa^{-1}-s)} d\mathbb{W}_s \\ 0 \\ \int_{\tau\kappa^{-1}}^{2\tau\kappa^{-1}} e^{\bar{A}(2\tau\kappa^{-1}-s)} d\mathbb{W}_s \\ \vdots \\ \int_{(\kappa-1)\tau\kappa^{-1}}^{\tau} e^{\bar{A}(\tau-s)} d\mathbb{W}_s \end{bmatrix} \in \mathbb{R}^{p(\kappa-1) \times 1},$$

555 to get

$$\sum_{n=0}^{\kappa-1} \left(\sum_{m=1}^n \beta_{m,n}^2 \right) = \|H\Gamma\|^2 \leq \bar{\lambda} (H^\top H) \|\Gamma\|^2.$$

556 Now, for the matrix H , we have [42, 43]:

$$\bar{\lambda} (H^\top H) \lesssim \left(\sum_{n=1}^{\kappa-1} \|H_{n,1}\| \right)^2 \lesssim \left(\tau \kappa^{-1} \sum_{n=1}^{\kappa} e^{\bar{A}n\tau\kappa^{-1}} \right)^2 \lesssim \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right)^2.$$

557 Note that thanks to the independent increments of the Wiener process, the blocks of Γ are statistically
558 independent. Further, by Ito Isometry [31], every block of Γ is a mean-zero normally distributed
559 vector with the covariance matrix

$$\int_0^{\tau\kappa^{-1}} e^{\bar{A}(\tau\kappa^{-1}-s)\Sigma_{\mathbb{W}}} e^{\bar{A}^\top(\tau\kappa^{-1}-s)} ds.$$

560 So, according to the exponential inequalities for quadratic forms of normally distributed random
561 variables [39], it holds with probability at least $1 - \delta$, that

$$\|\Gamma\|^2 \lesssim p\kappa \bar{\lambda}(\Sigma_{\mathbb{W}}) (\tau\kappa^{-1}) \log \frac{1}{\delta}.$$

562 Thus, with probability at least $1 - \delta$, we have

$$\sum_{n=0}^{\kappa-1} \left(\sum_{m=1}^n \beta_{m,n}^2 \right) \lesssim \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right)^2 p\bar{\lambda}(\Sigma_{\mathbb{W}}) \tau \log \frac{1}{\delta}.$$

563 Similarly, the bound above can be shown for $\sum_{n=1}^{\kappa-1} \alpha_n^2$. Hence, we obtain the corresponding high
564 probability bound for a single entry $e_i^\top \Phi_2 e_j$ of Φ_2 , which together with a union bound, implies that

$$\|\Phi_2\| \lesssim \sigma_w p q^{1/2} \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right) \bar{\lambda}(\Sigma_{\mathbb{W}})^{1/2} \tau^{1/2} \log \left(\frac{pq}{\delta} \right), \quad (17)$$

565 with probability at least $1 - \delta$.

566 Next, according to Fubini Theorem, Φ_3 can also be written as

$$\Phi_3 = \int_0^\tau \int_0^s e^{\bar{A}(s-t)} B_0 v(t) v(s)^\top dt ds = \int_0^\tau \int_t^\tau e^{\bar{A}(s-t)} B_0 v(t) v(s)^\top ds dt.$$

567 Thus, we have

$$2\Phi_3 = \int_0^\tau \int_0^\tau e^{\bar{A}|t-s|} B_0 v(t \wedge s) v(s \vee t)^\top dt ds.$$

568 Recall that the signal $v(t)$ in (14) is piecewise-constant, with values determined by the randomization
569 sequence w_n . So, the above double integral can be written as a double sum

$$\begin{aligned} 2\Phi_3 &= \sum_{n=0}^{\kappa-1} \sum_{m=0}^{\kappa-1} \left(\int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} \int_{m\tau\kappa^{-1}}^{(m+1)\tau\kappa^{-1}} e^{\bar{A}|t-s|} ds dt \right) B_0 w_{m \wedge n} w_{m \vee n}^\top \\ &= \sum_{n=0}^{\kappa-1} \sum_{m=0}^{\kappa-1} \left(e^{\bar{A}|m-n|\tau\kappa^{-1}} \int_0^{\tau\kappa^{-1}} \int_0^{\tau\kappa^{-1}} e^{\bar{A}|t-s|} ds dt \right) B_0 w_{m \wedge n} w_{m \vee n}^\top. \end{aligned}$$

570 Thus, we have

$$2\Phi_3 - \frac{\tau^2}{\kappa^2} B_0 \sum_{n=0}^{\kappa-1} w_n w_n^\top = \Phi_4 + \Phi_5, \quad (18)$$

571 for

$$\begin{aligned} \Phi_4 &= \left(\int_0^{\tau\kappa^{-1}} \int_0^{\tau\kappa^{-1}} e^{\bar{A}|t-s|} ds dt - \tau^2 \kappa^{-2} I_q \right) B_0 \sum_{n=0}^{\kappa-1} w_n w_n^\top, \\ \Phi_5 &= 2 \left(\int_0^{\tau\kappa^{-1}} \int_0^{\tau\kappa^{-1}} e^{\bar{A}|t-s|} ds dt \right) \sum_{n=0}^{\kappa-1} \sum_{m=n+1}^{\kappa-1} \left(e^{\bar{A}(m-n)\tau\kappa^{-1}} B_0 w_n w_m^\top \right). \end{aligned}$$

572 To proceed, we use the following concentration inequality for random matrices with martingale
573 difference structures, titled as Matrix Azuma inequality [40].

574 **Theorem 3** Let $\{\Psi_n\}_{n=1}^k$ be a $d_1 \times d_2$ martingale difference sequence. That is, for some filtration
575 $\{\mathcal{F}_n\}_{n=0}^k$, the matrix Ψ_n is \mathcal{F}_n -measurable, and $\mathbb{E}[\Psi_n | \mathcal{F}_{n-1}] = 0$. Suppose that $\|\Psi_n\| \leq \sigma_n$, for
576 some fixed sequence $\{\sigma_n\}_{n=1}^k$. Then, with probability at least $1 - \delta$, we have

$$\left\| \sum_{n=1}^k \Psi_n \right\|^2 \lesssim \left(\sum_{n=1}^k \sigma_n^2 \right) \log \frac{d_1 + d_2}{\delta}.$$

577 So, to study Φ_4 , we apply Theorem 3 to the random matrices $\Psi_n = w_n w_n^\top - \sigma_w^2 I_q$, using the trivial
578 filtration and the high probability upper-bounds for $\|\Psi_n\| \leq \|w_n\|^2 + \sigma_w^2$;

$$\|\Psi_n\| \leq \sigma_n = \sigma_w^2 \left(1 + q \log \frac{q\kappa}{\delta} \right),$$

579 as well as the fact

$$\left\| \int_0^{\tau\kappa^{-1}} \int_0^{\tau\kappa^{-1}} \left(e^{\bar{A}|t-s|} - I_q \right) ds dt \right\| \lesssim \tau^3 \kappa^{-3},$$

580 to obtain the following bound, which holds with probability at least $1 - \delta$:

$$\|\Phi_4\| \lesssim \|B_0\| \sigma_w^2 \tau^3 \kappa^{-2} \left(1 + \frac{q}{\kappa^{1/2}} \log^{3/2} \frac{\kappa q}{\delta} \right). \quad (19)$$

581 On the other hand, to establish an upper-bound for Φ_5 , consider the random matrices

$$\Psi_n = \sum_{m=n+1}^{\kappa-1} \left(e^{\bar{A}(m-n)\tau\kappa^{-1}} B_0 w_n w_m^\top \right),$$

582 subject to the natural filtration they generate, and apply Theorem 3, using the bounds

$$\|\Psi_n\| \leq \sigma_n \lesssim \tau^{-1} \kappa \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right) \|B_0\| \sigma_w^2 q \log \frac{\kappa q}{\delta},$$

583 together with

$$\left\| \int_0^{\tau\kappa^{-1}} \int_0^{\tau\kappa^{-1}} e^{\bar{A}|t-s|} ds dt \right\| \lesssim \tau^2 \kappa^{-2}.$$

584 Therefore, Theorem 3 indicates that with probability at least $1 - \delta$, it holds that

$$\Phi_5 \lesssim \frac{\tau}{\kappa^{1/2}} \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right) \|B_0\| \sigma_w^2 q \log^{3/2} \frac{\kappa q}{\delta}. \quad (20)$$

585 Finally, put (16), (17), (18), (19), and (20) together, to get the desired result.

586

■

587 **A.2 Bounding cross products of state and Wiener process**

588 **Lemma 2** *In Algorithm 1, with probability at least $1 - \delta$, we have*

$$\left\| \int_0^t \mathbf{x}_s d\mathbb{W}_s^\top \right\| \lesssim \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right) (\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2) p(p+q)^{1/2} \tau^{1/2} \log\left(\frac{pq}{\delta}\right).$$

589 *Proof.* First, according to (15), we can write

$$\int_0^\tau \mathbf{x}_t d\mathbb{W}_t^\top = \Phi_1 + \Phi_2 + \Phi_3,$$

590 where

$$\Phi_1 = \int_0^\tau e^{\bar{A}t} \mathbf{x}_0 d\mathbb{W}_t^\top, \quad (21)$$

$$\Phi_2 = \int_0^\tau \int_0^t e^{\bar{A}(t-s)} B_0 v(s) ds d\mathbb{W}_t^\top, \quad (22)$$

$$\Phi_3 = \int_0^\tau \int_0^t e^{\bar{A}(t-s)} d\mathbb{W}_s d\mathbb{W}_t^\top. \quad (23)$$

591 Now, according to Ito Isometry [31], similar to (16), we have

$$\|\Phi_1\| \lesssim \bar{\lambda}(\Sigma_{\mathbb{W}})^{1/2} \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right) \|\mathbf{x}_0\| \sqrt{pq \log\left(\frac{pq}{\delta}\right)}, \quad (24)$$

592 with probability at least $1 - \delta$. Moreover, in a procedure similar to the one that lead to (17), one can
593 show that with probability at least $1 - \delta$, it holds that

$$\|\Phi_2\| \lesssim \left(\int_0^\tau \|e^{\bar{A}t}\| dt \right) \bar{\lambda}(\Sigma_{\mathbb{W}})^{1/2} \sigma_w pq^{1/2} \tau^{1/2} \log\left(\frac{pq}{\delta}\right). \quad (25)$$

594 Therefore, we need to find a similar upper-bound for Φ_3 . To that end, Ito formula provides

$$d\left(e^{-\bar{A}s} \mathbb{W}_s\right) = -\bar{A} e^{-\bar{A}s} \mathbb{W}_s ds + e^{-\bar{A}s} d\mathbb{W}_s.$$

595 Therefore, integration gives

$$\int_0^t e^{-\bar{A}s} d\mathbb{W}_s = e^{-\bar{A}t} \mathbb{W}_t + \bar{A} \int_0^t e^{-\bar{A}s} \mathbb{W}_s ds,$$

596 which after rearranging and letting $\Psi_t = \int_0^t e^{\bar{A}(t-s)} d\mathbb{W}_s$, leads to

$$\Psi_t \mathbb{W}_t^\top = \left(\int_0^t e^{\bar{A}(t-s)} d\mathbb{W}_s \right) \mathbb{W}_t^\top = \mathbb{W}_t \mathbb{W}_t^\top + \bar{A} \left(\int_0^t e^{\bar{A}(t-s)} \mathbb{W}_s ds \right) \mathbb{W}_t^\top.$$

597 Now, since $d\Psi_t = d\mathbb{W}_t$, Ito Isometry [31] implies that $d\Psi_t d\mathbb{W}_t^\top = \Sigma_{\mathbb{W}} dt$. So, apply integration by
 598 part and use the above equation to get

$$\begin{aligned}\Phi_3 &= \int_0^\tau \Psi_t d\mathbb{W}_t^\top = \int_0^\tau d(\Psi_t \mathbb{W}_t^\top) - \left(\int_0^\tau \mathbb{W}_t d\Psi_t^\top \right)^\top - \int_0^\tau d\Psi_t d\mathbb{W}_t^\top \\ &= \Psi_\tau \mathbb{W}_\tau^\top - \left(\int_0^\tau \mathbb{W}_t d\mathbb{W}_t^\top \right)^\top - \Sigma_{\mathbb{W}} \tau \\ &= \mathbb{W}_\tau \mathbb{W}_\tau^\top + \bar{A} \left(\int_0^\tau e^{\bar{A}(\tau-s)} \mathbb{W}_s ds \right) \mathbb{W}_\tau^\top - \left(\int_0^\tau \mathbb{W}_t d\mathbb{W}_t^\top \right)^\top - \Sigma_{\mathbb{W}} \tau.\end{aligned}$$

599 Therefore, every entry of Φ_3 is a quadratic function of the normally distributed random vectors
 600 $\mathbb{W}_\tau, \int_0^\tau e^{\bar{A}(\tau-s)} \mathbb{W}_s ds$. Note that we used the fact that

$$\mathbb{W}_\tau \mathbb{W}_\tau^\top = \int_0^\tau d(\mathbb{W}_t \mathbb{W}_t^\top) = \left(\int_0^\tau \mathbb{W}_t d\mathbb{W}_t^\top \right)^\top + \left(\int_0^\tau \mathbb{W}_t d\mathbb{W}_t^\top \right) + \Sigma_{\mathbb{W}} \tau.$$

601 Thus, exponential inequalities for quadratic forms of normal random vectors [39] imply that for all
 602 $i, j = 1, \dots, p$, it holds that

$$(e_i^\top \Phi_3 e_j)^2 \lesssim p \mathbb{E} \left[(e_i^\top \Phi_3 e_j)^2 \right] \log^2 \frac{1}{\delta}, \quad (26)$$

603 since $\mathbb{E} [e_i^\top \Phi_3 e_j] = 0$. So, it suffices to find the expectation in (26). For that purpose, we use Ito
 604 Isometry [31] to obtain:

$$\begin{aligned}\mathbb{E} \left[(e_i^\top \Phi_3 e_j)^2 \right] &= \mathbb{E} \left[\left(\int_0^\tau e_i^\top \Psi_t e_j^\top \Sigma_{\mathbb{W}}^{1/2} d(\Sigma_{\mathbb{W}}^{-1/2} \mathbb{W}_t) \right)^2 \right] = \mathbb{E} \left[\int_0^\tau \|e_i^\top \Psi_t \Sigma_{\mathbb{W}}^{1/2} e_j\|^2 dt \right] \\ &\leq e_j^\top \Sigma_{\mathbb{W}} e_j \mathbb{E} \left[\int_0^\tau (e_i^\top \Psi_t)^2 dt \right] = e_j^\top \Sigma_{\mathbb{W}} e_j \mathbb{E} \left[\int_0^\tau \left(e_i^\top \int_0^t e^{\bar{A}(t-s)} d\mathbb{W}_s \right)^2 dt \right].\end{aligned}$$

605 To proceed with the above expression, apply Fubini Theorem [31] to interchange the expected value
 606 with the integral, and then use Ito Isometry again:

$$\begin{aligned}\mathbb{E} \left[\int_0^\tau \left(e_i^\top \int_0^t e^{\bar{A}(t-s)} d\mathbb{W}_s \right)^2 dt \right] &= \int_0^\tau \mathbb{E} \left[\left(e_i^\top \int_0^t e^{\bar{A}(t-s)} \Sigma_{\mathbb{W}}^{1/2} d(\Sigma_{\mathbb{W}}^{-1/2} \mathbb{W}_s) \right)^2 \right] dt \\ &= \int_0^\tau e_i^\top \left(\int_0^t e^{\bar{A}(t-s)} \Sigma_{\mathbb{W}} e^{\bar{A}^\top(t-s)} ds \right) e_i dt \\ &\leq e_i^\top \left(\int_0^\tau e^{\bar{A}s} \Sigma_{\mathbb{W}} e^{\bar{A}^\top s} ds \right) e_i \tau.\end{aligned}$$

607 Therefore, (26) yields to

$$\begin{aligned}
\|\Phi_3\|^2 &\leq \sum_{i,j=1}^p (e_i^\top \Phi_3 e_j)^2 \lesssim \sum_{i,j=1}^p \left[e_j^\top \Sigma_{\mathbb{W}} e_j e_i^\top \left(\int_0^\tau e^{\bar{A}s} \Sigma_{\mathbb{W}} e^{\bar{A}^\top s} \mathbf{d}s \right) e_i \right] \tau p \log^2 \frac{p}{\delta} \\
&= \text{tr}(\Sigma_{\mathbb{W}}) \text{tr} \left(\int_0^\tau e^{\bar{A}s} \Sigma_{\mathbb{W}} e^{\bar{A}^\top s} \mathbf{d}s \right) p \tau \log^2 \frac{p}{\delta} \\
&\lesssim \text{tr}(\Sigma_{\mathbb{W}})^2 \left(\int_0^\tau \|e^{\bar{A}s}\| \mathbf{d}s \right)^2 p \tau \log^2 \frac{p}{\delta}. \tag{27}
\end{aligned}$$

608 Finally, putting (24), (25), and (27) together, we obtain the desired result. ■

609

610 A.3 Concentration of normal posterior distribution in Algorithm 1

611 **Lemma 3** *In Algorithm 1, letting $\bar{A} = A_0 + B_0 K$, suppose that*

$$\tau \gtrsim \left(\int_0^\tau \|\exp(\bar{A}s)\|^2 \mathbf{d}s \right) \left(\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2 \|B_0\|^2 \right) (p+q) \log \frac{1}{\delta}, \tag{28}$$

$$\frac{\kappa}{\tau} \gtrsim \frac{\sigma_w^2}{\sigma_w^2 \wedge \underline{\lambda}(\Sigma_{\mathbb{W}})} \|B_0\| (1 \vee \|K\|) q \log \frac{\kappa q}{\delta}. \tag{29}$$

612 Then, for the matrix $\hat{\Sigma}_\tau$ in (8), with probability at least $1 - \delta$ we have

$$\underline{\lambda}(\hat{\Sigma}_\tau) \gtrsim \tau \left(\underline{\lambda}(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2 \right) \left(1 + \|K\|^2 \right)^{-1}.$$

613 **Proof.** First, we can write the control action in (7) as $\mathbf{u}_t = K \mathbf{x}_t + v(t)$, for the piecewise-constant
614 signal $v(t)$ in (14). Then, the dynamics in (1) provides

$$\mathbf{d}\mathbf{x}_t = (\bar{A} \mathbf{x}_t + B_0 v(t)) \mathbf{d}t + \mathbf{d}\mathbb{W}_t.$$

615 Therefore, similar to (15), one can solve the above stochastic differential equation to get

$$\mathbf{x}_t = e^{\bar{A}t} \mathbf{x}_0 + \int_0^t e^{\bar{A}(t-s)} \mathbf{d}\mathbb{W}_s + \int_0^t e^{\bar{A}(t-s)} B_0 v(s) \mathbf{d}s.$$

616 So, using the exponential inequalities for quadratic forms [39], with probability at least $1 - \delta$, it holds
617 that

$$\|\mathbf{x}_\tau - e^{\bar{A}\tau} \mathbf{x}_0\|^2 \lesssim \bar{\lambda} \left(\int_0^\tau e^{\bar{A}s} \Sigma_{\mathbb{W}} e^{\bar{A}^\top s} \mathbf{d}s + \sigma_w^2 \sum_{n=0}^{\kappa-1} J_n B_0 B_0^\top J_n^\top \right) \left(p + p^{1/2} \log \frac{1}{\delta} \right), \tag{30}$$

618 where

$$J_n = \int_{n\tau \kappa^{-1}}^{(n+1)\tau \kappa^{-1}} e^{\bar{A}s} \mathbf{d}s.$$

619 Furthermore, an application of Ito calculus [31] leads to $\mathbf{d}\mathbf{x}_t \mathbf{d}\mathbf{x}_t^\top = \mathbf{d}\mathbb{W}_t \mathbf{d}\mathbb{W}_t^\top = \Sigma_{\mathbb{W}} \mathbf{d}t$. Now, by
620 defining the matrix valued processes

$$\Phi_t = \int_0^t \mathbf{x}_s \mathbf{x}_s^\top \mathbf{d}s, \quad M_t = \int_0^t \mathbf{x}_s \mathbf{d}\mathbb{W}_s^\top + \int_0^t \mathbf{x}_s v(s)^\top B_0^\top \mathbf{d}s,$$

621 we obtain

$$\begin{aligned}
d(\mathbf{x}_t \mathbf{x}_t^\top) &= \mathbf{x}_t d\mathbf{x}_t^\top + d\mathbf{x}_t \mathbf{x}_t^\top + d\mathbf{x}_t d\mathbf{x}_t^\top \\
&= \mathbf{x}_t ((\bar{A}\mathbf{x}_t + B_0 v(t)) dt + d\mathbb{W}_t)^\top \\
&\quad + ((\bar{A}\mathbf{x}_t + B_0 v(t)) dt + d\mathbb{W}_t) \mathbf{x}_t^\top + \Sigma_{\mathbb{W}} dt \\
&= d\Phi_t \bar{A}^\top + \bar{A} d\Phi_t + dM_t + dM_t^\top + \Sigma_{\mathbb{W}} dt.
\end{aligned}$$

622 Thus, after integrating both sides of the above equality, we obtain

$$\Phi_t \bar{A}^\top + \bar{A} \Phi_t + M_t + M_t^\top + t\Sigma_{\mathbb{W}} + \mathbf{x}_0 \mathbf{x}_0^\top - \mathbf{x}_t \mathbf{x}_t^\top = 0.$$

623 Because all eigenvalues of \bar{A} are in the open left half-plane, we can solve the above equation for Φ_t ,
624 to get

$$\Phi_t = \int_0^\infty \exp(\bar{A}s) [M_t + M_t^\top + t\Sigma_{\mathbb{W}} + \mathbf{x}_0 \mathbf{x}_0^\top - \mathbf{x}_t \mathbf{x}_t^\top] \exp(\bar{A}^\top s) ds. \quad (31)$$

625 Next, putting Lemma 1, Lemma 2, and (30) together, as long as (28) holds, with probability at least
626 $1 - \delta$ we have

$$\lambda(M_\tau + M_\tau^\top + \tau\Sigma_{\mathbb{W}} + \mathbf{x}_0 \mathbf{x}_0^\top - \mathbf{x}_\tau \mathbf{x}_\tau^\top) \gtrsim \tau \lambda(\Sigma_{\mathbb{W}}).$$

627 Thus, (31) implies that $\lambda(\Phi_\tau) \gtrsim \tau \lambda(\Sigma_{\mathbb{W}})$. To proceed, consider the matrix $\widehat{\Sigma}_\tau$ in (8), which
628 comprises two signals $\mathbf{x}_t, v(t)$. The empirical covariance matrix of the state signal is studied above,
629 while for the piecewise-constant randomization signal $v(t)$ in (14), we have

$$\int_0^t v(s)v(s)^\top ds = \sum_{n=0}^{\kappa-1} \int_{n\tau\kappa^{-1}}^{(n+1)\tau\kappa^{-1}} w_n w_n^\top ds = \tau \kappa^{-1} \sum_{n=0}^{\kappa-1} w_n w_n^\top.$$

630 Thus, according to Theorem 3, similar to (19) we have

$$\left\| \sum_{n=0}^{\kappa-1} w_n w_n^\top - \kappa \sigma_w^2 I_q \right\| \lesssim \kappa^{1/2} \sigma_w^2 q \log^{3/2} \frac{\kappa q}{\delta},$$

631 with probability at least $1 - \delta$, which for

$$H_\tau = \int_0^\tau \begin{bmatrix} 0_p \\ v(s) \end{bmatrix} \begin{bmatrix} 0_p \\ v(s) \end{bmatrix}^\top ds - \tau \sigma_w^2 \begin{bmatrix} 0_{p \times p} & 0_{p \times q} \\ 0_{q \times p} & I_q \end{bmatrix},$$

632 leads to

$$\|H_\tau\| \lesssim \sigma_w^2 q \log^{3/2} \frac{\kappa q}{\delta}, \quad (32)$$

633 because $\kappa \gtrsim \tau^2$.

634 Next, using $\mathbf{z}_s = [\mathbf{x}_s^\top, \mathbf{x}_s^\top K^\top + v(s)^\top]^\top$, the matrix $\widehat{\Sigma}_\tau$ can be written as

$$\widehat{\Sigma}_\tau = \begin{bmatrix} I_p \\ K \end{bmatrix} \Phi_\tau \begin{bmatrix} I_p \\ K \end{bmatrix}^\top + \tau \sigma_w^2 \begin{bmatrix} 0_{p \times p} & 0_{p \times q} \\ 0_{q \times p} & I_q \end{bmatrix} + F_\tau + H_\tau, \quad (33)$$

635 where

$$F_\tau = \int_0^\tau \left(\begin{bmatrix} I_p \\ K \end{bmatrix} \mathbf{x}_s \begin{bmatrix} 0_p \\ v(s) \end{bmatrix}^\top + \begin{bmatrix} 0_p \\ v(s) \end{bmatrix} \mathbf{x}_s^\top \begin{bmatrix} I_p \\ K \end{bmatrix}^\top \right) ds.$$

636 However, Lemma 1 and $\kappa \gtrsim \tau^2$ give a high probability upper-bound for the above matrix:

$$\|F_\tau\| \lesssim (1 \vee \|K\|) (\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2) \left(pq^{1/2} \tau^{1/2} \log \frac{pq}{\delta} + q \log^{3/2} \frac{\kappa q}{\delta} \right). \quad (34)$$

637 In the sequel, we show that with probability at least $1 - \delta$, it holds that

$$\lambda \left(\begin{bmatrix} I_p \\ K \end{bmatrix} \Phi_\tau \begin{bmatrix} I_p \\ K \end{bmatrix}^\top + \tau \sigma_w^2 \begin{bmatrix} 0_{p \times p} & 0_{p \times q} \\ 0_{q \times p} & I_q \end{bmatrix} \right) \gtrsim \tau (\lambda(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2) (1 + \|K\|^2)^{-1},$$

638 which, according to (32), (33), and (34), implies the desired result. To show the above least eigenvalue
639 inequality, we use $\lambda(\Phi_\tau) \gtrsim \tau \lambda(\Sigma_{\mathbb{W}})$ to obtain

$$\lambda \left(\begin{bmatrix} I_p \\ K \end{bmatrix} \Phi_\tau \begin{bmatrix} I_p \\ K \end{bmatrix}^\top + \tau \sigma_w^2 \begin{bmatrix} 0_{p \times p} & 0_{p \times q} \\ 0_{q \times p} & I_q \end{bmatrix} \right) \gtrsim \tau (\lambda(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2) \lambda \left(\begin{bmatrix} I_p & K^\top \\ K & KK^\top + I_q \end{bmatrix} \right).$$

640 However, block matrix inversion gives

$$\lambda \left(\begin{bmatrix} I_p & K^\top \\ K & KK^\top + I_q \end{bmatrix} \right) = \bar{\lambda} \left(\begin{bmatrix} I_p & K^\top \\ K & KK^\top + I_q \end{bmatrix}^{-1} \right)^{-1} = \bar{\lambda} \left(\begin{bmatrix} K^\top K + I_p & -K^\top \\ -K & I_q \end{bmatrix} \right)^{-1},$$

641 that is clearly at least $(1 + \|K\|^2)^{-1}$, apart from a constant factor. Therefore, we get the desired
642 result. \blacksquare

643 A.4 Approximation of true drift parameter by Algorithm 1

644 **Lemma 4** Suppose that $\hat{\theta}$ is given by Algorithm 1. Then, with probability at least $1 - \delta$, we have

$$\|\hat{\theta} - \theta_0\| \lesssim \frac{p(p+q)^{1/2} \bar{\lambda}(\Sigma_{\mathbb{W}}) \vee \sigma_w^2}{\tau^{1/2} \lambda(\Sigma_{\mathbb{W}}) \wedge \sigma_w^2} (1 + \|K\|)^3 \log \left(\frac{pq\kappa}{\delta} \right). \quad (35)$$

645 *Proof.* First, consider the mean matrix of the Gaussian posterior distribution. Using the data
646 generation mechanism $d\mathbf{x}_t = \theta_0^\top \mathbf{z}_t dt + d\mathbb{W}_t$, we have

$$\hat{\mu}_\tau = \hat{\Sigma}_\tau^{-1} \int_0^\tau \mathbf{z}_s d\mathbf{x}_s^\top = \hat{\Sigma}_\tau^{-1} \left(\int_0^\tau \mathbf{z}_s \mathbf{z}_s^\top ds \theta_0 + \int_0^\tau \mathbf{z}_s d\mathbb{W}_s^\top \right) = \theta_0 - \hat{\Sigma}_\tau^{-1} \left(\theta_0 - \int_0^\tau \mathbf{z}_s d\mathbb{W}_s^\top \right),$$

647 where we used the definition of $\hat{\Sigma}_\tau$ in (8). Now, the sample $\hat{\theta}$ from \mathcal{D}_τ can be written as $\hat{\theta}_\tau =$
648 $\hat{\mu}_\tau + \hat{\Sigma}_\tau^{-1/2} \Phi$, where $\Phi \sim \mathcal{N}(0_{(p+q) \times p}, I_{p+q})$ is a standard normal random matrix, as defined in
649 the notation. So, for the error matrix, it holds that

$$\|\hat{\theta} - \theta_0\| \leq \|\hat{\Sigma}_\tau^{-1}\| \left(\|\theta_0\| + \left\| \int_0^\tau \mathbf{z}_s d\mathbb{W}_s^\top \right\| \right) + \|\hat{\Sigma}_\tau^{-1}\|^{1/2} \|\Phi\|.$$

650 To proceed towards bounding the above error matrix, use

$$\int_0^\tau \mathbf{z}_s d\mathbb{W}_s^\top = \int_0^\tau \left(\begin{bmatrix} I_p \\ K \end{bmatrix} \mathbf{x}_s + \begin{bmatrix} 0 \\ v(s) \end{bmatrix} \right) d\mathbb{W}_s^\top,$$

651 to obtain

$$\left\| \int_0^\tau \mathbf{z}_s d\mathbb{W}_s^\top \right\| \leq (1 \vee \|K\|) \left\| \int_0^\tau \mathbf{x}_s d\mathbb{W}_s^\top \right\| + \left\| \int_0^\tau v(s) d\mathbb{W}_s^\top \right\|,$$

652 To proceed, note that with probability at least $1 - \delta$, we have

$$\|\Phi\|^2 \lesssim p(p+q) \log \frac{p(p+q)}{\delta}.$$

653 Now, by putting this together with the results of Lemma 2, Lemma 3, and (25), we get the desired
654 result. \blacksquare

655

656 **A.5 Eigenvalue ratio bound for sum of two matrices**

657 **Lemma 5** Suppose that M, E are $p \times p$ matrices, and let $M = \Gamma^{-1} \Lambda \Gamma$ be the Jordan diagonalization
 658 of M . So, for some positive integer k , we have $\Lambda \in \mathbb{C}^{p \times p} = \text{diag}(\Lambda_1, \dots, \Lambda_k)$, where the blocks
 659 $\Lambda_1, \dots, \Lambda_k$ are Jordan matrices of the form

$$\Lambda_i = \begin{bmatrix} \lambda_i & 1 & 0 & \cdots & 0 & 0 \\ 0 & \lambda_i & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & \lambda_i & 1 \\ 0 & 0 & 0 & \cdots & 0 & \lambda_i \end{bmatrix} \in \mathbb{C}^{r_i \times r_i}.$$

660 Then, let $\mathbf{r} = \max_{1 \leq i \leq k} r_i \leq p$, and define $\Delta_M(E)$ as the difference between the largest real-part of the
 661 eigenvalues of $M + E$ and that of M . Then, it holds that

$$\Delta_M(E) \leq \left(1 \vee \mathbf{r}^{1/2} \|E\| \text{cond}(\Gamma)\right)^{1/\mathbf{r}}, \quad (36)$$

662 where $\text{cond}(\Gamma)$ is the condition number of Γ : $\text{cond}(\Gamma) = \bar{\lambda}(\Gamma^\top \Gamma)^{1/2} \underline{\lambda}(\Gamma^\top \Gamma)^{-1/2}$.

663 **Proof.** Since the expression on the right-hand-side of (36) is positive, it is enough to consider an
 664 eigenvalue λ of $M + E$ which is not an eigenvalue of M , and show that $\Re(\lambda) - \log \bar{\lambda}(\exp(M))$
 665 is less than the expression on the RHS of (36). So, for such λ , the matrix $M - \lambda I_p$ is non-singular,
 666 while $M + E - \lambda I_p$ is singular. Let the vector $v \neq 0$ be such that $(M + E - \lambda I_p)v = 0$, which by
 667 Jordan diagonalization above implies that

$$v = -\Gamma^{-1}(\Lambda - \lambda I)^{-1} \Gamma E v. \quad (37)$$

668 Then, $\Lambda = \text{diag}(\Lambda_1, \dots, \Lambda_k)$ indicates that $\Lambda - \lambda I$ and $(\Lambda - \lambda I)^{-1}$ are block diagonal, the latter
 669 consisting of the blocks $\text{diag}\left((\Lambda_1 - \lambda I_{r_1})^{-1}, \dots, (\Lambda_k - \lambda I_{r_k})^{-1}\right)$.

670 Now, multiplications show that

$$(\Lambda_i - \lambda I_{r_i})^{-1} = - \begin{bmatrix} (\lambda - \lambda_i)^{-1} & (\lambda - \lambda_i)^{-2} & \cdots & (\lambda - \lambda_i)^{-r_i} \\ 0 & (\lambda - \lambda_i)^{-1} & \cdots & (\lambda - \lambda_i)^{-r_i+1} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & (\lambda - \lambda_i)^{-1} \end{bmatrix}.$$

671 Therefore, according to the definition of matrix operator norms in Section 1, we obtain

$$\left\| (\Lambda_i - \lambda I_{r_i})^{-1} \right\|^2 \leq \mathbf{r} \left(1 \vee |\lambda - \lambda_i|^{-\mathbf{r}}\right)^2.$$

672 Putting these bounds for the blocks of $(\Lambda - \lambda I)^{-1}$ together, (37) leads to

$$\begin{aligned} 1 &\leq \left\| (\Lambda - \lambda I)^{-1} \right\| \|\Gamma\| \|\Gamma^{-1}\| \|E\| \\ &\leq \mathbf{r}^{1/2} \text{cond}(\Gamma) \|E\| \max_{1 \leq i \leq k} \left(1 \wedge |\lambda - \lambda_i|^{\mathbf{r}}\right)^{-1} \\ &\leq \mathbf{r}^{1/2} \text{cond}(\Gamma) \|E\| \left(1 \wedge (\Re(\lambda) - \log \bar{\lambda}(\exp(M)))^{\mathbf{r}}\right)^{-1}. \end{aligned}$$

673 To see the last inequality above, note that if $\Re(\lambda) - \log \bar{\lambda}(\exp(M))$ is positive, then it is larger than
 674 all the terms $|\lambda - \lambda_i|$, for $i = 1, \dots, k$. Thus, for

$$\Re(\lambda) = \log \bar{\lambda}(\exp(M + E)),$$

675 we obtain (36). ■

676 **B Proof of Theorem 2**

677 To establish the rates of exploration Algorithm 2 performs, we utilize Lemma 8, which indicates that

$$\|\hat{\mu}_{\tau_n} - \theta_0\| \lesssim \left\| \hat{\Sigma}_{\tau_n}^{-1/2} \right\| \left\| \hat{\Sigma}_{\tau_n}^{-1/2} \int_0^{\tau_n} \mathbf{x}_t d\mathbb{W}_t^\top \right\| \lesssim \underline{\lambda} \left(\hat{\Sigma}_{\tau_n} \right)^{-1/2} \left(p(p+q) \bar{\lambda}(\Sigma_{\mathbb{W}}) \log \bar{\lambda} \left(\hat{\Sigma}_{\tau_n} \right) \right)^{1/2}.$$

678 Now, (51) gives $\log \bar{\lambda} \left(\hat{\Sigma}_{\tau_n} \right) \lesssim \log \tau_n$, while Lemma 9 provides $\underline{\lambda} \left(\hat{\Sigma}_{\tau_n} \right) \gtrsim \tau_n^{1/2} \underline{\lambda}(\Sigma_{\mathbb{W}})$. More-
679 over, since $\hat{\Sigma}_{\tau_n}^{1/2} \left(\hat{\theta}_n - \hat{\mu}_{\tau_n} \right)$ is a standard normal $(p+q) \times p$ matrix, we have

$$\left\| \hat{\theta}_n - \hat{\mu}_{\tau_n} \right\| \lesssim \tau_n^{-1/4} \underline{\lambda}(\Sigma_{\mathbb{W}})^{-1/2} (p(p+q) \log(pq))^{1/2}.$$

680 Thus, we obtain the desired result for the estimation error.

681 To proceed toward establishing the regret bound, Lemma 7 shows that we need to integrate

682 $\left\| \mathbf{u}_t + Q_u^{-1} \hat{B}_n^\top P \left(\hat{\theta}_n \right) \mathbf{x}_t \right\|^2$ over the stabilized period of Algorithm 2: $\tau_0 \leq t \leq T$:

$$\text{Reg}(T) \lesssim (\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2) \tau_0 + \int_{\tau_0}^T \left\| \mathbf{u}_t + Q_u^{-1} B_0^\top P(\theta_0) \mathbf{x}_t \right\|^2 dt.$$

683 Further, according to (51), for $\tau_{n-1} < T \leq \tau_n$, we have

$$\int_{\tau_0}^T \left\| \mathbf{u}_t + Q_u^{-1} B_0^\top P(\theta_0) \mathbf{x}_t \right\|^2 dt \lesssim \bar{\lambda}(\Sigma_{\mathbb{W}}) \sum_{i=0}^{n-1} (\tau_{i+1} - \tau_i) \left\| K \left(\hat{\theta}_i \right) - K(\theta_0) \right\|^2.$$

On the other hand, Lemma 12 implies that

$$\left\| K \left(\hat{\theta}_i \right) - K(\theta_0) \right\| \lesssim \frac{\|P(\theta_0)\|^3}{\underline{\lambda}(Q_x) \underline{\lambda}(Q_u)^2} \left\| \hat{\theta}_i - \theta_0 \right\|.$$

684 Thus, we have

$$\text{Reg}(T) \lesssim (\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2) \tau_0 + \frac{\bar{\lambda}(\Sigma_{\mathbb{W}})^2}{\underline{\lambda}(\Sigma_{\mathbb{W}})} \frac{\|P(\theta_0)\|^6}{\underline{\lambda}(Q_x)^2 \underline{\lambda}(Q_u)^4} p(p+q) \sum_{i=0}^{n-1} (\tau_{i+1} - \tau_i) \frac{\log \tau_i}{\tau_i^{1/2}}.$$

685 Thus, according to (12), we obtain the desired regret bound result in Theorem 2.

686 **B.1 Geometry of drift parameters and optimal policies**

687 **Lemma 6** For the drift parameter θ_1 , and for $X \in \mathbb{R}^{p \times p}, Y \in \mathbb{R}^{p \times q}$, define

$$\Delta_{\theta_1}(X, Y) = P(\theta_1) Y + \int_0^\infty e^{\bar{A}_1^\top t} \left[M(X, Y)^\top P(\theta_1) + P(\theta_1) M(X, Y) \right] e^{\bar{A}_1 t} B_1 dt,$$

688 where $\bar{A}_1 = A_1 - B_1 Q_u^{-1} B_1^\top P(\theta_1)$ and $M(X, Y) = X - Y Q_u^{-1} B_1^\top P(\theta_1)$. Then, $\Delta_{\theta_1}(X, Y)$ is
689 the directional derivative of $B^\top P(\theta)$ at θ_1 in the direction $[X, Y]$. Importantly, the tangent space
690 of the manifold of matrices $\theta \in \mathbb{R}^{p \times (p+q)}$ that satisfy $B^\top P(\theta) = B_1^\top P(\theta_1)$ at θ_1 contains all
691 matrices X, Y that $\Delta_{\theta_1}(X, Y) = 0$.

692 Proof. First, note that according to the Lipschitz continuity of $P(\theta)$ in Lemma 12, the directional
693 derivative exists and is well-defined, as long as $\|P(\theta_1)\| < \infty$. However, Lemma 11 provides that
694 $P(\theta_1)$ is finite in a neighborhood of θ_0 , and so the required condition holds. Below, we start by
695 establishing the second result to identify the tangent space, and then prove the general result on the
696 directional derivative.

697 To proceed, let $\boldsymbol{\theta} = \boldsymbol{\theta}_1 + \epsilon [X, Y]^\top$ be such that $B^\top P(\boldsymbol{\theta}) = B_1^\top P(\boldsymbol{\theta}_1)$, and denote $K(\boldsymbol{\theta}_1) =$
698 $-Q_u^{-1} B_1^\top P(\boldsymbol{\theta}_1)$. So, the directional derivative of $P(\boldsymbol{\theta}_1)$ along the matrix $[X, Y]^\top$ can be found as
699 follows. First, denoting the closed-loop transition matrix by $\bar{A} = A - BQ_u^{-1}B^\top P(\boldsymbol{\theta})$, since

$$\bar{A}^\top P(\boldsymbol{\theta}) + P(\boldsymbol{\theta}) \bar{A} + Q_x + K(\boldsymbol{\theta})^\top Q_u K(\boldsymbol{\theta}) = 0,$$

700 we have

$$\begin{aligned} & (\bar{A}_1 + \epsilon X + \epsilon Y K(\boldsymbol{\theta}_1))^\top P(\boldsymbol{\theta}) + P(\boldsymbol{\theta}) (\bar{A}_1 + \epsilon X + \epsilon Y K(\boldsymbol{\theta}_1)) \\ &= -Q_x - K(\boldsymbol{\theta}_1)^\top Q_u K(\boldsymbol{\theta}_1) = \bar{A}_1^\top P(\boldsymbol{\theta}_1) + P(\boldsymbol{\theta}_1) \bar{A}_1. \end{aligned}$$

701 For the matrix $E = \lim_{\epsilon \rightarrow 0} \epsilon^{-1} (P(\boldsymbol{\theta}) - P(\boldsymbol{\theta}_1))$, the latter result implies that

$$\bar{A}_1^\top E + E \bar{A}_1 + (X + Y K(\boldsymbol{\theta}_1))^\top P(\boldsymbol{\theta}_1) + P(\boldsymbol{\theta}_1) (X + Y K(\boldsymbol{\theta}_1)) = 0.$$

702 Then, since all eigenvalues of \bar{A}_1 are in the open left half-plane, the above Lyapunov equation for E
703 leads to the integral form

$$E = \int_0^\infty e^{\bar{A}_1^\top t} \left((X + Y K(\boldsymbol{\theta}_1))^\top P(\boldsymbol{\theta}_1) + P(\boldsymbol{\theta}_1) (X + Y K(\boldsymbol{\theta}_1)) \right) e^{\bar{A}_1 t} dt.$$

704 On the other hand, $K(\boldsymbol{\theta}) = -Q_u^{-1} B^\top P(\boldsymbol{\theta})$ gives

$$0 = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} (B^\top P(\boldsymbol{\theta}) - B_1^\top P(\boldsymbol{\theta}_1)) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [(B^\top - B_1^\top) P(\boldsymbol{\theta}) - B_1^\top (P(\boldsymbol{\theta}_1) - P(\boldsymbol{\theta}))],$$

705 which, according to the definitions of $E, M(X, Y)$, implies the desired result about the tangent space
706 of the manifold under consideration.

707 Next, to establish the more general result on the directional derivative, we use the directional derivative
708 of $P(\boldsymbol{\theta})$ in (65):

$$\int_0^\infty e^{\bar{A}_1^\top t} \left(P(\boldsymbol{\theta}_1) [X + Y K(\boldsymbol{\theta}_1)] + [X + Y K(\boldsymbol{\theta}_1)]^\top P(\boldsymbol{\theta}_1) \right) e^{\bar{A}_1 t} dt.$$

709 Finally, since the directional derivative for B^\top is Y , for $B^\top P(\boldsymbol{\theta})$, by the product rule it is $\Delta_{\boldsymbol{\theta}_1}(X, Y)$,
710 which finishes the proof. \blacksquare

711 B.2 Regret bounds in terms of deviations in control actions

712 **Lemma 7** Let \mathbf{u}_t be the action that Algorithm 2 takes at time t . Then, for the regret of Algorithm 2,
713 it holds that

$$\begin{aligned} \text{Reg}(T) &\lesssim (\bar{\lambda}(\Sigma_{\mathbb{W}}) + \sigma_w^2) \tau_0 \|K + Q_u^{-1} B_0^\top P(\boldsymbol{\theta}_0)\|^2 \\ &+ \int_{\tau_0}^T \|\mathbf{u}_t + Q_u^{-1} B_0^\top P(\boldsymbol{\theta}_0) \mathbf{x}_t\|^2 dt + x_T^*{}^\top P(\boldsymbol{\theta}_0) x_T^*, \end{aligned}$$

714 where x_T^* is the terminal state under the optimal trajectory $\boldsymbol{\pi}_{\text{opt}}$ in (6).

715 **Proof.** First, denote the optimal linear feedback of $\boldsymbol{\pi}_{\text{opt}}$ in (6) by $\mathbf{u}_t = K(\boldsymbol{\theta}_0) \mathbf{x}_t$, where
716 $K(\boldsymbol{\theta}_0) = -Q_u^{-1} B_0^\top P(\boldsymbol{\theta}_0)$. According to the episodic structure of Algorithm 2, for $\tau_n \leq t < \tau_{n+1}$,
717 denote

$$K_t = -Q_u^{-1} \hat{B}_n^\top P(\hat{\boldsymbol{\theta}}_n).$$

718 We first consider the regret of Algorithm 2 after finishing stabilization by running Algorithm 1; i.e.,
719 for $\tau_0 \leq t \leq T$. Fix some small $\epsilon > 0$, that we will let decay later. We proceed by finding an
720 approximation of the regret through sampling at times $\tau_0 + k\epsilon$, for non-negative integers k . To do
721 that, denote $N = \lceil (T - \tau_0)/\epsilon \rceil$, and define the sequence of policies $\{\hat{\boldsymbol{\pi}}_i\}_{i=0}^N$ according to

$$\hat{\boldsymbol{\pi}}_i = \begin{cases} \mathbf{u}_t = K_t \mathbf{x}_t & t < \tau_0 + i\epsilon \\ \mathbf{u}_t = K(\boldsymbol{\theta}_0) \mathbf{x}_t & t \geq \tau_0 + i\epsilon \end{cases}.$$

722 That is, the policy $\hat{\pi}_i$ switches to the optimal feedback at time $\tau_0 + i\epsilon$. So, the zeroth policy $\hat{\pi}_0$
723 corresponds to applying the optimal policy π_{opt} after stabilization at time τ_0 , while the last one $\hat{\pi}_N$
724 is nothing but the one in Algorithm 2, that we denote by $\hat{\pi}$, for the sake of brevity. As such, we have
725 $\text{Reg}_{\hat{\pi}_0}(T) = 0$, and the telescopic summation below holds true:

$$\text{Reg}_{\hat{\pi}}(T) = \sum_{i=0}^{N-1} \left(\text{Reg}_{\hat{\pi}_{i+1}}(T) - \text{Reg}_{\hat{\pi}_i}(T) \right). \quad (38)$$

726 Now, to consider the difference $\text{Reg}_{\hat{\pi}_{i+1}}(T) - \text{Reg}_{\hat{\pi}_i}(T)$, for a fixed i in the range $0 \leq i < N$,
727 denote $t_1 = \tau_0 + i\epsilon$ and let $\mathbf{x}_t^{\hat{\pi}_i}, \mathbf{x}_t^{\hat{\pi}_{i+1}}$ be the state trajectories under $\hat{\pi}_i, \hat{\pi}_{i+1}$, respectively. By
728 definition, we have $\mathbf{x}_t^{\hat{\pi}_i} = \mathbf{x}_t^{\hat{\pi}_{i+1}}$, for all $t \leq t_1$. So, we drop the policy superscript and use \mathbf{x}_{t_1} to
729 refer to the states of both of them at time t_1 . Therefore, as long as $t_1 \leq t < t_1 + \epsilon$, similar to (15),
730 the solutions of the stochastic differential equation are

$$\begin{aligned} \mathbf{x}_t^{\hat{\pi}_i} &= e^{\bar{A}_0(t-t_1)} \mathbf{x}_{t_1} + \int_{t_1}^t e^{\bar{A}_0(t-s)} d\mathbb{W}_s, \\ \mathbf{x}_t^{\hat{\pi}_{i+1}} &= e^{\bar{A}(t-t_1)} \mathbf{x}_{t_1} + \int_{t_1}^t e^{\bar{A}(t-s)} d\mathbb{W}_s, \end{aligned}$$

731 where $\bar{A}_0 = A_0 + B_0K(\boldsymbol{\theta}_0)$ and $\bar{A} = A_0 + B_0K_{t_1}$ are the closed-loop transition matrices under $\hat{\pi}_i$
732 and $\hat{\pi}_{i+1}$, respectively. To work with the above two state trajectories, we define some notations for
733 convenience:

$$\begin{aligned} M_0 &= Q_x + K(\boldsymbol{\theta}_0)^\top Q_u K(\boldsymbol{\theta}_0), \\ M_1 &= Q_x + K_{t_1} Q_u K_{t_1}, \\ y_t &= \mathbf{x}_t^{\hat{\pi}_{i+1}} - \mathbf{x}_t^{\hat{\pi}_i}, \\ E_t &= e^{\bar{A}(t-t_1)} - e^{\bar{A}_0(t-t_1)}. \end{aligned}$$

Thus, letting

$$Z_t = \int_{t_1}^t \left[e^{\bar{A}(t-s)} - e^{\bar{A}_0(t-s)} \right] d\mathbb{W}_s,$$

734 it holds that $y_t = E_t \mathbf{x}_{t_1} + Z_t + O(\epsilon^2)$. Further, for the observation signal z_t and the cost matrix Q
735 defined in Section 2, we have

$$\begin{aligned} & \int_{t_1}^{t_1+\epsilon} \left(\left\| Q^{1/2} z_t(\hat{\pi}_{i+1}) \right\|^2 - \left\| Q^{1/2} z_t(\hat{\pi}_i) \right\|^2 \right) dt \\ &= \int_{t_1}^{t_1+\epsilon} \left[\left(\mathbf{x}_t^{\hat{\pi}_i} + y_t \right)^\top M_1 \left(\mathbf{x}_t^{\hat{\pi}_i} + y_t \right) - \mathbf{x}_t^{\hat{\pi}_i}^\top M_0 \mathbf{x}_t^{\hat{\pi}_i} \right] dt \\ &= \int_{t_1}^{t_1+\epsilon} \left[\mathbf{x}_t^{\hat{\pi}_i}^\top S \mathbf{x}_t^{\hat{\pi}_i} + 2y_t^\top M_1 \mathbf{x}_t^{\hat{\pi}_i} + y_t^\top M_1 y_t \right] dt, \end{aligned} \quad (39)$$

736 where $S = M_1 - M_0 = K_{t_1}^\top Q_u K_{t_1} - K(\boldsymbol{\theta}_0)^\top Q_u K(\boldsymbol{\theta}_0)$.

737 On the other hand, for $t \geq t_1 + \epsilon$, the evolutions of the state vectors are the same for the two policies
738 and we have

$$\mathbf{x}_t^{\hat{\pi}_i} = e^{\bar{A}_0(t-t_1-\epsilon)} \mathbf{x}_{t_1+\epsilon}^{\hat{\pi}_i} + \int_{t_1+\epsilon}^t e^{\bar{A}_0(t-s)} d\mathbb{W}_s.$$

739 Therefore, the difference signal becomes

$$y_t = e^{\bar{A}_0(t-t_1+\epsilon)} \left[\mathbf{x}_{t_1+\epsilon}^{\hat{\pi}_{i+1}} - \mathbf{x}_{t_1+\epsilon}^{\hat{\pi}_i} \right] = e^{\bar{A}_0(t-t_1+\epsilon)} y_{t_1+\epsilon} = e^{\bar{A}_0(t-t_1+\epsilon)} [E_{t_1+\epsilon} \mathbf{x}_{t_1} + Z_{t_1+\epsilon}],$$

740 and we obtain

$$\begin{aligned} & \int_{t_1+\epsilon}^T \left(\left\| Q^{1/2} \mathbf{z}_t(\hat{\pi}_{i+1}) \right\|^2 - \left\| Q^{1/2} \mathbf{z}_t(\hat{\pi}_i) \right\|^2 \right) dt \\ &= \int_{t_1+\epsilon}^T \left[\left(\mathbf{x}_t^{\hat{\pi}_i} + y_t \right)^\top M_0 \left(\mathbf{x}_t^{\hat{\pi}_i} + y_t \right) - \mathbf{x}_t^{\hat{\pi}_i \top} M_0 \mathbf{x}_t^{\hat{\pi}_i} \right] dt \\ &= \int_{t_1+\epsilon}^T \left[2y_t^\top M_0 \mathbf{x}_t^{\hat{\pi}_i} + y_t^\top M_0 y_t \right] dt. \end{aligned} \quad (40)$$

741 Now, after doing some algebra, the expressions in (39) and (40) lead to the following for small ϵ :

$$\text{Reg}_{\hat{\pi}_{i+1}}(T) - \text{Reg}_{\hat{\pi}_i}(T) = (\mathbf{x}_{t_1}^\top F_{t_1} \mathbf{x}_{t_1} + 2\mathbf{x}_{t_1}^\top g_{t_1}) \epsilon + O(\epsilon^2),$$

742 where

$$\begin{aligned} F_{t_1} &= S_t + \int_{t_1}^T \left(2H_{t_1}^\top e^{\bar{A}_0^\top(s-t_1)} \left(Q_x + K(\boldsymbol{\theta}_0)^\top Q_u K(\boldsymbol{\theta}_0) \right) e^{\bar{A}_0(s-t_1)} \right) ds + O(\epsilon), \\ g_{t_1} &= \int_{t_1}^T \left(H_{t_1}^\top e^{\bar{A}_0^\top(s-t_1)} \left(Q_x + K(\boldsymbol{\theta}_0)^\top Q_u K(\boldsymbol{\theta}_0) \right) \int_{t_1}^s e^{\bar{A}_0(s-u)} d\mathbb{W}_u \right) ds + O(\epsilon), \\ S_{t_1} &= K_{t_1}^\top Q_u K_{t_1} - K(\boldsymbol{\theta}_0)^\top Q_u K(\boldsymbol{\theta}_0), \\ H_{t_1} &= B_0 (K_{t_1} - K(\boldsymbol{\theta}_0)). \end{aligned}$$

743 Thus, as ϵ tends to zero, by (38), we have

$$\text{Reg}_{\hat{\pi}}(T) - \text{Reg}_{\hat{\pi}}(\tau_0) = \int_{\tau_0}^T (\mathbf{x}_t^\top F_t \mathbf{x}_t + 2\mathbf{x}_t^\top g_t) dt, \quad (41)$$

744 where F_t, g_t are the above expressions, without the $O(\epsilon)$ terms.

745 Next, by (61), the quadratic expression in terms of the matrix F_t can be equivalently written with

$$F_t = S_t + H_t^\top P(\boldsymbol{\theta}_0) + P(\boldsymbol{\theta}_0) H_t - H_t^\top E_t - E_t H_t,$$

746 where

$$E_t = \int_T^\infty e^{\bar{A}_0^\top(s-t)} M_0 e^{\bar{A}_0(s-t)} ds = e^{\bar{A}_0^\top(T-t)} P(\boldsymbol{\theta}_0) e^{\bar{A}_0(T-t)}.$$

747 Note that in the last equality above, we again used (61). Now, after doing some algebra similar to the
748 expression in (59), we have

$$S_t + H_t^\top P(\boldsymbol{\theta}_0) + P(\boldsymbol{\theta}_0) H_t = (K_t - K(\boldsymbol{\theta}_0))^\top Q_u (K_t - K(\boldsymbol{\theta}_0)),$$

749 which in turn implies that

$$\int_{\tau_0}^T \mathbf{x}_t^\top F_t \mathbf{x}_t dt = \int_{\tau_0}^T \left\| Q_u^{1/2} (K_t - K(\boldsymbol{\theta}_0)) \mathbf{x}_t \right\|^2 dt - 2 \int_{\tau_0}^T \mathbf{x}_t^\top E_t H_t \mathbf{x}_t dt. \quad (42)$$

750 To study the latter integral, suppose that x_t^* is the state trajectory under the optimal policy π_{opt}
751 in (6), and define $\xi_t = \mathbf{x}_t - x_t^*$. Note that (1) gives $d\mathbf{x}_t = (\bar{A}_0 + H_t) \mathbf{x}_t dt + d\mathbb{W}_t$, as well as

752 $d\mathbf{x}_t^* = \bar{A}_0 \mathbf{x}_t^* dt + d\mathbb{W}_t$. Thus, we get $d\xi_t = H_t \mathbf{x}_t dt + \bar{A}_0 \xi_t dt$, using which, we have the following
 753 for $\varphi_t = e^{-\bar{A}_0 t} \xi_t$:

$$d\varphi_t = d\left(e^{-\bar{A}_0 t} \xi_t\right) = e^{-\bar{A}_0 t} d\xi_t - \bar{A}_0 e^{-\bar{A}_0 t} \xi_t dt = e^{-\bar{A}_0 t} H_t \mathbf{x}_t dt.$$

754 Above, we used the fact that the matrices $e^{-\bar{A}_0 t}$, \bar{A}_0 commute. So, it holds that

$$\begin{aligned} \mathbf{x}_t^\top E_t H_t \mathbf{x}_t dt &= \mathbf{x}_t^\top e^{\bar{A}_0^\top (T-t)} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t \\ &= \mathbf{x}_t^{*\top} e^{\bar{A}_0^\top (T-t)} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t + \varphi_t^\top e^{\bar{A}_0^\top T} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t \\ &= \mathbf{x}_t^{*\top} e^{\bar{A}_0^\top (T-t)} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t + \frac{1}{2} d\left[\varphi_t^\top e^{\bar{A}_0^\top T} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} \varphi_t\right]. \end{aligned}$$

755 In the above expression, writing the solution of the stochastic differential equation as in (15), we have

$$e^{\bar{A}_0(T-t)} \mathbf{x}_t^* = \mathbf{x}_T^* - \int_t^T e^{\bar{A}_0(T-s)} d\mathbb{W}_s,$$

756 which gives

$$\begin{aligned} 2\mathbf{x}_t^\top E_t H_t \mathbf{x}_t dt &= 2\mathbf{x}_t^{*\top} e^{\bar{A}_0^\top (T-t)} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t + d\left[\varphi_t^\top e^{\bar{A}_0^\top T} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} \varphi_t\right] \\ &= -2\left(\int_t^T e^{\bar{A}_0(T-s)} d\mathbb{W}_s\right)^\top P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t \\ &\quad + 2\mathbf{x}_T^{*\top} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t + d\left[\varphi_t^\top e^{\bar{A}_0^\top T} P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} \varphi_t\right] \\ &= -2\left(\int_t^T e^{\bar{A}_0(T-s)} d\mathbb{W}_s\right)^\top P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t \\ &\quad + d\left[\left(\mathbf{x}_T^* + e^{\bar{A}_0 T} \varphi_t\right)^\top P(\boldsymbol{\theta}_0) \left(\mathbf{x}_T^* + e^{\bar{A}_0 T} \varphi_t\right)\right], \end{aligned}$$

757 where the latest equality holds since the differential of the constant term $\mathbf{x}_T^* P(\boldsymbol{\theta}_0) \mathbf{x}_T^*$ is zero. Next,
 758 integration by part yields to

$$\begin{aligned} \int_{\tau_0}^T \left(\int_t^T e^{\bar{A}_0(T-s)} d\mathbb{W}_s\right)^\top P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} d\varphi_t &= -\left(\int_{\tau_0}^T e^{\bar{A}_0(T-s)} d\mathbb{W}_s\right)^\top P(\boldsymbol{\theta}_0) e^{\bar{A}_0 T} \varphi_{\tau_0} \\ &\quad + \int_{\tau_0}^T \varphi_t^\top e^{\bar{A}_0^\top T} P(\boldsymbol{\theta}_0) e^{\bar{A}_0(T-t)} d\mathbb{W}_t \end{aligned}$$

759 Now, note the following simplifying expressions: First, by definition, we have $\mathbf{x}_T^* + e^{\bar{A}_0 T} \varphi_T =$
 760 $\mathbf{x}_T^* + \xi_T = \mathbf{x}_T$ and

$$\mathbf{x}_T^* + e^{\bar{A}_0 T} \varphi_{\tau_0} = \mathbf{x}_T^* + e^{\bar{A}_0(T-t)} (\mathbf{x}_{\tau_0} - \mathbf{x}_{\tau_0}^*) = e^{\bar{A}_0(T-t)} \mathbf{x}_{\tau_0} + \int_{\tau_0}^T e^{\bar{A}_0(T-s)} d\mathbb{W}_s,$$

761 is the terminal state vector under the policy $\hat{\pi}_0$ that switches to the optimal policy π_{opt} after the time
 762 τ_0 , because $\int_{\tau_0}^T e^{\bar{A}_0(T-s)} d\mathbb{W}_s = \mathbf{x}_T^* - e^{\bar{A}_0(T-\tau_0)} \mathbf{x}_{\tau_0}^*$. Finally, according to Lemma 8, we have

$$\left\| \int_{\tau_0}^T \varphi_t^\top e^{\bar{A}_0^\top T} P(\boldsymbol{\theta}_0) e^{\bar{A}_0(T-t)} d\mathbb{W}_t \right\| \lesssim \left(\int_{\tau_0}^T \left\| e^{\bar{A}_0(T-t)} \xi_t \right\|^2 dt \right)^{1/2} \log \int_{\tau_0}^T \left\| e^{\bar{A}_0(T-t)} \xi_t \right\|^2 dt.$$

763 Putting the above bounds together, we obtain

$$-2 \int_{\tau_0}^T \mathbf{x}_t^\top E_t H_t \mathbf{x}_t dt - x_T^*{}^\top P(\boldsymbol{\theta}_0) x_T^* \lesssim \int_{\tau_0}^T \left\| e^{\bar{A}_0(T-t)}(\mathbf{x}_t) \right\|^2 dt. \quad (43)$$

764 To proceed toward working with the integration of $\mathbf{x}_t^\top g_t$, employ Fubini Theorem [31] to obtain

$$\begin{aligned} \int_{\tau_0}^T \mathbf{x}_t^\top \tilde{g}_t dt &= \int_{\tau_0}^T \int_t^T \int_t^s \left(\mathbf{x}_t^\top H_t^\top e^{\bar{A}_0^\top(s-t)} M e^{\bar{A}_0(s-u)} \right) d\mathbb{W}_u ds dt \\ &= \int_{\tau_0}^T \int_{\tau_0}^u \int_u^T \left(\mathbf{x}_t^\top H_t^\top e^{\bar{A}_0^\top(s-t)} M e^{\bar{A}_0(s-u)} \right) ds dt d\mathbb{W}_u. \end{aligned}$$

765 Now, denote the inner double integral by y_u^\top :

$$y_u^\top = \int_{\tau_0}^u \int_u^T \left(\mathbf{x}_t^\top H_t^\top e^{\bar{A}_0^\top(s-t)} M_0 e^{\bar{A}_0(s-u)} \right) ds dt = \int_0^u \left(\mathbf{x}_t^\top (K_t - K(\boldsymbol{\theta}_0))^\top P_{t,u}^\top \right) dt,$$

766 where

$$P_{t,u}^\top = B_0^\top \int_u^T e^{\bar{A}_0^\top(s-t)} M_0 e^{\bar{A}_0(s-u)} ds.$$

767 Now, let $\beta_T = \int_{\tau_0}^T \|y_u\|^2 du$, and employ Lemma 8 to get

$$\int_{\tau_0}^T \mathbf{x}_t^\top \tilde{g}_t dt = \int_{\tau_0}^T y_u^\top d\mathbb{W}_u = O\left(\beta_T^{1/2} \log^{1/2} \beta_T\right). \quad (44)$$

768 Thus, we can work with β_T to bound the portion of the regret the integral of $\mathbf{x}_t^\top g_t$ captures. For that
769 purpose, the triangle inequality and Fubini Theorem [31] lead to

$$\begin{aligned} \beta_T &\leq \int_0^T \int_0^u \|P_{t,u}(K_t - K(\boldsymbol{\theta}_0)) \mathbf{x}_t\|^2 dt du \\ &= \int_0^T \left(\mathbf{x}_t^\top (K_t - K(\boldsymbol{\theta}_0))^\top \left[\int_t^T P_{t,u}^\top P_{t,u} du \right] (K_t - K(\boldsymbol{\theta}_0)) \mathbf{x}_t \right) dt \\ &\leq \bar{\lambda} \left(\int_t^T P_{t,u}^\top P_{t,u} du \right) \int_0^T \|(K_t - K(\boldsymbol{\theta}_0)) \mathbf{x}_t\|^2 dt. \end{aligned}$$

770 We can show that $\bar{\lambda} \left(\int_t^T P_{t,u}^\top P_{t,u} \mathbf{d}u \right) \lesssim 1$:

$$\begin{aligned}
\bar{\lambda} \left(\int_t^T P_{t,u}^\top P_{t,u} \mathbf{d}u \right) &\leq \int_t^T \|P_{t,u}^\top\|^2 \mathbf{d}u \\
&\lesssim \int_t^T \left\| \int_u^T e^{\bar{A}_0^\top(s-t)} M_0 e^{\bar{A}_0(s-u)} \mathbf{d}s \right\|^2 \mathbf{d}u \\
&\leq \int_t^T \|e^{\bar{A}_0^\top(u-t)}\|^2 \left\| \int_u^T e^{\bar{A}_0^\top(s-u)} M_0 e^{\bar{A}_0(s-u)} \mathbf{d}s \right\|^2 \mathbf{d}u \\
&\leq \|P(\boldsymbol{\theta}_0)\|^2 \int_t^T \|e^{\bar{A}_0^\top(u-t)}\|^2 \mathbf{d}u \lesssim 1.
\end{aligned}$$

771 Above, in the last inequality we use (61). Note that the last expression is a bounded constant, since
772 all eigenvalues of \bar{A}_0 are in the open left half-plane.

773 Thus, according to (44), it is enough to consider

$$\beta_T \lesssim \int_0^T \|(K_t - K(\boldsymbol{\theta}_0)) \mathbf{x}_t\|^2 \mathbf{d}t, \quad (45)$$

774 in order to bound the portion of the regret that the integration of $\mathbf{x}_t^\top g_t$ contributes.

775 While the above discussions apply to the regret during the time interval $\tau_0 \leq t \leq T$, we can similarly
776 bound the regret during the stabilization period $0 \leq t \leq \tau_0$. The difference is in the randomization
777 sequence $w_n, n = 0, 1, \dots$, which is reflected through the piece-wise constant signal $v(t)$ in (14).
778 Therefore, it suffices to add the effect of $v(t)$ to the one of the Wiener process \mathbb{W}_t , and so $\Sigma_{\mathbb{W}}$ will be
779 replaced with $(\Sigma_{\mathbb{W}} + \sigma_w^2)$:

$$\text{Reg}_{\hat{\pi}}(\tau_0) \leq (\Sigma_{\mathbb{W}} + \sigma_w^2) \tau_0 \|K - K(\boldsymbol{\theta}_0)\|^2. \quad (46)$$

780 Finally, (41), (42), (43), (44), (45), and (46) together, we get the desired result. \blacksquare

781 B.3 Stochastic inequality for continuous-time self-normalized martingales

782 **Lemma 8** Let $\mathbf{z}_t = [\mathbf{x}_t^\top, \mathbf{u}_t^\top]^\top$ be the observation signal and $\hat{\Sigma}_t$ be as in (8). Then, for the
783 stochastic integral $\Phi_t = \int_0^t \mathbf{z}_s \mathbf{d}\mathbb{W}_s^\top$, we have

$$\bar{\lambda} \left(\Phi_t^\top \hat{\Sigma}_t^{-1} \Phi_t \right) \lesssim p \bar{\lambda}(\Sigma_{\mathbb{W}}) \left[\log \det \hat{\Sigma}_t - \log \det \hat{\Sigma}_0 \right]. \quad (47)$$

784 *Proof.* We approximate the integrals over the interval $[0, t]$ through n equally distanced points in the
785 interval, and then let $n \rightarrow \infty$. So, let $\epsilon = \lfloor t/n \rfloor$, and for $k = 0, 1, \dots, n-1$, consider the matrix

$$M_k = \frac{1}{\epsilon} \hat{\Sigma}_0 + \sum_{i=0}^k \mathbf{z}_{i\epsilon} \mathbf{z}_{i\epsilon}^\top.$$

786 Using the above matrices, for $k = 1, \dots, n$, define $\alpha_k = \mathbf{z}_{k\epsilon}^\top M_{k-1}^{-1} \mathbf{z}_{k\epsilon}$. Thus, we have

$$\det M_k = \det [M_{k-1} (I + M_{k-1}^{-1} \mathbf{z}_{k\epsilon} \mathbf{z}_{k\epsilon}^\top)] = \det(M_{k-1}) \det(I + M_{k-1}^{-1} \mathbf{z}_{k\epsilon} \mathbf{z}_{k\epsilon}^\top).$$

787 Now, $M_{k-1}^{-1} \mathbf{z}_{k\epsilon} \mathbf{z}_{k\epsilon}^\top$ is a rank-one matrix, and so $p+q-1$ eigenvalues of $I + M_{k-1}^{-1} \mathbf{z}_{k\epsilon} \mathbf{z}_{k\epsilon}^\top$ except
788 one are 1, and one eigenvalue is $1 + \alpha_k$. So, it holds that

$$\frac{\det M_k}{\det M_{k-1}} = 1 + \alpha_k.$$

789 Next, it is straightforward to show that

$$M_k^{-1} = (M_{k-1} + \mathbf{z}_{k\epsilon} \mathbf{z}_{k\epsilon}^\top)^{-1} = M_{k-1}^{-1} - \frac{M_{k-1}^{-1} \mathbf{z}_{k\epsilon} \mathbf{z}_{k\epsilon}^\top M_{k-1}^{-1}}{1 + \mathbf{z}_{k\epsilon}^\top M_{k-1}^{-1} \mathbf{z}_{k\epsilon}}.$$

790 Therefore, we have

$$\mathbf{z}_{k\epsilon}^\top M_k^{-1} \mathbf{z}_{k\epsilon} = \mathbf{z}_{k\epsilon}^\top (M_{k-1} + \mathbf{z}_{k\epsilon} \mathbf{z}_{k\epsilon}^\top)^{-1} \mathbf{z}_{k\epsilon} = \alpha_k - \frac{\alpha_k^2}{1 + \alpha_k} = \frac{\det M_k - \det M_{k-1}}{\det M_k}.$$

791 However, since for all $\alpha \in \mathbb{R}$ we have $1 + \alpha \leq e^\alpha$, we obtain

$$\mathbf{z}_{k\epsilon}^\top M_k^{-1} \mathbf{z}_{k\epsilon} \leq \log \det M_k - \log \det M_{k-1}. \quad (48)$$

792 To proceed, let \mathcal{F}_k be the sigma-field generated by the Wiener process up to time $k\epsilon$:

$$\mathcal{F}_k = \mathcal{F}(\mathbb{W}_s, 0 \leq s \leq k\epsilon).$$

793 Further, define $L_k = \sum_{i=0}^k \mathbf{z}_{i\epsilon} (\mathbb{W}_{(i+1)\epsilon} - \mathbb{W}_{i\epsilon})^\top$.

794 So, we have

$$\mathbb{E}[L_k^\top M_k^{-1} L_k] = \mathbb{E}\left[\mathbb{E}\left[L_k^\top M_k^{-1} L_k \mid \mathcal{F}_k\right]\right] = \mathbb{E}\left[\mathbb{E}\left[\Psi_k^\top M_k^{-1} \Psi_k \mid \mathcal{F}_k\right]\right],$$

795 where $\Psi_k = L_{k-1} + \mathbf{z}_{k\epsilon} (\mathbb{W}_{(k+1)\epsilon} - \mathbb{W}_{k\epsilon})^\top$. Since L_{k-1} is \mathcal{F}_k -measurable, we get

$$\begin{aligned} & \mathbb{E}[L_k^\top M_k^{-1} L_k] \\ &= \mathbb{E}\left[L_{k-1}^\top M_k^{-1} L_{k-1} + \mathbb{E}\left[(\mathbb{W}_{(k+1)\epsilon} - \mathbb{W}_{k\epsilon}) \mathbf{z}_{k\epsilon}^\top M_k^{-1} \mathbf{z}_{k\epsilon} (\mathbb{W}_{(k+1)\epsilon} - \mathbb{W}_{k\epsilon})^\top \mid \mathcal{F}_k\right]\right] \\ &= \mathbb{E}\left[L_{k-1}^\top M_k^{-1} L_{k-1} + (\mathbf{z}_{k\epsilon}^\top M_k^{-1} \mathbf{z}_{k\epsilon}) \epsilon \Sigma_{\mathbb{W}}\right], \end{aligned}$$

796 where in the last line above we used \mathcal{F}_k -measurability of $\mathbf{z}_{k\epsilon}$, M_k , as well as the independent
797 increments property and the covariance matrix of the Wiener process. So, (48) implies that

$$\bar{\lambda}(\mathbb{E}[L_k^\top M_k^{-1} L_k]) - \bar{\lambda}(\mathbb{E}[L_{k-1}^\top M_{k-1}^{-1} L_{k-1}]) \leq \epsilon \bar{\lambda}(\Sigma_{\mathbb{W}}) (\log \det(\epsilon M_k) - \log \det(\epsilon M_{k-1})).$$

798 Thus, summing over $k = 1, \dots, n$, we get

$$\bar{\lambda}\left(\mathbb{E}\left[L_n^\top (\epsilon M_n)^{-1} L_n\right]\right) \leq \bar{\lambda}(\Sigma_{\mathbb{W}}) (\log \det(\epsilon M_n) - \log \det(\epsilon M_0)).$$

799 Now, consider $\bar{\lambda}\left(L_n^\top (\epsilon M_n)^{-1} L_n\right)$. Since $L_n^\top (\epsilon M_n)^{-1} L_n$ is positive semidefinite, its largest
800 eigenvalue can be upper-bounded by its trace, which implies that

$$\begin{aligned} \mathbb{E}\left[\bar{\lambda}\left(L_n^\top (\epsilon M_n)^{-1} L_n\right)\right] &\leq \mathbb{E}\left[\mathbf{tr}\left(L_n^\top (\epsilon M_n)^{-1} L_n\right)\right] \\ &= \mathbf{tr}\left(\mathbb{E}\left[L_n^\top (\epsilon M_n)^{-1} L_n\right]\right) \\ &\leq p \bar{\lambda}\left(\mathbb{E}\left[L_n^\top (\epsilon M_n)^{-1} L_n\right]\right) \\ &\leq p \bar{\lambda}(\Sigma_{\mathbb{W}}) (\log \det(\epsilon M_n) - \log \det(\epsilon M_0)), \end{aligned}$$

801 where we used the fact that the linear operators of trace and expected value interchange.

802 Thus, Martingale Convergence Theorem [31] implies that

$$\bar{\lambda}\left(L_n^\top (\epsilon M_n)^{-1} L_n\right) \lesssim p \bar{\lambda}(\Sigma_{\mathbb{W}}) (\log \det(\epsilon M_n) - \log \det(\epsilon M_0))$$

803 Finally, as n tends to infinity, ϵ shrinks and we obtain the desired result. ■

804 **B.4 Anti-concentration of the posterior precision matrix in Algorithm 2**

805 **Lemma 9** *In Algorithm 2, we have the following for the matrix $\widehat{\Sigma}_{\tau_n}$ that is defined in (8):*

$$\liminf_{n \rightarrow \infty} \tau_n^{-1/2} \underline{\lambda} \left(\widehat{\Sigma}_{\tau_n} \right) \gtrsim \underline{\lambda}(\Sigma_{\mathbb{W}}).$$

806 *Proof.* First, we define some notation. Recall that during the time interval $\tau_i \leq t < \tau_{i+1}$
 807 corresponding to episode i , Algorithm 2 uses a single parameter estimate $\widehat{\theta}_i$. So, for $i = 0, 1, \dots$,
 808 we use Φ_i, K_i, \bar{A}_i to denote the sample covariance matrix of the state vectors of episode i , and the
 809 feedback and closed-loop matrices during episode i :

$$\begin{aligned} \Phi_i &= \int_{\tau_i}^{\tau_{i+1}} \mathbf{x}_t \mathbf{x}_t^\top dt, \\ K_i &= -Q_u^{-1} \widehat{B}_i^\top P \left(\widehat{\theta}_i \right), \\ \bar{A}_i &= A_0 + B_0 K_i. \end{aligned}$$

810 So, it holds that

$$\widehat{\Sigma}_{\tau_n} = \widehat{\Sigma}_{\tau_0} + \sum_{i=0}^{n-1} L_i \Phi_i L_i^\top, \quad (49)$$

811 where $L_i = \begin{bmatrix} I_p \\ K_i \end{bmatrix}$.

812 Now, consider the matrix Φ_i . Note that according to the bounded grows rates of the episode (from
 813 both above and below) in (12), both $\tau_{i+1} - \tau_i$ and τ_i tend to infinity as i grows. Thus, in the sequel,
 814 we suppose that the indices n, i, j, k that are used for denoting the episodes, are large enough. Similar
 815 to (31), we have

$$\Phi_i = \int_0^\infty e^{\bar{A}_i s} \left[(\tau_{i+1} - \tau_i) \Sigma_{\mathbb{W}} + M_i + M_i^\top + \mathbf{x}_{\tau_i} \mathbf{x}_{\tau_i}^\top - \mathbf{x}_{\tau_{i+1}} \mathbf{x}_{\tau_{i+1}}^\top \right] e^{\bar{A}_i^\top s} ds,$$

816 where

$$M_i = \int_{\tau_i}^{\tau_{i+1}} \mathbf{x}_t d\mathbb{W}_t^\top.$$

817 So, using the fact that the real-parts of all eigenvalues of \bar{A}_i are negative and so $\mathbf{x}_{\tau_{i+1}}$ can be bounded
 818 with $\exp(\bar{A}_i(\tau_{i+1} - \tau_i)) \mathbf{x}_{\tau_i}$ similar to (30), as well as Lemma 2, we obtain the following bounds
 819 for the largest and smallest eigenvalues of Φ_i

$$\underline{\lambda}(\Phi_i) \gtrsim (\tau_{i+1} - \tau_i) \underline{\lambda}(\Sigma_{\mathbb{W}}) \underline{\lambda} \left(\int_0^\infty e^{\bar{A}_i s} e^{\bar{A}_i^\top s} ds \right), \quad (50)$$

$$\bar{\lambda}(\Phi_i) \lesssim (\tau_{i+1} - \tau_i) \bar{\lambda}(\Sigma_{\mathbb{W}}) \int_0^\infty \left\| e^{\bar{A}_i s} \right\|^2 ds. \quad (51)$$

820 On the other hand, for the parameter estimates at the end of episodes, similar to (35), we have

$$\widehat{\Sigma}_{\tau_i}^{1/2} \left(\widehat{\theta}_i - \theta_0 \right) = \widehat{\Sigma}_{\tau_i}^{1/2} \left(\widehat{\theta}_i - \widehat{\mu}_{\tau_i} \right) + \widehat{\Sigma}_{\tau_i}^{-1/2} \left(-\theta_0 + \int_0^{\tau_i} \mathbf{z}_s d\mathbb{W}_s^\top \right).$$

821 Note that by the construction of the posterior \mathcal{D}_{τ_i} in (9), for the first term we have $\widehat{\Sigma}_{\tau_i}^{1/2} \left(\widehat{\theta}_i - \widehat{\mu}_{\tau_i} \right) \sim$
 822 $\mathcal{N}(0, I_{p+q})$. Further, for the second term, Lemma 8 together with (51) lead to

$$\left\| \widehat{\Sigma}_{\tau_i}^{-1/2} \left(-\theta_0 + \int_0^{\tau_i} \mathbf{z}_s d\mathbb{W}_s^\top \right) \right\| \lesssim (p+q) \log^{1/2} \tau_i.$$

823 Therefore, we have

$$\left\| \widehat{\Sigma}_{\tau_i}^{1/2} \left(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_0 \right) \right\| \lesssim (p+q) \log^{1/2} \tau_i.$$

824 However, using the relationship between $\widehat{\Sigma}_{\tau_i}$ and $\Phi_0, \dots, \Phi_{i-1}$ in (49), we can write

$$(p+q)^2 \log \tau_i \gtrsim \left(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_0 \right)^\top \widehat{\Sigma}_{\tau_i} \left(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_0 \right) \geq \left(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_0 \right)^\top \left[\sum_{j=0}^{i-1} L_j \Phi_j L_j^\top \right] \left(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_0 \right),$$

825 which according to the bound in (50) implies that

$$\lambda(\Sigma_{\mathbb{W}}) \sum_{j=0}^{i-1} (\tau_{j+1} - \tau_j) \left\| L_j^\top \left(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_0 \right) \right\|^2 \lesssim (p+q)^2 \log \tau_i.$$

826 Clearly, the above result indicates that for $j < i$, it holds that

$$\left\| \left(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_0 \right)^\top L_j \right\|^2 \lesssim \frac{(p+q)^2 \log \tau_i}{\lambda(\Sigma_{\mathbb{W}}) (\tau_{j+1} - \tau_j)}. \quad (52)$$

827 Next, we employ Lemma 6 to study how Algorithm 2 utilizes Thompson sampling to diversify the
828 matrices L_1, L_2, \dots . To do so, we consider the randomization the posterior \mathcal{D}_{τ_i} applies to the
829 sub-matrix of the parameter estimate corresponding to the input matrix \widehat{B}_i . That is, we aim to find

830 the distribution of the random $p \times q$ matrix $\left(\widehat{\boldsymbol{\theta}}_i - \widehat{\boldsymbol{\mu}}_{\tau_i} \right)^\top \begin{bmatrix} 0_{p \times q} \\ I_q \end{bmatrix}$. Since $\widehat{\boldsymbol{\theta}}_i - \widehat{\boldsymbol{\mu}}_{\tau_i} \sim \mathcal{N} \left(0, \widehat{\Sigma}_{\tau_i}^{-1} \right)$,

831 we have

$$E_i = \begin{bmatrix} 0_{p \times q} \\ I_q \end{bmatrix}^\top \left(\widehat{\boldsymbol{\theta}}_i - \widehat{\boldsymbol{\mu}}_{\tau_i} \right) \sim \mathcal{N} \left(0, \begin{bmatrix} 0_{p \times q} \\ I_q \end{bmatrix}^\top \widehat{\Sigma}_{\tau_i}^{-1} \begin{bmatrix} 0_{p \times q} \\ I_q \end{bmatrix} \right) = \mathcal{N} \left(0, \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{22} \right), \quad (53)$$

832 where $\left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{22}$ is the $q \times q$ lower-left block in $\widehat{\Sigma}_{\tau_i}^{-1}$:

$$\widehat{\Sigma}_{\tau_i}^{-1} = \begin{bmatrix} \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{11} & \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{12} \\ \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{21} & \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{22} \end{bmatrix}.$$

833 Note that $\widehat{\Sigma}_{\tau_0}$ is a positive semi-definite matrix. Therefore, it suffices to show the desired result for
834 $\widehat{\Sigma}_{\tau_n} - \widehat{\Sigma}_{\tau_0}$, and so in the sequel we remove the effect of $\widehat{\Sigma}_{\tau_0}$ by treating τ_0 as 0. So, to calculate
835 the inverse $\widehat{\Sigma}_{\tau_i}^{-1}$, we apply block matrix inversion to

$$\widehat{\Sigma}_{\tau_i} = \begin{bmatrix} \left[\widehat{\Sigma}_{\tau_i} \right]_{11} & \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \\ \left[\widehat{\Sigma}_{\tau_i} \right]_{21} & \left[\widehat{\Sigma}_{\tau_i} \right]_{22} \end{bmatrix} = \begin{bmatrix} \sum_{j=0}^{i-1} \Phi_j & \sum_{j=0}^{i-1} \Phi_j K_j^\top \\ \sum_{j=0}^{i-1} K_j \Phi_j & \sum_{j=0}^{i-1} K_j \Phi_j K_j^\top \end{bmatrix},$$

836 to obtain

$$\begin{aligned} \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{11} &= \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} + \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \Omega_i^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{21} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1}, \\ \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{12} &= - \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \Omega_i^{-1}, \\ \left[\widehat{\Sigma}_{\tau_i}^{-1} \right]_{22} &= \Omega_i^{-1}, \\ \Omega_i &= \left[\widehat{\Sigma}_{\tau_i} \right]_{22} - \left[\widehat{\Sigma}_{\tau_i} \right]_{21} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12}. \end{aligned}$$

837 The smallest eigenvalue of $\widehat{\Sigma}_{\tau_i}$ is related to that of Ω_i . On one hand, since Ω_i^{-1} is a sub-matrix
838 of $\widehat{\Sigma}_{\tau_i}^{-1}$; i.e., $\bar{\lambda}(\Omega_i^{-1}) \leq \bar{\lambda}(\widehat{\Sigma}_{\tau_i}^{-1})$, which implies that $\lambda(\Omega_i) \geq \lambda(\widehat{\Sigma}_{\tau_i})$. Now, we show that the

839 inequality holds in the opposite direction as well, modulo a constant factor. Suppose that $\nu \in \mathbb{R}^{p+q}$
840 is a unit vector, $\nu = [\nu_1^\top, \nu_2^\top]^\top$, $\nu_1 \in \mathbb{R}^p$, and $\nu_2 \in \mathbb{R}^q$. So, after doing some algebra as follows, we
841 have

$$\begin{aligned}
\nu^\top \widehat{\Sigma}_{\tau_i} \nu &= \nu_1^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{11} \nu_1 + 2\nu_1^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \nu_2 + \nu_2^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{22} \nu_2 \\
&= \nu_1^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{11} \nu_1 + 2\nu_1^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{11} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \nu_2 \\
&\quad + \nu_2^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{21} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{11} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \nu_2 \\
&\quad + \nu_2^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{22} \nu_2 - \nu_2^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{21} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \nu_2 \\
&= \left(\nu_1 + \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \nu_2 \right)^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{11} \left(\nu_1 + \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \nu_2 \right) \\
&\quad + \nu_2^\top \Omega_i \nu_2.
\end{aligned}$$

842 For the matrix $\left[\widehat{\Sigma}_{\tau_i} \right]_{11} = \sum_{j=0}^{i-1} \Phi_j$, the smallest eigenvalue lower bounds in (50) lead to
843 $\lambda \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11} \right) \gtrsim \tau_i \lambda(\Sigma_{\mathbb{W}})$. Thus, in order to show the desired smallest eigenvalue result for $\widehat{\Sigma}_{\tau_n}$, it
844 suffices to consider unit vectors ν for which $\left\| \nu_1 + \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \nu_2 \right\| \lesssim \tau_i^{-1/4}$ holds. For such
845 unit vectors ν , the expressions $\left[\widehat{\Sigma}_{\tau_i} \right]_{11} = \sum_{j=0}^{i-1} \Phi_j$ and $\left[\widehat{\Sigma}_{\tau_i} \right]_{12} = \sum_{j=0}^{i-1} \Phi_j K_j^\top$, as well as Lemma 11
846 that indicates that the matrices K_j are bounded, $\|\nu_2\|$ needs to be bounded away from zero since
847 $\|\nu_1\|^2 + \|\nu_2\|^2 = \|\nu\|^2 = 1$. Thus, we have

$$\lambda(\Omega_i) \geq \lambda \left(\widehat{\Sigma}_{\tau_i} \right) \gtrsim \lambda(\Omega_i). \quad (54)$$

848 Otherwise, the desired result about the eigenvalue of $\widehat{\Sigma}_{\tau_n}$ holds true.

849 By simplifying the following expression, we get

$$\begin{aligned}
&\sum_{j=0}^{i-1} \left(K_j^\top - \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \right)^\top \Phi_j \left(K_j^\top - \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \right) \\
&= \sum_{j=0}^{i-1} K_j \Phi_j K_j^\top - \sum_{j=0}^{i-1} K_j \Phi_j \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \\
&\quad - \sum_{j=0}^{i-1} \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \right)^\top \Phi_j K_j^\top + \sum_{j=0}^{i-1} \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \right)^\top \Phi_j \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \\
&= \left[\widehat{\Sigma}_{\tau_i} \right]_{22} - \left[\widehat{\Sigma}_{\tau_i} \right]_{21} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} - \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \right)^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{21} \\
&\quad + \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \right)^\top \left[\widehat{\Sigma}_{\tau_i} \right]_{11} \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} \\
&= \Omega_i.
\end{aligned}$$

850 However, we have

$$K_j^\top - \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left[\widehat{\Sigma}_{\tau_i} \right]_{12} = \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11} K_j^\top - \sum_{k=0}^{i-1} \Phi_k K_k^\top \right) = \left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \sum_{k=0}^{i-1} \Phi_k (K_j - K_k)^\top,$$

851 i.e.,

$$\Omega_i = \sum_{j=0}^{i-1} \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \sum_{k=0}^{i-1} \Phi_k (K_j - K_k)^\top \right)^\top \Phi_j \left(\left[\widehat{\Sigma}_{\tau_i} \right]_{11}^{-1} \sum_{k=0}^{i-1} \Phi_k (K_j - K_k)^\top \right). \quad (55)$$

852 We use the above expression to relate the matrices $\Omega_0, \Omega_1, \dots$ to each others. First, let Ψ_0, Ψ_1, \dots
 853 be a sequence of independent random $q \times p$ matrices with standard normal distribution

$$\Psi_i \sim \mathcal{N}(0_{q \times p}, I_q). \quad (56)$$

854 Then, since $[\widehat{\Sigma}_{\tau_i}^{-1}]_{22} = \Omega_i^{-1}$ and (53), we can let $E_i = \Omega_i^{-1/2} \Psi_i$. Further, for $j, k = 0, 1, \dots$,
 855 denote the B -part of the differences $\widehat{\mu}_k - \widehat{\mu}_j$ by

$$H_{kj} = [0_{q \times p}, I_q] (\widehat{\mu}_k - \widehat{\mu}_j).$$

856 Note that the above result together with (53) give

$$[0_{q \times p}, I_q] (\widehat{\theta}_k - \widehat{\theta}_j) = H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j.$$

857 We will show in the sequel that the above normally distributed random matrices are the effective
 858 randomizations that Thompson sampling Algorithm 2 applies for exploration. For that purpose, using
 859 the directional derivatives and the optimality manifolds in Lemma 6, we calculate $K_k - K_j$ according
 860 to $H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j$. Plugging (52) in the expression for $\Delta_{\theta_1}(X, Y)$ in Lemma 6 for

$$[X, Y] = \widehat{\theta}_k^\top - \widehat{\theta}_j^\top,$$

861 we have

$$\left\| \int_0^\infty e^{\bar{A}_j^\top t} \left[L_j^\top (\widehat{\theta}_k - \widehat{\theta}_j) P(\widehat{\theta}_j) + P(\widehat{\theta}_j) (\widehat{\theta}_k - \widehat{\theta}_j)^\top L_j \right] e^{\bar{A}_j t} dt \right\|^2 \lesssim \frac{(p+q)^2 \log \tau_k}{\underline{\lambda}(\Sigma_{\mathbb{W}}) (\tau_{j+1} - \tau_j)},$$

862 and

$$P(\widehat{\theta}_j) Y = P(\widehat{\theta}_j) (\widehat{\theta}_k - \widehat{\theta}_j)^\top \begin{bmatrix} 0_{p \times q} \\ I_q \end{bmatrix} = P(\widehat{\theta}_j) (H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j)^\top.$$

863 Putting the above two portions of $\Delta_{\theta_1}(X, Y)$ together, since Ψ_k, Ψ_j are independent and standard
 864 normal random matrices, (54) implies that the latter portion of $\Delta_{\theta_1}(X, Y)$ is the dominant one. Thus,
 865 according to Lemma 6 and the expression for the optimal feedbacks in (6), we can approximate
 866 $K_k - K_j$ in (55) by

$$-Q_u^{-1} (H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j) P(\widehat{\theta}_j).$$

867 We use the above approximation for the matrix $[\widehat{\Sigma}_{\tau_i}]_{11}^{-1} \sum_{k=0}^{i-1} \Phi_k (K_j - K_k)^\top$ in (55), letting the
 868 episode number i grow. So, the following expression captures the limit behavior of the least eigenvalue
 869 of $\widehat{\Sigma}_{\tau_n}$ in Algorithm 2:

$$\lim_{n \rightarrow \infty} \frac{Q_u \Omega_n Q_u}{\tau_n^{1/2} \underline{\lambda}(\Sigma_{\mathbb{W}})} = \lim_{n \rightarrow \infty} \sum_{j=0}^{n-1} \left(\sum_{k=0}^{n-1} \widetilde{\Phi}_k P(\widehat{\theta}_j) \left(\frac{H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j}{\tau_n^{-1/4}} \right)^\top \right)^\top \frac{\tau_{j+1} - \tau_j}{\tau_n \underline{\lambda}(\Sigma_{\mathbb{W}})} \frac{\Phi_j}{\tau_{j+1} - \tau_j} \left(\sum_{k=0}^{n-1} \widetilde{\Phi}_k P(\widehat{\theta}_j) \left(\frac{H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j}{\tau_n^{-1/4}} \right)^\top \right), \quad (57)$$

870 where

$$\widetilde{\Phi}_k = \left[\sum_{i=0}^{n-1} \Phi_i \right]^{-1} \Phi_k.$$

871 The equation in (57) provides the limit behavior of the randomized exploration Algorithm 2 performs
 872 for learning to control the diffusion process. More precisely, it shows the roles of the random samples
 873 from the posteriors through the random matrices $\Omega_k^{-1/2} \Psi_k$, for $k = 0, \dots, n-1$, which render the
 874 limit matrix in (57) a positive definite one, as describe below.

Note that since $\sum_{i=0}^{n-1} \widetilde{\Phi}_i = I_p$, the expression

$$\sum_{k=0}^{n-1} \widetilde{\Phi}_k P(\widehat{\theta}_j) \left(\frac{H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j}{\tau_n^{-1/4}} \right)^\top$$

875 is a weighted average of the random matrices $\tau_n^{1/4} \left(H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j \right)^\top$. Moreover,
 876 according to the discussions leading to (50) and (51), the matrix $(\tau_{j+1} - \tau_j)^{-1} \Phi_j$ converge as
 877 j grows to a positive definite matrix, for which all eigenvalues are larger than $\underline{\lambda}(\Sigma_{\mathbb{W}})$, modulo
 878 a constant factor. On the other hand, because the lengths of the episodes satisfies the bounded
 879 growth rates in (12), the ratios $\tau_n^{-1}(\tau_{j+1} - \tau_j)$ are bounded from above and below by $\underline{\alpha}^{n-j}$ and
 880 $\overline{\alpha}(\overline{\alpha} + 1)^{n-j-1}$, and their sum over $j = 0, \dots, n-1$ is 1. A similar property of boundedness from
 881 above and below applies to $\tau_j^{-1/4} \tau_n^{-1/4}$. So, the expression on the right-hand-side of (57) is in fact a
 882 weighted average of

$$\sum_{k=0}^{n-1} \widetilde{\Phi}_k P(\widehat{\theta}_j) \left(\frac{H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j}{\tau_n^{-1/4}} \right)^\top,$$

883 for $j = 0, \dots, n-1$.

884 Note that by the distribution of the random matrices in (56), all rows of
 885 $\tau_n^{1/4} \left(H_{kj} + \Omega_k^{-1/2} \Psi_k - \Omega_j^{-1/2} \Psi_j \right)^\top$ are independent normal random vectors, implying
 886 that these random matrices are almost surely full-rank. Therefore, $\tau_n^{-1/2} \Omega_n$ converges to a positive
 887 definite random matrix, which according to (54) implies the desired result.

888 ■

889 C Auxiliary Lemmas

890 C.1 Behaviors of diffusion processes under non-optimal feedback

891 **Lemma 10** Let \hat{A}, \hat{B} be an arbitrary pair of stabilizable system matrices. Suppose that for the
892 closed-loop matrix $\bar{A} = \hat{A} + \hat{B}K$, we have $\bar{\lambda}(\exp(\bar{A})) < 1$, and P satisfies

$$\bar{A}^\top P + P\bar{A} + Q_x + K^\top Q_u K = 0.$$

893 Then, it holds that

$$P = P(\hat{\theta}) + \int_0^\infty e^{\bar{A}^\top t} \left(K + Q_u^{-1} \hat{B}^\top P(\hat{\theta}) \right)^\top \left(Q_u K + \hat{B}^\top P(\hat{\theta}) \right) e^{\bar{A} t} dt.$$

894 Proof. Denote $K(\hat{\theta}) = -Q_u^{-1} \hat{B}^\top P(\hat{\theta})$ and $\hat{\bar{A}} = \hat{A} + \hat{B}K(\hat{\theta})$. So, after doing some algebra, it
895 is easy to show that the algebraic Riccati equation in (5) gives

$$\hat{\bar{A}}^\top P(\hat{\theta}) + P(\hat{\theta}) \hat{\bar{A}} + Q_x + K(\hat{\theta})^\top Q_u K(\hat{\theta}).$$

896 Now, let $\Phi = K^\top Q_u K - K(\hat{\theta})^\top Q_u K(\hat{\theta})$, and subtract the above equation that $P(\hat{\theta})$ solves,
897 from the similar one in the statement of the lemma that P satisfies, to get

$$\left(\bar{A} - \hat{\bar{A}} \right)^\top P(\hat{\theta}) + P(\hat{\theta}) \left(\bar{A} - \hat{\bar{A}} \right) + \bar{A}^\top \left(P - P(\hat{\theta}) \right) + \left(P - P(\hat{\theta}) \right) \bar{A} + \Phi = 0. \quad (58)$$

898 Because $\bar{\lambda}(\exp(\bar{A})) < 1$, by solving (58) for $P - P(\hat{\theta})$, we have

$$P - P(\hat{\theta}) = \int_0^\infty e^{\bar{A}^\top t} \left(\Phi + \left[K - K(\hat{\theta}) \right]^\top \hat{B}^\top P(\hat{\theta}) + P(\hat{\theta}) \hat{B} \left[K - K(\hat{\theta}) \right] \right) e^{\bar{A} t} dt,$$

899 where the fact $\bar{A} - \hat{\bar{A}} = \hat{B} \left[K - K(\hat{\theta}) \right]$ is used above. Then, using $\hat{B}^\top P(\hat{\theta}) = -Q_u K(\hat{\theta})$, it
900 is straightforward to see

$$\begin{aligned} & \left(K - K(\hat{\theta}) \right)^\top Q_u \left(K - K(\hat{\theta}) \right) \\ &= \Phi + \left[K - K(\hat{\theta}) \right]^\top \hat{B}^\top P(\hat{\theta}) + P(\hat{\theta}) \hat{B} \left[K - K(\hat{\theta}) \right], \end{aligned} \quad (59)$$

901 which leads to the desired result. ■

902 C.2 Behaviors of diffusion processes in a neighborhood of the truth

903 **Lemma 11** Letting ζ_0 be as defined in (10), assume that

$$\left\| \hat{\theta} - \theta_0 \right\| \lesssim \frac{\underline{\lambda}(Q_u)}{\|B_0\| \|P(\theta_0)\|} \left(\frac{[\zeta_0 \wedge 1]^p}{p^{1/2}} \wedge \frac{\underline{\lambda}(Q_x)}{\|P(\theta_0)\|} \right). \quad (60)$$

904 Then, for the Riccati equation in (5) which is denoted by $P(\theta)$, we have $\|P(\hat{\theta})\| \lesssim \|P(\theta_0)\|$. Fur-
905 thermore, for any eigenvalue λ of $\hat{A} - \hat{B}Q_u^{-1} \hat{B}^\top P(\hat{\theta})$, it holds that $\Re(\lambda) \lesssim -\underline{\lambda}(Q_x) \|P(\theta_0)\|^{-1}$.

906 Proof. First, let us write $\bar{A}_0 = A_0 - B_0 Q_u^{-1} B_0^\top P(\theta_0)$ and

$$\bar{A}_1 = \hat{A} - \hat{B} Q_u^{-1} B_0^\top P(\theta_0) = \bar{A}_0 + E_1,$$

907 where $E_1 = \hat{A} - A_0 - \left(\hat{B} - B_0 \right) Q_u^{-1} B_0^\top P(\theta_0)$. Since $r \leq p$, (60) implies that E_1 satisfies

$$\|E_1\| \lesssim r^{-1/2} [\zeta_0 \wedge 1]^r.$$

908 So, letting $M = \bar{A}_0$ in (36), Lemma 5 leads to the fact that all eigenvalues of $\exp(\bar{A}_1)$ are inside
 909 the unit-circle. Therefore, all eigenvalues of \bar{A}_1 are on the open left half-plane of the complex plane.
 910 Now, in Lemma 10, let $K = -Q_u^{-1}B_0^\top P(\theta_0)$ and $\bar{A} = \bar{A}_1$, to obtain the matrix denoted by P in
 911 the lemma. Since P satisfies

$$Q_x + K^\top Q_u K = -\bar{A}_1^\top P - P\bar{A}_1 = -\bar{A}_0^\top P - P\bar{A}_0 - E_1^\top P - PE_1,$$

912 writing Lemma 10 for $\bar{A} = \bar{A}_0$, but replacing Q_x with $Q_x + E_1^\top P + PE_1$, we have

$$P = \int_0^\infty e^{\bar{A}_0^\top t} [Q_x + K^\top Q_u K + E_1^\top P + PE_1] e^{\bar{A}_0 t} dt.$$

913 However, according to (5), we have

$$P(\theta_0) = \int_0^\infty e^{\bar{A}_0^\top t} [Q_x + K^\top Q_u K] e^{\bar{A}_0 t} dt. \quad (61)$$

914 Thus, it holds that

$$P = P(\theta_0) + \int_0^\infty e^{\bar{A}_0^\top t} [E_1^\top P + PE_1] e^{\bar{A}_0 t} dt,$$

915 which leads to

$$\|P\| \leq \|P(\theta_0)\| + 2\|E_1\| \|P\| \int_0^\infty \|e^{\bar{A}_0 t}\|^2 dt.$$

916 We will shortly show that $2\|E_1\| \int_0^\infty \|e^{\bar{A}_0 t}\|^2 dt < 1$. So, by Lemma 10, we have $\|P(\hat{\theta})\| \leq \|P\| \lesssim$
 917 $\|P(\theta_0)\|$, which is the desired result.

918 To proceed, denote the closed-loop matrix by $\hat{A} = \hat{A} - \hat{B}Q_u^{-1}\hat{B}^\top P(\hat{\theta})$, and let the p dimensional
 919 unit vector ν attain the maximum of $\|\exp(\hat{A})\nu\|$, i.e., $\|\exp(\hat{A})\nu\| = \|\exp(\hat{A})\|$. Then, (61) for $\hat{\theta}$
 920 (instead of θ_0) implies that

$$\|P(\hat{\theta})\| \geq \nu^\top P(\hat{\theta})\nu = \int_0^\infty \nu^\top e^{\hat{A}^\top t} \left[Q_x + P(\hat{\theta})^\top \hat{B}Q_u^{-1}\hat{B}^\top P(\hat{\theta}) \right] e^{\hat{A} t} \nu dt.$$

921 Therefore, $\lambda \left(Q_x + P(\hat{\theta})^\top \hat{B}Q_u^{-1}\hat{B}^\top P(\hat{\theta}) \right) \geq \lambda(Q_x)$, together with the fact that the magni-
 922 tudes of all eigenvalues are smaller than the operator norm, imply that for an arbitrary eigenvalue λ
 923 of \hat{A} , we have

$$\|P(\theta_0)\| \gtrsim \|P(\hat{\theta})\| \geq \lambda(Q_x) \int_0^\infty e^{2\Re(\lambda)t} dt, \quad (62)$$

924 which leads to the second desired result of the lemma. To complete the proof, we need to establish
 925 that $\|E_1\| \int_0^\infty \|e^{\bar{A}_0 t}\|^2 dt < 1/2$. For that purpose, if we write (62) for θ_0 instead of $\hat{\theta}$, the condition
 926 in (60) implies the above bound. ■

927 C.3 Perturbation analysis for algebraic Riccati equation in (5)

928 **Lemma 12** Assume that (60) holds. Then, we have

$$\|P(\hat{\theta}) - P(\theta_0)\| \lesssim \frac{\|P(\theta_0)\|^2}{\lambda(Q_x)} (1 \vee \|Q_u^{-1}B_0^\top P(\theta_0)\|) \|\hat{\theta} - \theta_0\|.$$

929 Proof. First, fix the dynamics matrix $\widehat{\theta}$, and let \mathcal{C} be a linear segment connecting θ_0 and $\widehat{\theta}$:

$$\mathcal{C} = \left\{ (1 - \alpha)\theta_0 + \alpha\widehat{\theta} \right\}_{0 \leq \alpha \leq 1}.$$

930 Let $\theta_1 \in \mathcal{C}$ be arbitrary. Then, the derivative of $P(\theta)$ at θ_1 in the direction of \mathcal{C} can be found by
 931 using the difference matrices $E_A = \widehat{A} - A_0$, $E_B = \widehat{B} - B_0$. Denote $E = [E_A, E_B]^\top$. Then, we
 932 find $P(\theta_2)$, where $\theta_2 = \theta_1 + \epsilon E$, for an infinitesimal value of ϵ . So, we have

$$\begin{aligned} & P(\theta_1) B_2 Q_u^{-1} B_2^\top P(\theta_1) \\ &= \epsilon P(\theta_1) E_B Q_u^{-1} B_2^\top P(\theta_1) + P(\theta_1) B_1 Q_u^{-1} B_2^\top P(\theta_1) \\ &= O(\epsilon^2) + \epsilon P(\theta_1) E_B Q_u^{-1} B_1^\top P(\theta_1) + \epsilon P(\theta_1) B_1 Q_u^{-1} E_B^\top P(\theta_1) + P(\theta_1) B_1 Q_u^{-1} B_1^\top P(\theta_1). \end{aligned}$$

933 Therefore, we can calculate $P(\theta_2) B_2 Q_u^{-1} B_2^\top P(\theta_2)$. To that end, let $P = P(\theta_2) - P(\theta_1)$, write
 934 $P(\theta_2)$ in terms of $P, P(\theta_1)$, and use the above result to get

$$\begin{aligned} & P(\theta_2) B_2 Q_u^{-1} B_2^\top P(\theta_2) \\ &= P(\theta_2) B_2 Q_u^{-1} B_2^\top P + P(\theta_2) B_2 Q_u^{-1} B_2^\top P(\theta_1) \\ &= O(\|P\|^2) + P(\theta_1) B_2 Q_u^{-1} B_2^\top P + P B_2 Q_u^{-1} B_2^\top P(\theta_1) + P(\theta_1) B_2 Q_u^{-1} B_2^\top P(\theta_1) \\ &= O(\|P\|^2) + P(\theta_1) B_2 Q_u^{-1} B_2^\top P + P B_2 Q_u^{-1} B_2^\top P(\theta_1) \\ &+ O(\epsilon^2) + \epsilon P(\theta_1) E_B Q_u^{-1} B_1^\top P(\theta_1) + \epsilon P(\theta_1) B_1 Q_u^{-1} E_B^\top P(\theta_1) \\ &+ P(\theta_1) B_1 Q_u^{-1} B_1^\top P(\theta_1). \end{aligned} \tag{63}$$

935 Again, expanding $A_2 = A_1 + E_A$ and $P(\theta_2) = P(\theta_1) + P$, it yields to

$$\begin{aligned} & A_2^\top P(\theta_2) + P(\theta_2) A_2 \\ &= A_2^\top P(\theta_1) + A_2^\top P + P(\theta_1) A_2 + P A_2 \\ &= A_1^\top P(\theta_1) + \epsilon E_A^\top P(\theta_1) + A_2^\top P \\ &+ P(\theta_1) A_1 + \epsilon P(\theta_1) E_A + P A_2. \end{aligned}$$

936 To proceed, plug in the continuous-time algebraic Riccati equation in (5) for θ_1, θ_2 below in the
 937 above expression:

$$\begin{aligned} A_2^\top P(\theta_2) + P(\theta_2) A_2 &= P(\theta_2) B_2 Q_u^{-1} B_2^\top P(\theta_2) + Q_x, \\ A_1^\top P(\theta_1) + P(\theta_1) A_1 &= P(\theta_1) B_1 Q_u^{-1} B_1^\top P(\theta_1) + Q_x. \end{aligned}$$

938 So, we obtain

$$\begin{aligned} & A_2^\top P(\theta_2) + P(\theta_2) A_2 - A_1^\top P(\theta_1) - P(\theta_1) A_1 \\ &= \epsilon E_A^\top P(\theta_1) + \epsilon P(\theta_1) E_A + P A_2 + A_2^\top P \\ &= P(\theta_2) B_2 Q_u^{-1} B_2^\top P(\theta_2) - P(\theta_1) B_1 Q_u^{-1} B_1^\top P(\theta_1) \\ &= O(\|P\|^2) + P(\theta_1) B_2 Q_u^{-1} B_2^\top P + P B_2 Q_u^{-1} B_2^\top P(\theta_1) \\ &+ O(\epsilon^2) + \epsilon P(\theta_1) E_B Q_u^{-1} B_1^\top P(\theta_1) + \epsilon P(\theta_1) B_1 Q_u^{-1} E_B^\top P(\theta_1), \end{aligned}$$

939 where in the last equality above, we used (63). Now, rearrange the terms in the above statement to
 940 get an equation that does not contain any expression in term of θ_2 . So, it becomes

$$\begin{aligned} 0 &= [A_2^\top - P(\theta_1) B_2 Q_u^{-1} B_2^\top] P + P [A_2 - B_2 Q_u^{-1} B_2^\top P(\theta_1)] - O(\|P\|^2) - O(\epsilon^2) \\ &+ \epsilon E_A^\top P(\theta_1) + \epsilon P(\theta_1) E_A - \epsilon P(\theta_1) E_B Q_u^{-1} B_1^\top P(\theta_1) - \epsilon P(\theta_1) B_1 Q_u^{-1} E_B^\top P(\theta_1). \end{aligned}$$

941 Next, to simplify the above equality, define the followings:

$$\begin{aligned} D &= A_2 - B_2 Q_u^{-1} B_2^\top P(\theta_1), \\ K(\theta_1) &= -Q_u^{-1} B_1^\top P(\theta_1), \\ R &= \epsilon P(\theta_1) [E_A + E_B K(\theta_1)] + \epsilon [K(\theta_1)^\top E_B^\top + E_A^\top] P(\theta_1) - O(\epsilon^2). \end{aligned}$$

942 So, writing our equation in terms of $D, K(\boldsymbol{\theta}_1), R$, it gives

$$0 = D^\top P + PD - O(\|P\|^2) + R. \quad (64)$$

943 The discussion after (6) states that all eigenvalues of $\bar{A}_1 = A_1 - B_1 Q_u^{-1} B_1^\top P(\boldsymbol{\theta}_1)$ lie in the open left
 944 half-plane. Therefore, if ϵ is small enough, real-parts of all eigenvalues of D are negative, according
 945 to Lemma 5. Therefore, (64) implies that

$$\bar{\lambda}(P) \leq \bar{\lambda} \left(\int_0^\infty e^{D^\top t} R e^{Dt} dt \right) \leq \|R\| \int_0^\infty \|e^{Dt}\|^2 dt.$$

946 So, as ϵ decays, R vanishes, which by the above inequality shows that P shrinks as ϵ tends to zero.
 947 Further, as ϵ decays, D converges to \bar{A}_1 . Thus, by (64), we have

$$\lim_{\epsilon \rightarrow 0} \epsilon^{-1} P = \int_0^\infty e^{\bar{A}_1^\top t} \left(P(\boldsymbol{\theta}_1) [E_A + E_B K(\boldsymbol{\theta}_1)] + [E_A + E_B K(\boldsymbol{\theta}_1)]^\top P(\boldsymbol{\theta}_1) \right) e^{\bar{A}_1 t} dt. \quad (65)$$

948 Recall that the above expression is the derivative of $P(\boldsymbol{\theta})$ at $\boldsymbol{\theta}_1$, along the linear segment \mathcal{C} . Thus,
 949 integrating along \mathcal{C} , (65) and Cauchy-Schwarz Inequality imply that

$$\|P(\hat{\boldsymbol{\theta}}) - P(\boldsymbol{\theta}_0)\| \lesssim \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| \sup_{\boldsymbol{\theta}_1 \in \mathcal{C}} \|P(\boldsymbol{\theta}_1)\| (1 \vee \|K(\boldsymbol{\theta}_1)\|) \int_0^\infty \|e^{\bar{A}_1 t}\|^2 dt.$$

950 Finally, using Lemma 11, (61), and (62), we obtain the desired result. ■

951 **D Numerical Results**

952 In this section, we provide further empirical results illustrating the performance of Algorithm 2 in the
 953 settings of flight control, as well and blood glucose control. First, we provide box plots depicting
 954 the distribution of the normalized squared estimation error and the normalized regret of Algorithm 2
 955 for X-29A airplane. Note that the corresponding worst- and average-case curves are presented in
 956 Figure 1. Then, Figures 3 and 4 provide the corresponding curves of estimation and regret versus
 957 time as well as the box-plots, for Boeing 747. Finally, we present similar empirical result for learning
 958 to control blood glucose level. As shown in the presented figures, Thompson sampling Algorithm 2
 959 clearly outperforms the competing reinforcement learning policy.

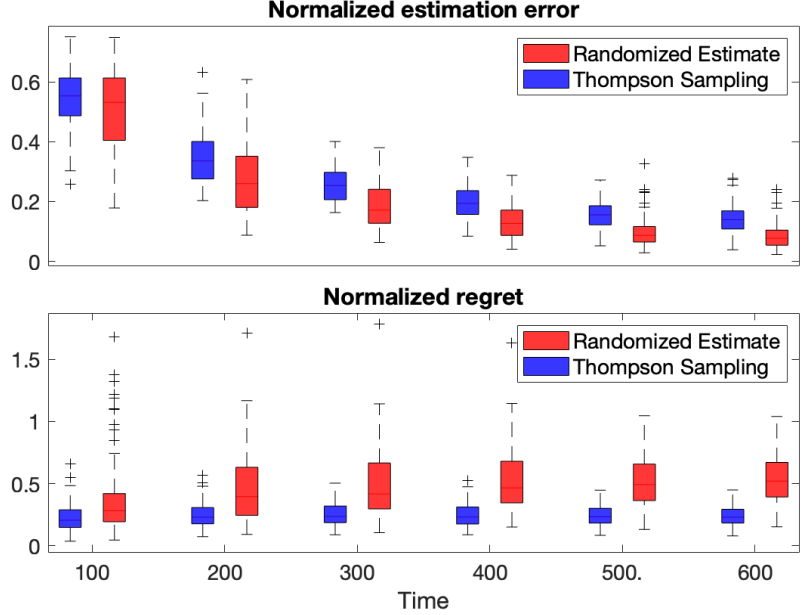


Figure 2: The performance of Algorithm 2 (blue) compared to Randomized Estimate policy (red) [2] for flight control of X-29A airplane. The top box-plots are for the normalized squared estimation error, $\|\hat{\theta}_n - \theta_0\|^2$ divided by $p(p+q)\tau_n^{-1/2} \log \tau_n$, at times 100, 200, 300, 400, 500, and 600 for 100 replications. Similarly, the lower graph showcases the distribution of the regret $\text{Reg}(T)$, normalized by $p(p+q)T^{1/2} \log T$.

960 Figure 2 depicts the box plot corresponding to Figure 1 that is for the flight control of X-29A airplane
 961 at 2000 ft. In the following experiments, we keep the setting given in Section 6 for the cost and noise
 962 covariance matrices, and compare Algorithm 2 to Randomized Estimate policy [2].

963 Next, the empirical results of the flight control problem in Boeing 747 airplane at 20000 ft altitude
 964 are provided [37]. The true drift matrices of the Boeing 747 are

$$A_0 = \begin{bmatrix} -0.199 & 0.003 & -0.980 & 0.038 \\ -3.868 & -0.929 & 0.471 & -0.008 \\ 1.591 & -0.015 & -0.309 & 0.003 \\ -0.198 & 0.958 & 0.021 & 0.000 \end{bmatrix}, \quad B_0 = \begin{bmatrix} -0.001 & 0.058 \\ 0.296 & 0.153 \\ 0.012 & -0.908 \\ 0.015 & 0.008 \end{bmatrix}.$$

965 Then, the blood glucose control problem is studied [41, 47]. The true drift matrices are

$$A_0 = \begin{bmatrix} 1.91 & -2.82 & 0.91 \\ 1.00 & -1.00 & 0.00 \\ 0.00 & 1.00 & -1.00 \end{bmatrix}, \quad B_0 = \begin{bmatrix} -0.0992 \\ 0.0000 \\ 0.0000 \end{bmatrix}.$$

966 Note that from a practical point of view, worst-case behavior are of crucial importance in this problem.

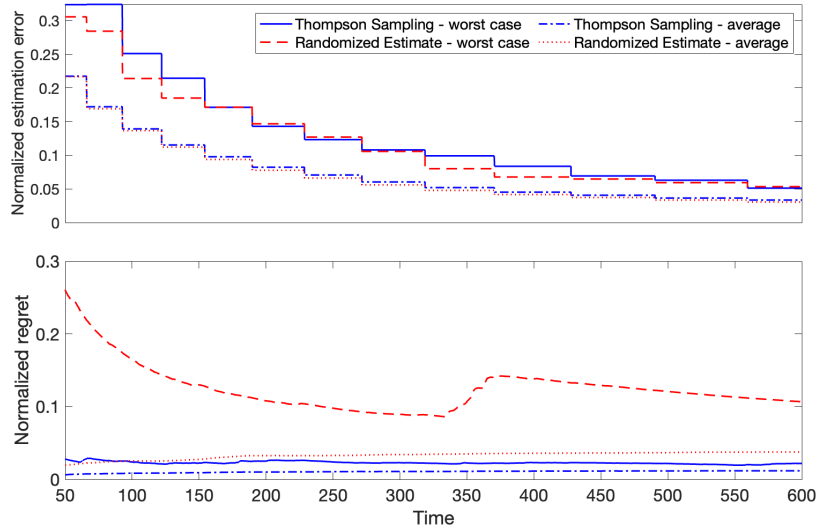


Figure 3: The performance of Algorithm 2 (blue) compared to Randomized Estimate policy (red) [2] for the flight control of Boeing 747 airplane. The top graph plots the normalized squared estimation error, $\left\| \hat{\theta}_n - \theta_0 \right\|^2$ divided by $p(p+q)\tau_n^{-1/2} \log \tau_n$, for 100 replications. Similarly, the lower graph showcases the regret $\text{Reg}(T)$, normalized by $p(p+q)T^{1/2} \log T$.

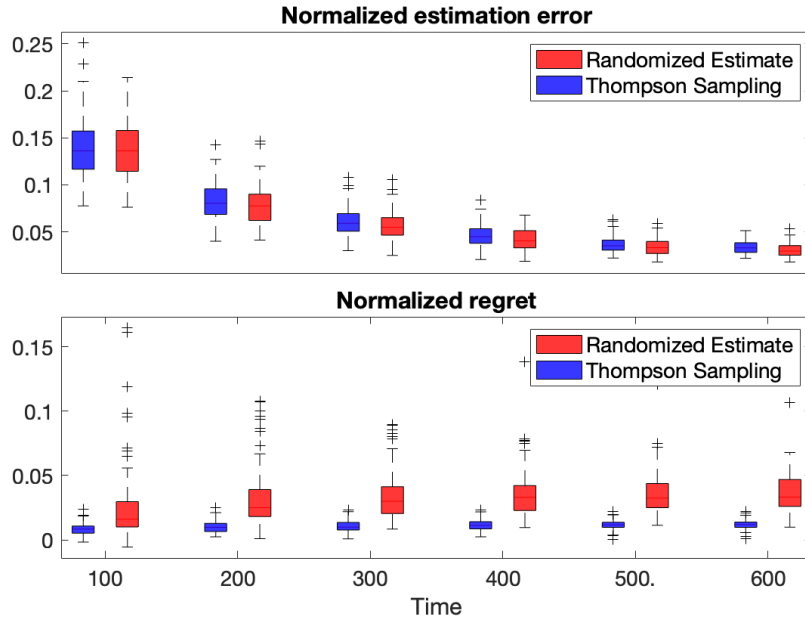


Figure 4: The performance of Algorithm 2 (blue) compared to Randomized Estimate policy (red) [2] for the flight control of Boeing 747 airplane. The top graph plots the normalized squared estimation error, $\left\| \hat{\theta}_n - \theta_0 \right\|^2$ divided by $p(p+q)\tau_n^{-1/2} \log \tau_n$, at times 100, 200, 300, 400, 500, and 600 for 100 replications. Similarly, the lower graph showcases the regret $\text{Reg}(T)$, normalized by $p(p+q)T^{1/2} \log T$.

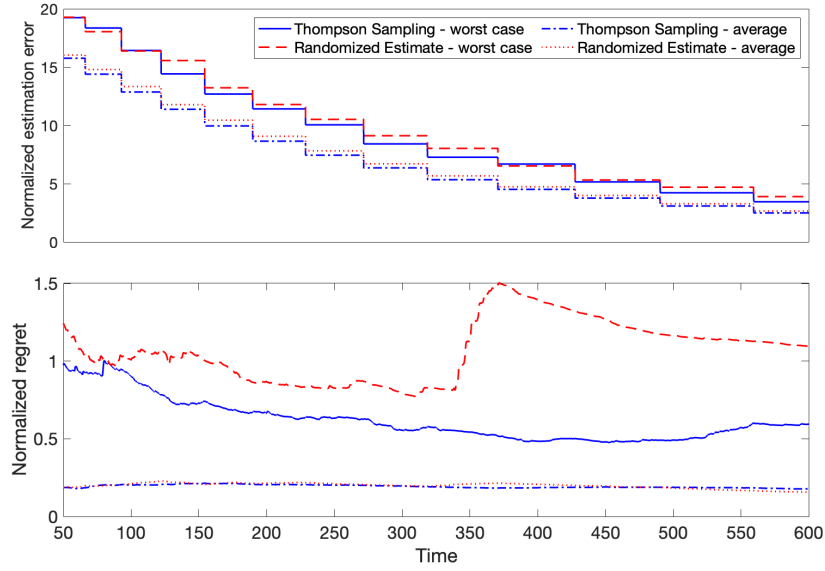


Figure 5: The performance of Algorithm 2 (blue) compared to Randomized Estimate policy (red) [2] for the the blood glucose control. The top graph plots the normalized squared estimation error, $\left\| \hat{\theta}_n - \theta_0 \right\|^2$ divided by $p(p+q)\tau_n^{-1/2} \log \tau_n$, for 100 replications. Similarly, the lower graph showcases the regret $\text{Reg}(T)$, normalized by $p(p+q)T^{1/2} \log T$.

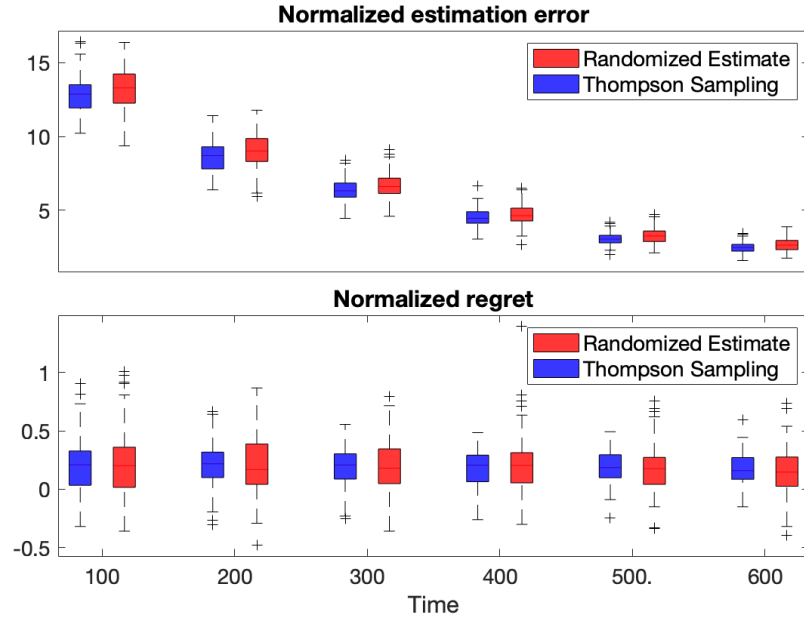


Figure 6: The performance of Algorithm 2 (blue) compared to Randomized Estimate policy (red) [2] for the the blood glucose control. The top graph plots the normalized squared estimation error, $\left\| \hat{\theta}_n - \theta_0 \right\|^2$ divided by $p(p+q)\tau_n^{-1/2} \log \tau_n$, at times 100, 200, 300, 400, 500, and 600 for 100 replications. Similarly, the lower graph showcases the regret $\text{Reg}(T)$, normalized by $p(p+q)T^{1/2} \log T$.