

Statistical Learning with Sublinear Regret of Propagator Models

Eyal Neuman and Yufei Zhang

Department of Mathematics, Imperial College London

January 23, 2025

Abstract

We consider a class of learning problems in which an agent liquidates a risky asset while creating both transient price impact driven by an unknown convolution propagator and linear temporary price impact with an unknown parameter. We characterize the trader's performance as maximization of a revenue-risk functional, where the trader also exploits available information on a price predicting signal. We present a trading algorithm that alternates between exploration and exploitation phases and achieves sublinear regrets with high probability. For the exploration phase we propose a novel approach for non-parametric estimation of the price impact kernel by observing only the visible price process and derive sharp bounds on the convergence rate, which are characterised by the singularity of the propagator. These kernel estimation methods extend existing methods from the area of Tikhonov regularisation for inverse problems and are of independent interest. The bound on the regret in the exploitation phase is obtained by deriving stability results for the optimizer and value function of the associated class of infinite-dimensional stochastic control problems. As a complementary result we propose a regression-based algorithm to estimate the conditional expectation of non-Markovian signals and derive its convergence rate.

Mathematics Subject Classification (2010): 62L05, 60H30, 91G80, 68Q32, 93C73, 93E35, 62G08

JEL Classification: C02, C61, G11

Keywords: optimal portfolio liquidation, price impact, propagator models, predictive signals, Volterra stochastic control, non-parametric estimation, reinforcement learning, regret analysis

Contents

1	Introduction	3
2	Problem formulation and main results	7
2.1	Episodic learning for optimal liquidation problems	7
2.2	A least-squares estimator and its convergence rate	11
2.3	Phased-based learning algorithm and its regret bound	16
3	Numerical implementation	20
4	Analytic solution to the control problem	24
4.1	Function spaces, integral operators	24
4.2	Essential operators and processes	25
5	Bound on the performance gap	26
6	Analysis of regularised least-squares estimators	29
7	Proof of Theorem 2.19	37
8	Proof of Proposition 5.3	39
9	Proofs of Lemma 8.1 and Proposition 8.2	43
A	Regression-based algorithm for signal estimation	46

1 Introduction

Price impact refers to the empirical fact that execution of a large order affects the risky asset's price in an adverse and persistent manner leading to less favourable prices. Propagator models are a central tool in describing this phenomena mathematically. This class of models provides deep insight into the nature of price impact and price dynamics. It expresses price moves in terms of the influence of past trades, which gives reliable reduced form view on the limit order book. It also provides interesting insights on liquidity, price formation and on the interaction between different market participants through price impact. The model's tractability provides a convenient formulation for stochastic control problems arising from optimal execution [11, 23]. More precisely, if the trader's holdings in a risky asset is denoted by $Q = \{Q_t\}_{t \geq 0}$, then the asset price S_t is given by

$$S_t = S_0 + \lambda \dot{Q}_t + \int_0^t G(t-s) dQ_s + P_t, \quad (1.1)$$

where P is a semimartingale, λ is the positive temporary price impact coefficient and the price impact kernel G is also known as the *propagator*. We refer to \dot{Q} as the execution trading speed. Since $G(t)$ typically decays for large values of t , the convolution on the right-hand-side of (1.1) is referred to as transient price impact (see e.g. Bouchaud et al. [11, Chapter 13]). A well-known example à la Almgren and Chriss, introduces to the case where G is a constant, then the above convolution represents permanent price impact (see [3, 4]).

In the aforementioned setting, the trader can only observe the visible price process S and her own inventory Q . In order to quantify the price impact and hence the trading costs, the trader needs a good estimation of G and λ . Some estimators for discrete-time versions of the model were proposed in [10, 20, 51, 52] and in Chapter 13.2 of [11], where only a finite amount of values $\{G(t_n)\}_{n=1}^N$ are estimated for a pre-determined grid $0 \leq t_1 < \dots < t_n$. However, even in this finite dimensional projection of the problem, the convergence of the estimators remains unproved, hence rigorous results on the estimation of G are considered as a long-standing open problem. In one of the main results of this paper we propose a novel approach for non-parametric estimation of the price impact kernel by observing only the visible price process and we derive sharp bounds on the convergence rate of our estimators, which are characterised by the singularity of the kernel.

Precise quantification of price impact is a crucial ingredient in portfolio liquidation problems. Considering the adverse effect of the price impact on the execution price, a trader who wishes to minimize her trading costs has to split her order into a sequence of smaller orders which are executed over a finite time horizon. At the same time, the trader also has an incentive to execute these split orders rapidly because she does not want to carry the risk of an adverse price move far away from her initial decision price. This trade-off between price impact and market risk is usually translated into

a stochastic optimal control problem where the trader aims to minimize a risk-cost functional over a suitable class of execution strategies, see [13, 24, 26, 29, 37, 39] among others. In addition, many traders and trading algorithms also strive for using short term price predictors in their dynamic order execution schedules, which are often related to order book dynamics as discussed in [35, 36, 38, 45]. From the modelling point of view, incorporating signals into execution problems translates into taking a non-martingale price process P , in contrast to a martingale price in the classical setting (see [12, 6, 40, 41, 1]). This changes the problem significantly as the resulting optimal strategies are often random and in particular signal-adaptive, in contrast to deterministic strategies, which are typically obtained in the martingale price case [7].

The main goal of this paper is to estimate the price impact kernel while trading a risky asset in a cost-effective manner. In order to do that we propose a learning algorithm that alternates between exploration and exploitation phases. In the exploration phase we propose a novel approach for non-parametric estimation of the price impact kernel by observing only the visible price process S and the signal process, which is the non-martingale component of P in (1.1) (see a note about observables in Section 2.1). Our estimation method extends the existing theory of Tikhonov regularisation for inverse problems and is of independent interest (see Remark 2.11). Specifically, we propose a regularised least-squares estimator for a squared integrable kernel G , where samples of the visible price process S are generated by a deterministic trading strategy executed by the trader. We derive *sharp bounds* on the convergence rate of the estimator with arbitrary high probability under two different assumptions. For a regular kernel, which has a squared integrable weak derivative, we prove that the convergence rate is of order $N^{-1/6}$. For a singular kernel with a decay rate $G(t) \sim t^{-\alpha}$ for some $\alpha \in (0, 1/2)$ we find that the convergence rate is of order $N^{-\frac{1-2\alpha}{2(3-2\alpha)}}$ (see Theorem 2.10). Here N is the sample size for the least-squares estimation. Moreover, from Proposition 6.5 it follows that the convergence rates given in Theorem 2.10 are optimal. More precisely, the outlined upper bounds for the rate of convergence match the lower bound rates of convergence under the assumption of regular kernel and of a power law kernel $G^*(t) = t^{-\alpha}$, with $\alpha \in (0, 1/2)$. See Remark 2.12 for specific details. Theorem 2.10 is the first result that proves convergence of any estimator for a propagator, which is based on market price quotes. These results answer an open question that arises from [10, 11, 20, 51, 52] among others.

After each exploration episode, we perform a number of exploitation episodes in which we use the current estimation of the kernel and temporary price impact coefficient $\theta^n = (\lambda^n, G^n)$ in a framework of optimal execution (see (1.1)). We execute the optimal trading strategy subject to the estimated parameter θ^n and derive the performance gap between the revenues of the aforementioned strategy and the optimal strategy with the real choice of θ^* . The bound on the performance gap, which is presented in Theorem 2.16, is derived by proving a stability result (see Proposition 5.3) for the optimizer of the associated infinite-dimensional stochastic control problem

proposed in [1].

Finally, by combining the exploration and exploitation schemes we propose Algorithm 1 which achieves sublinear regrets in high probability. In Theorem 2.19 and Corollary 2.21 we derive a bound on the regret after N trading episodes, which is of order $N^{3/4}$ for a regular kernel and of order $N^{\frac{3-2\alpha}{4-4\alpha}}$ for a singular kernel. These sublinear regret bounds underperform the square-root (or logarithmic) regret for reinforcement learning problems with finite-dimensional parametric models (see e.g., [5, 30, 48, 21, 22, 49]), due to the present infinite-dimensional non-parametric kernel estimation (see Remark 2.20).

In order to complete our argument, we provide a regression-based algorithm for signal estimation which is performed off-line, that is, independently from the trading algorithm. Specifically we decompose the semimartingale price process P in (1.1) to a martingale and to a finite variation process A which has the interpretation of a trading signal (see e.g. [41]). We observe that the optimal trading strategy which is used in the algorithm (see (2.8)) involves the conditional process $(t, s) \mapsto \mathbb{E}[A_s | \mathcal{F}_t]$. As the conditional distribution of A is in general not observable, we propose a regression-based algorithm to estimate it, based on observed signal trajectories. Since the agent's trading strategy does not affect the signal, the signal estimation can be carried out separately from the learning algorithm for (λ^*, G^*) . The convergence rate of this algorithm is derived in Theorem A.6.

Our main results which were outlined above significantly extend the work on reinforcement learning for continuous-time parametric models which were studied by [5, 30, 48, 21, 22, 49] among others. We outline our main contributions that correspond to each part of the learning algorithm.

Non-parametric kernel estimation: The main component of the exploration phase is to estimate the kernel function G in the *non-Markovian* model (1.1) in non-parametric manner. This results in a learning problem with infinite dimensional parameter, input, and output spaces, which stands in contrast to existing theoretical works on discrete Markov decision processes (see e.g. [43, 15, 32]) and on continuous time parametric Markov processes [5, 30, 48, 21, 22, 49]. The sample complexity bounds therein depend explicitly on the dimensions of the parameter space, the input space, and the output space, and hence cannot be applied in the present infinite dimensional setting.

It is also challenging to apply existing functional linear regression (FLR) frameworks to our model (1.1). Standard FLR frameworks directly estimate the mapping from the input Q to the response S as an unknown regression operator, instead of estimating λ and G individually. Moreover, most existing FLR works characterise the convergence rate of the proposed estimators under the so-called source conditions, which assume that the unknown regression operator lies in the range of a suitable fractional power of the input covariance operator (see, e.g., [8, Assumption 4] and also [8, 46]). It is well known that identifying explicit conditions on G and \dot{Q} such that

these source conditions are met is challenging, as it typically requires computing the spectral decomposition of the input covariance and the unknown regression operator, which cannot be performed analytically for general G and \dot{Q} .

In this work, we propose a novel method for the convergence rate analysis of the estimator of G which applies even for singular kernels. For a suitable deterministic trading speed \dot{Q} , the proposed method estimates λ using a classical Monte Carlo estimator and estimates G using a regularised least-squares estimator involving the estimated λ . We further quantify the convergence rate of the estimator for G by using *appropriate source conditions*, instead of the standard source conditions found in [8, 9, 46]. These appropriate source conditions quantify the degree to which the true kernel G violates the assumption of being in the range of the input operator (see (6.7)) and were introduced in [33] to study deterministic inverse problems. We extend these ideas to the present setting with stochastic observations (see Theorem 2.10 and Remark 2.11).

In particular, we identify explicit regularity conditions on the kernel G and the trading speed \dot{Q} such that the approximate source conditions hold and optimise the convergence rates of the estimator accordingly. The resulting convergence rates are optimal under the regularity conditions of the kernel G and are better than the worst-case convergence rates given in the statistical inverse problem literature [9, 46]; see Remark 2.12. The method is general and can be applied for non-parametric estimation in other classes of infinite-dimensional stochastic control problems.

Lipschitz stability of infinite-dimensional optimal control: In [5, 30, 48], Lipschitz stability of optimal controls for stochastic control problems was derived, by showing that the optimiser is continuously differentiable with respect to finite-dimensional model parameters, hence establishing Lipschitz continuity. The Lipschitz stability of controls is crucial for quantifying the precise performance gap between controls derived from estimated and true models, and for characterizing the regret order of learning algorithms.

The performance gap in Theorem 2.16 entails proving the Lipschitz stability of the optimal control with respect to the (infinite-dimensional) kernel function G (see Proposition 5.3), for which the preceding argument developed for finite-dimensional parameters does not apply. Moreover, in our setting the non-Markovianity introduced by G turns the problem to infinite-dimensional stochastic control, in contrast to the finite-dimensional control problems studied in the aforementioned references. This major difference is reflected in the ingredients of optimiser. For example standard Riccati equations become operator-valued Riccati equations and solution to a BSDE becomes a solution to infinite dimensional BSDE (see Sections 6.2-6.3 of [1] for additional details). In this work, we establish uniform boundedness and Lipschitz stability for all components of the optimiser in suitable norms (see Remark 2.17).

Also note that for the control problem studied in this paper the running cost is not strongly concave in the control variable, and is not concave in the state variable

(see (2.7)). Such a (strong) concavity assumption is assumed in the aforementioned references in order to establish the required Lipschitz stability. We overcome this issue by imposing a nearly non-negativeness condition of the estimated kernel G (see Definition 2.14) in order to prove stability.

Organisation of the paper: In Section 2 we describe the reinforcement learning problem and present our main results on the convergence of the estimator for propagator and on the regret bounds. Section 3 is dedicated to a numerical implementation of our propagator estimation results. In Section 4 we recall some essential results on the associated optimal liquidation problem. Section 5 is dedicated to the proof Theorem 2.16, which provides the bound on the performance gap. Section 6 deals with the analysis of the regularised least-squares estimator. Section 7 dedicated to the proof of Theorem 2.19, which derives regret rate. Sections 8 and 9 contain proofs for some auxiliary results. Finally in Appendix A we provide a regression-based algorithm for signal estimation and derive its convergence rate.

2 Problem formulation and main results

This section studies the optimal liquidation problem with unknown transient price impact kernel and temporary price coefficient $\theta^* = (\lambda^*, G^*)$. The agent's objective is to search for the optimal trading strategy while simultaneously learn the price dynamics, that is to learn θ^* . We first propose a least-squares estimator for these coefficients and derive its convergence rate. Then we present a phased-based learning algorithm and establish its regret bound.

2.1 Episodic learning for optimal liquidation problems

Optimal liquidation with known price impacts. We first recall the optimal liquidation framework which was presented in [1].

Let $T > 0$ denote a finite deterministic time horizon and fix a filtered probability space $(\Omega, \mathcal{F}, (\mathcal{F}_t)_{0 \leq t \leq T}, \mathbb{P})$ satisfying the usual conditions of right continuity and completeness. We consider a semimartingale unaffected price process $P = (P_t)_{0 \leq t \leq T}$ with a canonical decomposition

$$P_t = A_t + M_t, \quad 0 \leq t \leq T, \quad (2.1)$$

into a predictable finite-variation signal process $A = (A_t)_{0 \leq t \leq T}$ and an independent martingale M satisfying $\mathbb{E}[M_0] = 0$ and

$$E[\langle M \rangle_T] + E\left[\left(\int_0^T |dA_s|\right)^2\right] < \infty. \quad (2.2)$$

Let $\theta^* = (\lambda^*, G^*) \in (0, \infty) \times L^2([0, T], \mathbb{R})$ be fixed coefficients such that for every $f \in L^2([0, T], \mathbb{R})$,

$$\int_0^T \int_0^T G^*(|t-s|) f(s) f(t) ds dt \geq 0. \quad (2.3)$$

Remark 2.1. *Note that in (2.3) we consider a class of non-negative kernels. An important subclass of these kernels is the class of bounded non-increasing convex functions (see Example 2.7 in [25]). See Remark 2.7 for further examples.*

We consider a trader with an initial position of $q > 0$ shares in a risky asset. The number of shares the trader holds at time $t \in [0, T]$ is prescribed as

$$Q_t^u = q - \int_0^t u_s ds, \quad (2.4)$$

where $(u_s)_{s \in [0, T]}$ denotes the trading speed which is chosen from the set of admissible strategies

$$\mathcal{A} \triangleq \left\{ u : u \text{ progressively measurable s.t. } \mathbb{E} \left[\int_0^T u_s^2 ds \right] < \infty \right\}. \quad (2.5)$$

We assume that the trader's trading activity causes price impact on the risky asset's execution price. In order to define the price impact effects we introduce some additional definitions. For any trading speed $u \in \mathcal{A}$, the price process S^u satisfies the following dynamics: for all $t \in [0, T]$,

$$S_t^u := P_t - \lambda^* u_t - Z_t^{\theta^*, u}, \quad \text{with } Z_t^{\theta^*, u} = \int_0^t G^*(t-s) u_s ds. \quad (2.6)$$

Note that λ^* is the temporary price impact coefficient and $Z_t^{\theta^*, u}$ is the transient price impact term, which is associated with the price impact kernel G^* , also known as the propagator.

Consider maximising the following risk-revenue functional over $u \in \mathcal{A}$:

$$J^{\theta^*}(u) := \mathbb{E} \left[\int_0^T S_t^u u_t dt + Q_T^u P_T - \phi \int_0^T (Q_t^u)^2 dt - \varrho (Q_T^u)^2 \right]. \quad (2.7)$$

The first two terms in (2.7) represent the trader's terminal wealth; that is, her final cash position including the accrued revenue as well as her remaining final risky asset position's book value. The third and fourth terms in (2.7) implement a penalty $\phi \geq 0$ and $\varrho \geq 0$ on her running and terminal inventory, respectively. Observe that $J(u) < \infty$ for any strategy $u \in \mathcal{A}$.

If the agent knows θ^* , then (2.6)-(2.7) is a special case of the Volterra stochastic control problem studied in [1]. By Proposition 4.5 therein, the optimal trading strategy u^{θ^*} is given by

$$u_t^{\theta^*} = a_t^{\theta^*} + \int_0^t B^{\theta^*}(t, s) u_s^{\theta^*} ds, \quad 0 \leq t \leq T, \quad (2.8)$$

where a^{θ^*} is a stochastic process satisfying (2.5), depending on A but not on M in (2.1), and B^{θ^*} is a function satisfying $\sup_{t \leq T} \int_0^t (B^{\theta^*}(t, s))^2 ds < \infty$, see (4.13) for the precise definition. We emphasise the dependence of u^{θ^*} in (2.8) by writing

$$u^{\theta^*} = \text{Greedy}(A, \theta^*).$$

A note about observables. Recall that the visible price process S^u was introduced in (2.6). In addition to this observable the agent clearly knows her own trading rate u , which impacts S^u . Recall that the fundamental price process P was defined in (2.1). While P is unobserved by the trader, it is a common practice that the short term price predicting signal A (also called *alpha*) is an observable, typically obtained from limit order books real-time data. We briefly survey some well known examples for such signals which impact the price at different time scales. In Section 4 of [36] a detailed statistical analysis of the limit order book imbalance signal was performed. The effect of this signal on future price moves was demonstrated in time intervals of the 10 next trades. The usage of this signal by high frequency proprietary traders was also proved statistically. The order flow imbalance signal has been extensively studied in the literature, in particular the correlation between the current order flow and the future price move in 10 seconds intervals was studied by R. Cont and Stoikov [45]. More examples of observed trading signals which are used in optimal execution can be found in a practitioners presentation by Robert Almgren [2]. In reality the agent also determines the penalty parameters ϕ, ϱ in the quadratic costs (2.7), however the parameters $\theta^* = (\lambda^*, G^*)$ are unknown and are subject to estimation in this paper.

Optimal liquidation with unknown price impacts. In this work, we consider an agent who repeatedly liquidates the risky asset in (2.6) without knowing the price impact coefficient θ^* . This is often referred to as the episodic (also known as reset or restart) learning framework in the reinforcement learning literature. The agent will improve her knowledge of θ^* through successive learning episodes, while simultaneously optimise the objective (2.7). In reality the agent knows the dynamics of A in (2.1), the form of (2.6) (excluding the coefficient θ^*) and the penalty parameters ϕ, ϱ in the quadratic costs (2.7). We will therefore assume that these are known features of the model in the following. For each episode, the agent observes (a realisation of) the price S and the signal A , but not the noise M . Additional regularity properties of θ^* will be assumed in order to optimise the learning algorithm (see Assumptions 2.6 and 2.13).

Mathematically, the learning problem is described as follows. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(A^m)_{m \in \mathbb{N}}$ and $(M^m)_{m \in \mathbb{N}}$ be mutually independent copies of A and M on $(\Omega, \mathcal{F}, \mathbb{P})$, respectively, and for each $m \in \mathbb{N}$, let $P^m = A^m + M^m$. Here A^m and M^m correspond to the observed signal and unobserved martingale noise for the m -th learning episode, respectively. For each episode, the agent interacts with (2.6) by choosing controls that are adapted to available observations. These admissible controls and observation filtrations are defined recursively as follows. The observation before the first episode is given by the σ -algebra $\mathcal{F}_0 = \mathcal{N}$, where \mathcal{N} is the σ -algebra generated by the \mathbb{P} -null set. For the m -th episode with $m \in \mathbb{N}$, taking the σ -algebra \mathcal{F}_{m-1} , the agent executes a square-integrable control u^m that is progressively measurable with respect to the filtration $(\mathcal{G}_t^m)_{t \in [0, T]}$ with $\mathcal{G}_t^m := \mathcal{F}_{m-1} \vee \sigma\{A_s^m \mid s \in [0, t]\}$, and observes a trajectory of the price process S^m governed by the following dynamics (cf. (2.6)):

$$S_t^m = A_t^m + M_t^m - \lambda^* u_t^m - \int_0^t G^*(t-s) u_s^m ds. \quad (2.9)$$

The available information for the agent before the $(m+1)$ -th episode is $\mathcal{F}_m := \mathcal{F}_{m-1} \vee \sigma\{S_t^m, A_t^m \mid t \in [0, T]\}$.

To measure the performance of the controls $(u^m)_{m \in \mathbb{N}}$ (also referred to as a learning algorithm) in this setting, one widely adopted criteria is the regret of learning [30, 5, 21, 22]: for each $N \in \mathbb{N}$, the regret of learning up to N -th episode is given by

$$R(N) = \sum_{m=1}^N (J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^m)), \quad (2.10)$$

where $J^{\theta^*}(u^{\theta^*})$ is the optimal value that agent can achieve knowing the parameter θ^* , and $J^{\theta^*}(u^m)$ is the expected performance of the control u^m for the m -episode¹:

$$J^{\theta^*}(u^m) := \mathbb{E} \left[\int_0^T S_t^m u_t^m dt + Q_T^{u^m} P_T^m - \phi \int_0^T (Q_t^{u^m})^2 dt - \varrho (Q_T^{u^m})^2 \middle| \mathcal{F}_{m-1} \right], \quad (2.11)$$

with $Q_t^{u^m} = q - \int_0^t u_s^m ds$ being the corresponding inventory (cf. (2.4)). The expectation in (2.11) is only taken with respect to P^m , and hence $J^{\theta^*}(u^m)$ is a random variable depending on the realisations of the signals $(A^n)_{n=1}^{m-1}$ and noises $(M^n)_{n=1}^{m-1}$. Intuitively, the regret $R(N)$ characterises the cumulative expected loss from taking sub-optimal controls up to the N -th episode. Agent's aim is to construct a learning algorithm for which the regret $R(N)$ grows sublinearly in N in high probability.

Note that the above setting assumes the algorithm runs indefinitely without a prescribed maximal number of learning episodes. The agent will then derive an *anytime*

¹With a slight abuse of notation, we denote by $J^{\theta^*}(\cdot)$ the performance functional for all episodes, without specifying its dependence on m . It is possible as P^m is independent of \mathcal{F}_{m-1} .

learning algorithm (see e.g., [34, 48]), i.e., an algorithm whose implementation does not require advance knowledge of the algorithm termination time and whose performance guarantee holds for all learning episodes; see Remark 2.18 for more details.

2.2 A least-squares estimator and its convergence rate

In this section we derive the identifiability of $\theta^* = (\lambda^*, G^*)$ under suitable exploratory strategies. We propose a regularised least-squares estimator based on observed trajectories and analyse its finite sample accuracy. The estimator will be employed in Section 2.3 to design a regret optimal learning algorithm for (2.7). By an abuse of notation, we will index the observed trajectories for the estimator by m , which is typically different from the number of learning episodes in Section 2.3.

More precisely, let A and M the processes in (2.6), and let $(A^m)_{m \in \mathbb{N}}$ and $(M^m)_{m \in \mathbb{N}}$ be mutually independent copies of A and M , respectively, defined on the probability space $(\Omega, \mathcal{F}, \mathbb{P})$. The agent executes a trading strategy $u^e \in L^2([0, T], \mathbb{R})$, and estimates $\theta^* = (\lambda^*, G^*)$ using the corresponding price trajectories $(S^m, A^m)_{m \in \mathbb{N}}$, where for all $m \in \mathbb{N}$, $(S^m, A^m)_{t \in [0, T]}$ satisfies for all $t \in [0, T]$,

$$\begin{aligned} S_t^m &= A_t^m + M_t^m - \lambda^* u^e(t) - \int_0^t G^*(t-s) u^e(s) ds \\ &= A_t^m + M_t^m - \lambda^* u^e(t) - (\mathbf{u}^e G^*)(t), \end{aligned} \quad (2.12)$$

with $\mathbf{u}^e : L^2([0, T], \mathbb{R}) \rightarrow L^2([0, T], \mathbb{R})$ being the integral operator defined by

$$(\mathbf{u}^e f)(t) := \int_0^t u^e(t-s) f(s) ds, \quad f \in L^2([0, T], \mathbb{R}). \quad (2.13)$$

Note that in (2.12), G^* plays the role of an unknown function instead of a kernel.

The following regularity condition on u^e is imposed for the identifiability of θ^* . Recall that $H^1([0, T], \mathbb{R})$ is the space of absolute continuous functions $f : [0, T] \rightarrow \mathbb{R}$ whose derivative (which exists a.e.) belongs to $L^2([0, T], \mathbb{R})$.

Assumption 2.2. $u^e \in H^1([0, T], \mathbb{R})$ and $u^e(0) \neq 0$.

Remark 2.3. Assumption 2.2 holds for any nonzero constant strategy or classical trading strategies in the Almgren–Chriss framework (see e.g. [13, Chapter 6]). Unfortunately, the trajectories of the greedy strategy u^{θ^*} in (2.8) may not satisfy Assumption 2.2. Indeed, $u_0^{\theta^*}$ could be zero, due to the randomness of the signal process A . Moreover, the time regularity of u^{θ^*} relies on the regularity of a_t and $B(t, \cdot)$ with respect to t , which subsequently depends on the path regularity of the conditional expectations of the signal A (cf. (4.13)). Even for the special case with $A \equiv 0$, it is still challenging to obtain explicit conditions for the differentiability of $t \mapsto B(t, \cdot)$ and $t \mapsto a_t$ to ensure that $u^{\theta^*} \in H^1([0, T], \mathbb{R})$.

Under Assumption 2.2, the operator $\mathbf{u}^e : L^2([0, T], \mathbb{R}) \rightarrow L^2([0, T], \mathbb{R})$ is injective as shown in Lemma 6.1. This indicates that θ^* can be uniquely identified based on sufficiently many trajectories $(S^m, A^m)_{m \in \mathbb{N}}$. In the sequel, we propose a regularised least-squares estimator for θ^* and analyse its finite sample accuracy.

By (2.6) and $\mathbb{E}[M_0] = 0$, $\lambda^* u^e(0) = -\mathbb{E}[S_0 - A_0]$. Replacing the expectation by an empirical mean yields the following estimation for λ^* :

$$\lambda^N := -\frac{1}{Nu^e(0)} \sum_{m=1}^N (S_0^m - A_0^m), \quad \text{for } N \in \mathbb{N}, \quad (2.14)$$

which is well-defined as $u^e(0) \neq 0$. Given the estimators $(\lambda^N)_{N \in \mathbb{N}}$, we then introduce a sequence of projected least-squares estimators for the kernel G^* . To this end, let $\mathcal{P}_{[0, T]}$ be the collection of all partitions of $[0, T]$, and let $(\pi_N)_{N \in \mathbb{N}} \subset \mathcal{P}_{[0, T]}$ be such that $\pi_N = \{0 = t_0^{(N)} < \dots < t_N^{(N)} = T\}$ for all $N \in \mathbb{N}$ and $\lim_{N \rightarrow \infty} |\pi_N| = 0$, where $|\pi_N| := \max_{i=0, \dots, N-1} (t_{i+1}^{(N)} - t_i^{(N)})$ is the mesh size of π_N . For each $N \in \mathbb{N}$, let V_N be the space of piecewise constant functions on π_N :

$$V_N = \left\{ f \in L^2([0, T], \mathbb{R}) \mid f_t = \sum_{i=0}^{N-1} f_{t_i} \mathbb{1}_{[t_i^{(N)}, t_{i+1}^{(N)})}(t), \text{ for all } t \in [0, T] \right\}. \quad (2.15)$$

Then for each regularising weight $\tau_N > 0$, consider minimising the following L^2 -regularised least-squares estimation error over V_N (cf. (2.12)):

$$G^N := \arg \min_{G \in V_N} \left(\frac{1}{N} \sum_{m=1}^N \|S^m - A^m + \lambda^N u^e + \mathbf{u}^e G\|_{L^2([0, T])}^2 + \tau_N \|G\|_{L^2([0, T])}^2 \right), \quad (2.16)$$

which is derived by replacing λ^* in (2.12) with λ^N . As $\tau_N > 0$, it is easy to see that the quadratic functional in (2.16) has a unique minimum and hence G^N is well-defined.

Remark 2.4. Here, we take $(V_N)_{N \in \mathbb{N}}$ to be spaces of piecewise constant functions for the clarity of presentation, but the estimator (2.16) and its convergence analysis can be extended to any subspaces $(V_N)_{N \in \mathbb{N}}$ of $L^2([0, T], \mathbb{R})$ such that $\overline{\bigcup_{N=1}^{\infty} V_N} = L^2([0, T], \mathbb{R})$; see Section 2.3 for details.

For notational simplicity, we write $\theta^N = (\lambda^N, G^N)$ in (2.14) and (2.16) as

$$\theta^N = \text{LSE}((S^m, A^m)_{1 \leq m \leq N}, \tau_N, \pi_N), \quad (2.17)$$

which emphasises the dependence on the data $(S^m, A^m)_{1 \leq m \leq N}$, the regularising weight τ_N and the mesh size of π_N .

Remark 2.5. *The fact that $\tau_N > 0$ is critical for the well-posedness of (2.16). Indeed, consider $\tau_N = 0$, $V_N = L^2([0, T], \mathbb{R})$ and $u^e \equiv 1$. Then (2.12) and (2.16) suggest that*

$$G^N(t) = -\frac{1}{N} \sum_{m=1}^N (dS_t^m - dA_t^m) = -\frac{1}{N} \sum_{m=1}^N dM_t^m + G^*(t), \quad \forall t \in [0, T]. \quad (2.18)$$

As a non-constant continuous martingale has infinite variation, $G^N \in L^2([0, T], \mathbb{R})$ satisfying (2.18) does not exist in general.

The above observation also indicates that a proper scaling of the regularising weight τ_N with respect to the sample size N is crucial for the smoothness and convergence of $(G^N)_{N \in \mathbb{N}}$. Reducing the weight τ_N too fast essentially fits the time fluctuation of $(M^m)_{m=1}^N$, and hence leads to an irregular estimate G^N . This is in contrast to the regularised least-squares estimator for parametric models as in [5]. The regularising weights $(\tau_N)_{N \in \mathbb{N}}$ therein can be chosen as any vanishing sequence such that $\limsup_{N \rightarrow \infty} \sqrt{N} \tau_N < \infty$.

The dependence of τ_N on N results in a slower convergence of $(G^N)_{N \in \mathbb{N}}$ compared with the $\mathcal{O}(N^{-1/2})$ order for classical Monte-Carlo methods. It is known that the optimal choice of $(\tau_N)_{N \in \mathbb{N}}$ depends on the regularity of the true kernel G^* (also known as the “source condition” in inverse problem literature [33, 9]). We impose the following regularity conditions on the kernel G^* .

Assumption 2.6. *$G^* \in L^2([0, T], \mathbb{R})$ is differentiable a.e., and is one of the two types:*

- (1) *Regular kernel: $G^* \in H^1([0, T], \mathbb{R})$ and $G^*(T) \neq 0$.*
- (2) *Power-type singular kernel: there exists $\alpha \in (0, 1/2)$, $t_0 \in (0, T)$ and $C_0 > 0$ such that $|\frac{d}{dt} G^*(t)| \leq C_0 t^{-\alpha-1}$ for a.e. $t \in (0, t_0)$, and $G^* \in H^1([t_0, T], \mathbb{R})$.*

Remark 2.7. *Note that Assumption 2.6(1) requires G^* to be continuous on $[0, T]$, due to Morrey’s inequality. It is satisfied by the exponential kernel $G^*(t) = e^{-\beta t}$ for $\beta > 0$ proposed by [42] and the truncated power law kernel $G^*(t) = (c_0 + t)^{-\beta}$ for some $\beta, c_0 > 0$ studied in [10, 23]. On the other hand, Assumption 2.6(2) allows for a power-type singularity at $t = 0$. It includes as a special case the power law kernel $G^*(t) = t^{-\beta}$, for any $0 < \beta \leq \alpha$ proposed in [23]. Note that the constant α , which determines the range of power law singularities allowed in the kernel G^* , is well documented in the literature, both by non-rigorous empirical estimates using historical data (see e.g. [10, 50]) and from theoretical arguments (see [23] and Chapter 13 of [11]).*

In the sequel, we assume that the agent knows the precise type of G^* as in Assumption 2.6, i.e., G^* is regular on $[0, T]$ or admits a power-type singularity at zero with

known component α . This allows for specifying the precise decay rate of $(\tau_N)_{N \in \mathbb{N}}$ in (2.16) and then establishing the convergence rate of $(G^N)_{N \in \mathbb{N}}$. To quantify the convergence rate of $(\lambda^N, G^N)_{N \in \mathbb{N}}$ in high probability, we impose the following concentration condition on the martingale process M .

Assumption 2.8. *There exists $C_M > 0$ such that for all $N \in \mathbb{N}$ and $\eta > 0$,*

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{m=1}^N M_0^m \right|^2 + \left\| \frac{1}{N} \sum_{m=1}^N M^m \right\|_{L^2([0,T])}^2 \geq C_M^2 (\log(2\eta^{-1}))^2 N^{-1} \right) \leq \eta.$$

The following lemma provides a sufficient condition for Assumption 2.8.

Lemma 2.9. *There exists $L, \sigma > 0$ such that for all $p \geq 2$, $\mathbb{E}[(|M_0|^2 + \|M\|_{L^2([0,T])}^2)^{p/2}] \leq \frac{1}{2} p! \sigma^2 L^{p-2}$. Then Assumption 2.8 holds with $C_M = 2(L + \sigma)$.*

The proof of Lemma 2.9 follows by applying [9, Proposition A.1] to the random variable $Z := (M_0, M)$ taking values in the Hilbert space $\mathbb{R} \times L^2([0, T])$. The moment condition in Lemma 2.9 is commonly referred to as a Bernstein-type assumption. It is often imposed on the observation noise distribution in the statistical inverse problem literature for conducting complexity analysis [9, 46]. For martingales given by stochastic integrals with respect to Brownian motions or Poisson measures, this moment condition can be verified by Burkholder's inequality as in [30]. In the sequel, we will directly work with Assumption 2.8, which is more general than the Bernstein assumption, and is sufficient for obtaining the optimal convergence rate of (2.17) in high probability.

Under Assumptions 2.6 and 2.8, the following theorem chooses the optimal regularising weights $(\tau_N)_{N \in \mathbb{N}}$ and mesh sizes $(|\pi_N|)_{N \in \mathbb{N}}$, and quantifies the convergence rate of $(\lambda^N, G^N)_{N \in \mathbb{N}}$ in high probability. It follows as a special case of Theorem 6.3 in Section 6.

Theorem 2.10. *Suppose that Assumptions 2.2 and 2.8 hold. Let $C \geq 1$.*

- (1) *If Assumption 2.6(1) holds, then for all $\eta \in (0, 1)$, by setting $(\tau_N)_{N \in \mathbb{N}} \subset (0, \infty)$ and $(\pi_N)_{N \in \mathbb{N}} \subset \mathcal{P}_{[0,T]}$ such that for all $N \in \mathbb{N}$,*

$$\frac{1}{C} \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{4}{3}} \leq \tau_N \leq C \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{4}{3}}, \quad |\pi_N| \leq C \tau_N^{\frac{1}{2}}, \quad (2.19)$$

it holds with probability at least $1 - \eta$ that, for all $N \in \mathbb{N} \cap [2, \infty)$,

$$|\lambda^N - \lambda^*| \leq C' \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right), \quad \|G^N - G^*\|_{L^2([0,T])} \leq C' \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{1}{3}}. \quad (2.20)$$

(2) If Assumption 2.6(2) holds, then for all $\eta \in (0, 1)$, by setting $(\tau_N)_{N \in \mathbb{N}} \subset (0, \infty)$ and $(\pi_N)_{N \in \mathbb{N}} \subset \mathcal{P}_{[0, T]}$ such that for all $N \in \mathbb{N}$,

$$\frac{1}{C} \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{4}{3-2\alpha}} \leq \tau_N \leq C \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{4}{3-2\alpha}}, \quad |\pi_N| \leq C \tau_N^{\frac{1}{2}}, \quad (2.21)$$

it holds with probability at least $1 - \eta$ that, for all $N \in \mathbb{N} \cap [2, \infty)$,

$$|\lambda^N - \lambda^*| \leq C' \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right), \quad \|G^N - G^*\|_{L^2([0, T])} \leq C' \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{1-2\alpha}{3-2\alpha}}. \quad (2.22)$$

The constant $C' > 0$ appearing in (2.20) and (2.22) is independent of η and N .

The proof of Theorem 2.10 is given in Section 6.

Remark 2.11. Theorem 2.10 is proved by first interpreting (2.16) as Tikhonov regularisation of (2.12), and then adapting existing theoretical frameworks of Tikhonov regularisation with deterministic observations (see [28, 33] and references therein) to the present setting with random observations. The crucial step in our argument is to quantify the distance between the true kernel G^* to the range of the operator \mathbf{u}^e in (2.13), which is characterised by the behavior of a distance function $R \mapsto \mathcal{D}(R)$ for large R (cf. (6.7)). We prove in Theorem 6.4 that the function \mathcal{D} decays as a power function as $R \rightarrow \infty$, whose exponent depends explicitly on Assumption 2.6. To the best of our knowledge, such a power-type decay of \mathcal{D} has only been established in the literature for $\mathbf{u}^e \equiv 1$ and $G^* \equiv 1$ (see e.g., [33]). We further prove in Proposition 6.5 that the exponents of these power functions are optimal, i.e., they are the maximal power-type decay rates of \mathcal{D} under Assumption 2.6. Specifically, Proposition 6.5 considers the power law kernel $G^*(t) = t^{-\alpha}$, which satisfies Assumption 2.6(1) if $\alpha = 0$, and Assumption 2.6(2) if $\alpha \in (0, 1/2)$.

Remark 2.12. The convergence rates given in Theorem 2.10 are optimal (up to a logarithmic order in N) under Assumption 2.6. In order to see this, note that the estimation of G^* can be interpreted as an inverse problem with the forward operator \mathbf{u}^e based on noisy observations, where the noise level is of the magnitude $\mathcal{O}(1/\sqrt{N})$ with high probability; see Theorem 6.3. By the order optimality result [19, Proposition 3.15], under the source condition that $G^* = ((\mathbf{u}^e)^* \mathbf{u}^e)^\mu w$ for some $w \in L^2([0, T]; \mathbb{R})$ and $\mu > 0$, no estimation algorithm can recover G^* with a rate faster than $\mathcal{O}(N^{-\frac{\mu}{2\mu+1}})$ as $N \rightarrow \infty$. We then characterise the worst μ for a kernel G^* satisfying Assumption 2.6. Recall that by [19, Proposition 3.13], G^* satisfies the source condition with $\mu > 0$ if and only if $\sum_{n=1}^{\infty} \frac{1}{\sigma_n^{4\mu}} \langle G^*, \mathbf{u}_n \rangle_{L^2([0, T])}^2 < \infty$, where $\sigma_1 \geq \sigma_2 \geq \dots > 0$ is the singular values of \mathbf{u}^e , and $(\mathbf{u}_n)_{n \in \mathbb{N}}$ is the orthonormal system of eigenfunctions of $(\mathbf{u}^e)^* \mathbf{u}^e$. Now consider the power law kernel $G^*(t) = t^{-\alpha}$, which satisfies Assumption 2.6(1) if $\alpha = 0$, and Assumption 2.6(2) if $\alpha \in (0, 1/2)$. By (6.17) and the above criterion,

the power law kernel $G^*(t) = t^{-\alpha}$ satisfies the source condition for $\mu < \frac{1}{2}(\frac{1}{2} - \alpha)$ (but not $\mu = \frac{1}{2}(\frac{1}{2} - \alpha)$). This suggests that under Assumption 2.6, the optimal rate is not greater than $\mathcal{O}(N^{-\frac{1}{2}\frac{1-\alpha}{3-2\alpha}})$ as $N \rightarrow \infty$. This lower rate of convergence matches the upper rate of convergence in Theorem 2.10 (up to a logarithmic term), which indicates the parameter choices in Theorem 2.10 are order optimal under Assumption 2.6.

Note that the convergence rates in Theorem 2.10 are better than the lower rates for general statistical inverse problems with random input and output variables [9, 46]. This is due to the usage of a deterministic trading strategy u^e in (2.12), which allows for applying deterministic inverse problem theory to achieve an improved rate.

2.3 Phased-based learning algorithm and its regret bound

By leveraging Theorem 2.10, in this section we propose a phased-based algorithm for learning (2.6)-(2.7). The algorithm alternates between exploration and exploitation phases, and achieves sublinear regrets with high probability.

Admissible estimated models. We first introduce a class of estimated models based on which the greedy policies are constructed during the learning process. To facilitate the regret analysis, we assume that the agent knows the order of magnitude of the true parameter θ^* , from heuristic estimations using historical data (see [10], [11, Chapter 13], [13, Chapter 6.2] and [14] among others). Note that the constant L which is defined below is a known parameter of the problem along with ϕ, ϱ in (2.11).

Assumption 2.13. *There exists a known constant $L > 0$ such that $L^{-1} < \lambda^* < L$ and $\|G^*\|_{L^2([0,T])} < L$.*

Definition 2.14 (Class of admissible parameters Ξ_ε). *Let $L > 0$ be the constant in Assumption 2.13. For each $\varepsilon \in (0, L^{-1}/2)$, define Ξ_ε to be the set containing all $(\lambda, G) \in \mathbb{R} \times L^2([0, T]; \mathbb{R})$ such that $L^{-1} \leq \lambda \leq L$, $\|G\|_{L^2([0,T])} \leq L$, and*

$$\int_0^T \int_0^T G(|t-s|)f(s)f(t)dsdt \geq -\varepsilon\|f\|_{L^2([0,T])}^2, \quad \text{for all } f \in L^2([0, T], \mathbb{R}). \quad (2.23)$$

Remark 2.15. *Recall that by Theorem 2.10, the estimator (2.16) only approximates the true kernel G^* in the L^2 sense, and hence may not be non-negative definite. Thus, Definition 2.14 only requires the estimated kernel G to be nearly non-negative definite relative to the estimated λ , as reflected by (2.23) and $\varepsilon \in (0, L^{-1}/2)$. Since $\lambda^* > 0$ and G^* is non-negative definite (see (2.3)), this condition can be satisfied by estimated models with sufficiently many samples, as shown in Lemma 7.1.*

Definition 2.14 ensures that the greedy policy u^θ is well-defined for any admissible model $\theta \in \Xi_\varepsilon$. Moreover, Theorem 2.16 shows that the performance gap of the greedy policy of an estimated model depends quadratically on the model error. The proof of Theorem 2.16 is given in Section 5.

Theorem 2.16. *Let $\varepsilon \in (0, L^{-1}/2)$ with $L > 0$ as in Assumption 2.13. For each $\theta \in \Xi_\varepsilon$, let $u^\theta = \text{Greedy}(A, \theta)$ be defined by (2.8). Then there exists a constant $C > 0$, depending on L and ε , such that*

$$|J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^\theta)| \leq C \left(|\lambda^* - \lambda|^2 + \|G^* - G\|_{L^2([0,T])}^2 \right), \quad \text{for all } \theta, \theta' \in \Xi_\varepsilon.$$

Remark 2.17. *The performance gap in Theorem 2.16 relies on the Lipschitz stability of the optimal control u^θ with respect to the parameter θ , in particular, the kernel function G (see Proposition 5.3). In [5, 30, 48], Lipschitz stability of optimal controls has been derived for finite-dimensional parametric control problems, where the parameter varies in a compact subset of a finite-dimensional space. The stability analysis for Theorem 2.16 is more technically involved, as G takes value in the infinite-dimensional space $L^2([0, T]; \mathbb{R})$, and the control problem (2.7) is non-Markovian due to the kernel G . For instance, one can no longer prove the Lipschitz continuity of u^θ by simply arguing a continuous differentiability of $\theta \mapsto u^\theta$ as in [5, 48], since a bounded subset of $L^2([0, T]; \mathbb{R})$ may not be compact.*

To overcome these difficulties, we exploit an explicit representation of u^θ given in (4.13), and establish uniform bounds of a and B in (2.8) with suitable norms over all $\theta \in \Xi_\varepsilon$. These a-priori bounds further allow for proving the Lipschitz stability of the non-Markovian controls u^θ . Note that the explicit form of a and B is given in (4.13), and its main ingredients are operators and stochastic processes. We point out that as the running cost in (2.7) is not strongly concave with respect to the control variable, the control problem is an indefinite linear-quadratic control problem (see [53, Chapter 6] and references therein). Hence the condition (2.23) is essential for the well-definedness and the stability of u^θ .

Phased-based algorithm and its regret. The algorithm goes as follows. The input of the algorithm includes $\varepsilon > 0$ which satisfies Definition 2.14, and a deterministic exploration strategy $u^e \in L^2([0, T], \mathbb{R})$ as in Assumption 2.2. The algorithm starts with an initial exploration phase, where the agent exercises u^e for \mathbf{m}_0^e episodes, and forms an estimate θ^0 of θ^* according to (2.17):

$$\theta^0 = \text{LSE} \left((S^m, A^m)_{1 \leq m \leq \mathbf{m}_0^e}, \tau_{\mathbf{m}_0^e}, \pi_{\mathbf{m}_0^e} \right). \quad (2.24)$$

Here $\mathbf{m}_0^e \in \mathbb{N}$ is a prescribed number such that θ^0 is guaranteed to be in Ξ_ε (with high probability). Note that such initial exploration phase has a constant weight in the regret bounds, hence it has no impact on the results established in Theorem 2.19 and Corollary 2.21. One can alternatively use estimators based on historical datasets in order to get θ^0 which is in Ξ_ε . Some references for these preliminary estimates are available in [10, 50] and in Chapter 13 of [11]. Additional approach for getting θ^0 from historical trading data, uses the offline learning approach which was developed in [17]. Indeed a continuous version of Theorem 2.10 therein allows us to get a candidate for

G which satisfies the properties in Definition 2.14. Finally, as shown in [25, Example 2.7], any bounded, non-increasing and convex G satisfies (2.23) with $\varepsilon = 0$. If one assumes the true kernel G^* satisfies these shape constraints (as suggested by the empirical studies in [10, 42]), then one can enforce these constraints in the estimated kernels by minimising (2.16) over shaped constrained functions. This approach could potentially avoid the initial exploration phase and improve the sample complexity. These shape-constrained estimators have been analysed for discrete-time propagator models in [17], and extending the analysis therein to continuous-time propagator models is left for future work.

After this initial exploration, the algorithm then operates in cycles, and each cycle consists of exploitation and exploration phases. The exploitation phase of the k -th cycle, $k \in \mathbb{N}$, contains $\mathbf{n}(k)$ consecutive episodes for some prescribed $\mathbf{n}(k) \in \mathbb{N}$. At each exploitation episode, the agent executes the optimal strategy (2.8) defined using the current estimate θ^{k-1} and the signal trajectory observed in this episode. During the exploration phase of the k -th cycle, the agent exercises the exploration strategy u^e for one episode, and constructs an updated estimate θ^k by (2.17) using data from previous exploration episodes:

$$\theta^k = \text{LSE} \left((S^m, A^m)_{m \in \mathcal{E}_k}, \tau_{\mathbf{m}_0^e + k}, \pi_{\mathbf{m}_0^e + k} \right), \quad (2.25)$$

where $\mathcal{E}_k = \{1, \dots, \mathbf{m}_0^e\} \cup \{\mathbf{m}_0^e + \sum_{i=1}^j \mathbf{n}(i) + j \mid 1 \leq j \leq k\}$ is the indices of all exploration episodes up to the k -th cycle. This parameter θ^k will be used in the exploitation phase of the $(k+1)$ -th cycle. The algorithm is summarised as follows.

Algorithm 1: Phased-based learning algorithm

Input: $\varepsilon > 0$, $u^e \in L^2([0, T], \mathbb{R})$, $\mathbf{m}_0^e \in \mathbb{N}$ and $\mathbf{n} : \mathbb{N} \rightarrow \mathbb{N}$.

- 1 Execute u^e for \mathbf{m}_0^e episodes, and set $\theta^0 \in \Xi_\varepsilon$ as in (2.24).
- 2 **for** $k = 1, 2, \dots$ **do**
- 3 $\mathfrak{L}(k-1) = \mathbf{m}_0^e + \sum_{i=1}^{k-1} \mathbf{n}(i) + k - 1$. /* last episode's index */
- 4 **for** $m = \mathfrak{L}(k-1) + 1, \dots, \mathfrak{L}(k-1) + \mathbf{n}(k)$ **do**
- 5 Execute the greedy strategy $u^m = \text{Greedy}(A^m, \theta^{k-1})$.
- 6 **end**
- 7 Execute u^e for one episode, and set $\theta^k \in \Xi_\varepsilon$ as in (2.25).
- 8 **end**

Remark 2.18. *Algorithm 1 is an anytime algorithm, as it does not restrict the maximum number of learning episodes (see the last paragraph of Section 2.1). It distributes the exploration episodes over the whole learning process according to the schedulers \mathbf{m}_0^e and \mathbf{n} , which are chosen to optimise the regret order for all episodes. This should be in contrast to the setting where the algorithm termination time is fixed and known by the agent. In this case, the agent can put all exploration episodes at the beginning, whose number depends explicitly on the prescribed maximal episode number (see [34, Chapter 6]).*

Compared with the algorithm in [48], Algorithm 1 introduces an initial exploration step, and updates $(\theta^k)_{k \geq 0}$ using trajectories generated by a fixed exploration strategy u^e . This allows for applying Theorem 2.10 to ensure that, with high probabilities, $(\theta^k)_{k \geq 0}$ stay in Ξ_ε (without an explicit projection as in [48]) and converge to θ^* as $k \rightarrow \infty$.

The following theorem chooses the learning schedulers $\mathbf{m}_0^e \in \mathbb{N}$ and $\mathbf{n} : \mathbb{N} \rightarrow \mathbb{N}$ such that Algorithm 1 achieves sublinear regrets in high probability. These hyper-parameters are optimised depending on the convergence rate of the regularised least-squares estimator (2.17) in Theorem 2.10. The proof of Theorem 2.19 is given in Section 7.

Theorem 2.19. *Suppose that the parameters $(\tau_N, \pi_N)_{N \in \mathbb{N}}$ for (2.17) are chosen such that $(\theta^N)_{N \in \mathbb{N}}$ satisfies the following error estimate: there exists $\tilde{C} > 0$ and $\kappa \in (0, 1)$ such that for all $\eta \in (0, 1)$, it holds with probability at least $1 - \eta$ that, for all $N \in \mathbb{N} \cap [2, \infty)$,*

$$|\lambda^N - \lambda^*| \leq \tilde{C} \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right), \quad \|G^N - G^*\|_{L^2([0, T])} \leq \tilde{C} \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^\kappa. \quad (2.26)$$

Assume further that Assumption 2.13 holds and let $\varepsilon \in (0, L^{-1}/2)$ as in Definition 2.14. Then there exists $C_0 > 0$ such that for all $\eta \in (0, 1)$ and $C \geq C_0$, if one sets $\mathbf{m}_0^e = \lceil C(\log(\eta^{-1})^2 + 1) \rceil$, and $\mathbf{n} : \mathbb{N} \rightarrow \mathbb{N}$ such that $\mathbf{n}(k) = \lfloor k^\kappa \rfloor$ for all $k \in \mathbb{N}$, then with probability at least $1 - \eta$, the regret of Algorithm 1 satisfies

$$R(N) \leq C' \left(N^{\frac{1}{1+\kappa}} (\log(\eta^{-1}) + \log N)^{2\kappa} + \log(\eta^{-1})^2 \right), \quad \text{for all } N \in \mathbb{N} \cap [2, \infty),$$

where $C' > 0$ is a constant independent of η and N .

As Algorithm 1 operates in cycles, for each $N \in \mathbb{N}$, the regret of learning $R(N)$ up to N episodes can be upper bounded by the accumulated regret at the end of the K -th cycle, with $K = \min\{k \in \mathbb{N} \cup \{0\} \mid \mathfrak{L}(k) = \mathbf{m}_0^e + \sum_{i=1}^k \mathbf{n}(i) + k \geq N\}$.

Remark 2.20. *If (2.26) holds with $\kappa = 1$, then Theorem 2.19 recovers the square-root regret bound in [48] for linear-convex RL problems with finite-dimensional unknown parameters. However, in the present non-parametric setting, (2.26) typically holds with $\kappa < 1$ (see Theorems 2.10), and this leads to a worse sublinear regret bound. Indeed, classical results for Tikhonov regularization indicate that $\kappa = 2/3$ is the best rate one can expect, even for a smooth kernel G^* (see [28, Section 3.2]). Employing other regularisation approaches to improve the sample efficiency of the kernel estimation is left for future research.*

By combining Theorems 2.10 and 2.19, the following corollary optimises the regret bounds of Algorithm 1 depending on the regularity of the true kernel G^* .

Corollary 2.21. *Suppose that Assumptions 2.2, 2.8 and 2.13 hold. Let $\varepsilon \in (0, L^{-1}/2)$ as in Definition 2.14.*

- (1) *If Assumption 2.6(1) holds, then there exists $C_0 > 0$ such that for all $\eta \in (0, 1)$ and $C \geq C_0$, by setting $(\tau_N, \pi_N)_{N \in \mathbb{N}}$ as (2.19) for (2.17), $\mathbf{m}_0^e = \lceil C(\log(\eta^{-1})^2 + 1) \rceil$, and $\mathbf{n} : \mathbb{N} \rightarrow \mathbb{N}$ such that $\mathbf{n}(k) = \lfloor k^{\frac{1}{3}} \rfloor$ for all $k \in \mathbb{N}$, then with probability at least $1 - \eta$, the regret of Algorithm 1 satisfies for all $N \in \mathbb{N} \cap [2, \infty)$,*

$$R(N) \leq C' \left(N^{\frac{3}{4}} (\log(\eta^{-1}) + \log N)^{\frac{2}{3}} + \log(\eta^{-1})^2 \right). \quad (2.27)$$

- (2) *If Assumption 2.6(2) holds, then there exists $C_0 > 0$ such that for all $\eta \in (0, 1)$ and $C \geq C_0$, by setting $(\tau_N, \pi_N)_{N \in \mathbb{N}}$ as (2.21) for (2.17), $\mathbf{m}_0^e = \lceil C(\log(\eta^{-1})^2 + 1) \rceil$, and $\mathbf{n} : \mathbb{N} \rightarrow \mathbb{N}$ such that $\mathbf{n}(k) = \lfloor k^{\frac{1-2\alpha}{3-2\alpha}} \rfloor$ for all $k \in \mathbb{N}$, then with probability at least $1 - \eta$, the regret of Algorithm 1 satisfies for all $N \in \mathbb{N} \cap [2, \infty)$,*

$$R(N) \leq C' \left(N^{\frac{3-2\alpha}{4-4\alpha}} (\log(\eta^{-1}) + \log N)^{\frac{2(1-2\alpha)}{3-2\alpha}} + \log(\eta^{-1})^2 \right). \quad (2.28)$$

The constant $C' > 0$ appearing in (2.27) and (2.28) is independent of η and N .

3 Numerical implementation

In this section, we numerically examine the performance of the least-squares estimator (2.17) developed in Section 2.2, which is the key component of the phased-based learning algorithm (Algorithm 1). We focus on estimating singular power law propagators, which are extensively used by practitioners (see e.g. Chapter 13 of [10]).

More precisely, let $T > 0$ and consider $\lambda^* > 0$ to be an unknown temporary price impact coefficient, and $G^* \in L^2([0, T], \mathbb{R})$ to be an unknown transient impact kernel. For each $n \in \mathbb{N}$, consider the following price process (cf. (2.9)):

$$S_t^n = A_t^n + M_t^n - \lambda^* u_t^e - \int_0^t G^*(t-s) u_s^e ds, \quad t \in [0, T], \quad (3.1)$$

where $u^e \in L^2([0, T], \mathbb{R})$ be a trading strategy specified by the agent, $(A^n)_{n \in \mathbb{N}}$ are observed signals, and $(M^n)_{n \in \mathbb{N}}$ are unobserved zero-mean noises. The agent estimates (λ^*, G^*) based on the observed trajectories $(S^n, A^n)_{n \in \mathbb{N}}$ and the trading strategy u^e .

Given sample trajectories $(S^n, A^n)_{n=1}^N$, we estimate the parameter λ^* by λ^N defined in (2.14), and estimate the kernel G^* by

$$\begin{aligned} G^N := \arg \min_{G \in V_N} & \left(\frac{1}{N} \sum_{n=1}^N \int_0^T \left| S_t^n - A_t^n + \lambda^N u_t^e + \int_0^t u^e(t-s) G_s ds \right|^2 dt \right. \\ & \left. + \tau_N \int_0^T (G_t - H_t)^2 dt \right). \end{aligned} \quad (3.2)$$

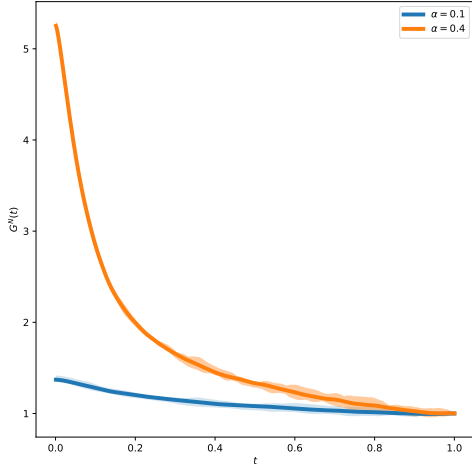
The estimator (3.2) extends (2.16) by allowing for a (non-zero) initial guess H of G^* in the regularisation term. It satisfies the convergence rates in Theorem 2.10 provided that $G^* - H$ satisfies Assumption 2.6 (see e.g., [46]). The estimator (3.2) can be numerically approximated by

$$G^N = \arg \min_{(G_k)_{k=0}^{K-1}} \left(\frac{1}{N} \sum_{n=1}^N \sum_{j=1}^K \left| S_{t_j}^n - A_{t_j}^n + \lambda^N u_{t_j}^e + \sum_{k=0}^{j-1} u_{t_j-k}^e G_k \Delta t \right|^2 \Delta t + \tau_N \sum_{k=0}^{K-1} (G_k - H_k)^2 \Delta t \right), \quad (3.3)$$

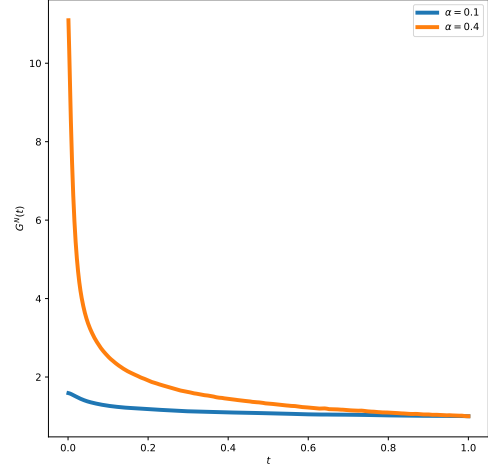
where $\Delta t = T/K$ and $t_i = iT/K$ for a sufficiently large $K \in \mathbb{N}$. Note that the minimiser of (3.3) can be computed analytically by a first-order condition.

For our numerical experiments, we fix $\lambda^* = 0.5$ and $G^*(t) = t^{-\alpha}$ for some $\alpha \in (0, 1/2)$. For each $N \in \mathbb{N}$, we generate observed trajectories $(S^m - A^m)_{m=1}^N$ according to (3.1) with $T = 1$, $u^e \equiv 1$ and $M_t^n = 0.5(B_t^n + \iota_0)$, $t \in [0, T]$, where $(B^n)_{n=1}^N$ are independent Brownian motions and ι_0 is an independent standard normal random variable. Using these trajectories, we evaluate λ^N as in (2.14) and G^N as in (3.3) with $H \equiv 1$, $K = 10^3$ and $\tau_N = N^{-2/(3-2\alpha)}$ as suggested by (2.21). To estimate statistical properties of the estimators, we carry out the experiments for 10 independent runs, where among different executions, the observed state trajectories are simulated based on independent noises. In the sequel, we only report the performance of G^N , as G^* is more challenging to estimate than λ^* .

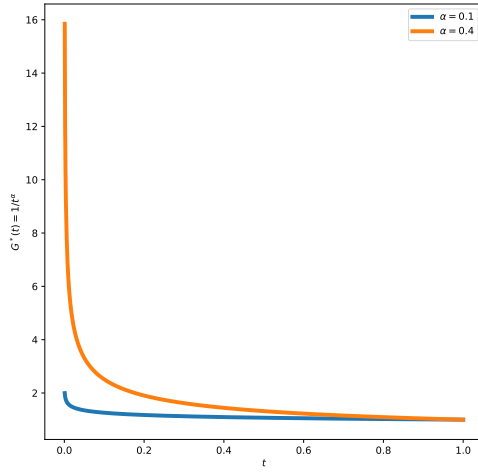
Figure 1 illustrates the effectiveness of G^N for estimating the power law kernel $G^*(t) = t^{-\alpha}$ with different $\alpha \in \{0.1, 0.4\}$ and different sample sizes N . Figure 1c plots the true kernels, showing that the degree of singularity at $t = 0$ increases for larger α . Figures 1a and 1b display the estimated kernels, where the solid lines represent the mean and the shaded areas indicate the extremes over 10 repeated experiments. One can see that the estimator G^N successfully recovers the overall behavior of the true kernels, even with a small sample size. The regularisation τ_N prevents oscillation in the estimated kernel, in contrast to an unregularised estimator as discussed in Remark 2.5. However, the estimator is not able to accurately capture the singularities of the true kernels without a sufficient number of samples. By increasing the sample sizes, these singularities can be better captured by the estimator.



(a) Estimated kernels with $N = 1024$



(b) Estimated kernels with $N = 65526$



(c) True kernels

Figure 1: Comparison between the true power law kernels $G^*(t) = t^{-\alpha}$, with $\alpha = 0.1$ (in blue) and $\alpha = 0.4$ (in orange), and the estimated kernels with different sample sizes N .

We further demonstrate the decay of the relative error $\|G^N - G^*\|_{L^2([0,T])} / \|G^*\|_{L^2([0,T])}$ in Figure 2 for sample sizes $N = 2^n$ with $n \in \{10, 11, \dots, 16\}$. The estimator makes larger errors for $\alpha = 0.4$ compared to $\alpha = 0.1$, due to the more severe singularity of the true kernel. A linear regression on the logarithms of relative errors and sample sizes reveals that the convergence rate of $(G^N)_{N \in \mathbb{N}}$ is of the order $\mathcal{O}(N^{-0.22})$ for $\alpha = 0.1$ and $\mathcal{O}(N^{-0.2})$ for $\alpha = 0.4$, which are better than the theoretical upper bounds presented

in Theorem 2.10. Note that this improved rate does not contradict the theorem, as the convergence rates in Theorem 2.10 pertain to the worst-case scenario over all kernels G^* satisfying Assumption 2.6 and all realizations of the noises M satisfying Assumption 2.8. As already pointed out in [19, Section 3.2], for particular realised sample trajectories, the error of the estimator (2.16) (or (3.3)) may be smaller than the bounds provided in Theorem 2.10.

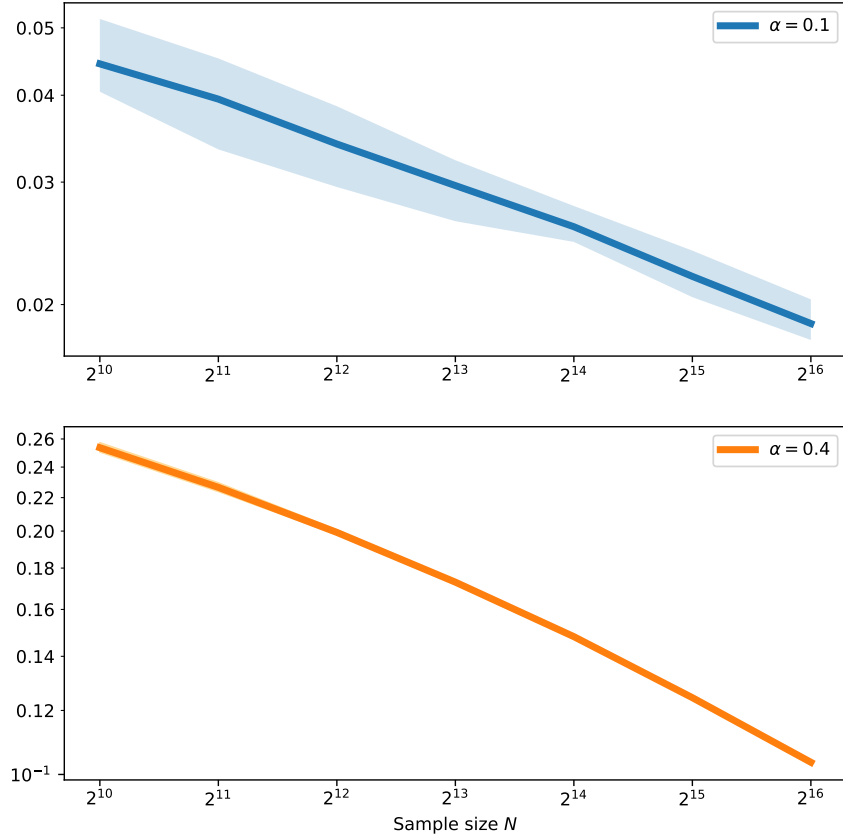


Figure 2: Mean relative errors of G^N for different sample sizes plotted in solid lines and the intervals containing the errors denoted by the lighter regions (the plot is in a log-log scale). The true power law kernels $G^*(t) = t^{-\alpha}$, with $\alpha = 0.1$ (upper panel) and $\alpha = 0.4$ (lower panel).

It is worth noting that in this experiment, we focus on the challenging task of estimating singular kernels, for which the proposed estimators are expected to produce larger estimation errors compared to more regular kernels. However, as shown in Theorem 2.16, the error of the greedy policy for an estimated model depends quadratically on the kernel estimation error. Therefore, a rough estimate of the kernel G is

often sufficient to design a nearly optimal trading strategy. In particular, as shown in Figure 2, even for the most singular kernel with $\alpha = 0.4$, the relative error of the resulting greedy policy is expected to be around 6% with a sample size of $N = 1024$.

4 Analytic solution to the control problem

In this section, we recall the explicit form on the optimiser of (2.7) from [1]. Before stating this result we introduce some essential definitions of function spaces, integral operators and stochastic processes.

4.1 Function spaces, integral operators

Let $T > 0$. We denote by $\langle \cdot, \cdot \rangle_{L^2}$ the inner product on $L^2([0, T], \mathbb{R}^2)$, that is

$$\langle f, g \rangle_{L^2} = \int_0^T f(s)^\top g(s) ds, \quad f, g \in L^2([0, T], \mathbb{R}^2). \quad (4.1)$$

We define $L^2([0, T]^2, \mathbb{R}^{2 \times 2})$ to be the space of measurable kernels $\Sigma : [0, T]^2 \rightarrow \mathbb{R}^2$ such that

$$\int_0^T \int_0^T |\Sigma(t, s)|^2 dt ds < \infty.$$

The notation $|\cdot|$ stands for a matrix norm, and in particular we have,

$$\int_0^T \int_0^T |\Sigma_{i,j}(t, s)|^2 dt ds < \infty, \quad \text{for all } i, j = 1, 2.$$

For any $\Sigma, \Lambda \in L^2([0, T]^2, \mathbb{R}^{2 \times 2})$ we define the \star -product as follows,

$$(\Sigma \star \Lambda)(s, u) = \int_0^T \Sigma(s, z) \Lambda(z, u) dz, \quad (s, u) \in [0, T]^2, \quad (4.2)$$

which is a well-defined kernel in $L^2([0, T]^2, \mathbb{R}^{2 \times 2})$ due to Cauchy-Schwarz inequality. For any kernel $\Sigma \in L^2([0, T]^2, \mathbb{R}^{2 \times 2})$, we denote by Σ the integral operator induced by the kernel Σ that is

$$(\Sigma g)(s) = \int_0^T \Sigma(s, u) g(u) du, \quad g \in L^2([0, T], \mathbb{R}^2). \quad (4.3)$$

Σ is a linear bounded operator from $L^2([0, T], \mathbb{R}^2)$ into itself. For Σ and Λ that are two integral operators induced by the kernels Σ and Λ in $L^2([0, T]^2, \mathbb{R}^{2 \times 2})$, we denote by $\Sigma \Lambda$ the integral operator induced by the kernel $\Sigma \star \Lambda$.

We denote by Σ^* the adjoint kernel of Σ for $\langle \cdot, \cdot \rangle_{L^2}$, that is

$$\Sigma^*(s, u) = \Sigma(u, s)^\top, \quad (s, u) \in [0, T]^2, \quad (4.4)$$

and by Σ^* the corresponding adjoint integral operator.

We recall that an operator Σ as above is said to be non-negative definite if $\langle \Sigma f, f \rangle_{L^2} \geq 0$ for all $f \in L^2([0, T], \mathbb{R}^2)$. It is said to be positive definite if $\langle \Sigma f, f \rangle_{L^2} > 0$ for all $f \in L^2([0, T], \mathbb{R}^2)$ not identically zero.

4.2 Essential operators and processes

The Γ_t^{-1} operator: We define

$$\tilde{G}(t, s) = (2\varrho + G(t - s)) \mathbb{1}_{\{s < t\}}, \quad 0 \leq s, t \leq T. \quad (4.5)$$

and $\tilde{\mathbf{G}}_t$ as the operator induced by the kernel

$$\tilde{G}_t(s, u) = \tilde{G}(s, u) \mathbb{1}_{\{u \geq t\}}. \quad (4.6)$$

We introduce

$$\mathbf{D}_t := 2\lambda \text{id} + (\tilde{\mathbf{G}}_t + \tilde{\mathbf{G}}_t^*) + 2\phi \mathbf{1}_t^* \mathbf{1}_t, \quad (4.7)$$

where id is the identity operator, i.e. $(\text{id}f)(t) = f(t)$, $\mathbf{1}_t$ is the integral operator induced by the kernel

$$\mathbb{1}_t(u, s) := \mathbb{1}_{\{u \geq s\}} \mathbb{1}_{\{s \geq t\}}. \quad (4.8)$$

In Lemma 3.1 of [1] it was proved that \mathbf{D}_t is invertible if $\lambda > 0$ and $\varrho, \phi \geq 0$. This will be necessary for upcoming definitions. We define the operator $\mathbf{\Gamma}_t^{-1}$ by

$$\mathbf{\Gamma}_t^{-1} = \begin{pmatrix} \mathbf{D}_t^{-1} & -2\phi \mathbf{D}_t^{-1} \mathbf{1}_t^* \\ -2\phi \mathbf{1}_t \mathbf{D}_t^{-1} & -2\phi \text{id} + 4\phi^2 \mathbf{1}_t \mathbf{D}_t^{-1} \mathbf{1}_t^* \end{pmatrix}. \quad (4.9)$$

We recall that $\mathbf{\Gamma}^{-1}$ is associated with a solution to an operator Riccati equation (see Lemmas 6.1 and 6.2 in [1]).

The process Θ : For convenience we introduce the following notation,

$$\mathbb{1}_t(s) = \mathbb{1}_{\{s \geq t\}}. \quad (4.10)$$

For A as in (2.1), we define $\Theta = \{\Theta_t(s) : t \in [0, s], s \in [0, T]\}$ as follows,

$$\Theta_t(s) = - \left(\mathbf{\Gamma}_t^{-1} \mathbb{1}_t \mathbb{E} \left[A - A_T \mid \mathcal{F}_t \right] e_1 \right) (s), \quad \text{with } e_1 = (1, 0)^\top. \quad (4.11)$$

Note that Θ solves the L^2 -valued BSDE (see Proposition 6.3 in [1]).

The optimal control u^* : Proposition 4.5 in [1] states that the optimiser of (2.7), u^* is a solution to the following equation,

$$u_t^* = a_t + \int_0^t B(t, s) u_s^* ds, \quad (4.12)$$

where the process $\{a_t\}_{t \in [0, T]}$ and the kernel B which are given by

$$\begin{aligned} a_t &= \frac{1}{2\lambda} \left(\mathbb{E}[A_t - A_T \mid \mathcal{F}_t] + 2\varrho q + \langle \Theta_t, K_t \rangle_{L^2} + \langle \Gamma_t^{-1} K_t, \mathbf{1}_t(-2\varrho q, q)^\top \rangle_{L^2} \right), \\ B(t, s) &= \mathbf{1}_{\{s < t\}} \frac{1}{2\lambda} \left(\langle \Gamma_t^{-1} K_t, \mathbf{1}_t(\tilde{G}(\cdot, s), -1)^\top \rangle_{L^2} - \tilde{G}(t, s) \right). \end{aligned} \quad (4.13)$$

Here the function K_t is defined as follows

$$K_t(s) = (\tilde{G}(t, s), -\mathbf{1}_{\{s \leq t\}})^\top. \quad (4.14)$$

5 Bound on the performance gap

In this section we prove Theorem 2.16. Recall that the parameter space Ξ_ε was defined in (2.23) and that u^{θ^*} is the maximiser of J^{θ^*} in (2.7). In the following proposition we derive an upper bound on $J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u)$ for any admissible strategy u .

For any admissible strategy u as in (2.5) we define

$$\|u\|_{\mathcal{H}^2} = \left(\mathbb{E} \left[\int_0^T u_s^2 ds \right] \right)^{1/2}. \quad (5.1)$$

Proposition 5.1. *Let $\theta^* \in \Xi_\varepsilon$, then there exists a constant $C > 0$ such that*

$$0 \leq J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u) \leq C \|u^{\theta^*} - u\|_{\mathcal{H}^2}^2, \quad \text{for all } u \in \mathcal{A}.$$

We start with characterising the optimal strategy by using a variational approach. Note that for any $u \in \mathcal{A}$ the map $u \mapsto J(u)$ in (2.7) is strictly concave, which can be easily shown by repeating the same lines as in the proof of Lemma 9.1 of [41]. Therefore, it admits a unique maximizer characterized by the critical point at which the Gâteaux derivative

$$\langle (J^{\theta^*}(u))', \alpha \rangle \triangleq \lim_{\varepsilon \rightarrow 0} \frac{J^{\theta^*}(u + \varepsilon \alpha) - J^{\theta^*}(u)}{\varepsilon} \quad (5.2)$$

of the functional J^{θ^*} vanishes for any direction $\alpha = (\alpha_t)_{0 \leq t \leq T} \in \mathcal{A}$; see, e.g., [18]. The Gâteaux derivative in (5.2) can be derived along the same line as in [40]. The result is given in the following lemma.

Lemma 5.2. *Let J^{θ^*} as in (2.7). Let $u \in \mathcal{A}$, then we have*

$$\begin{aligned} \langle J^{\theta^*}(u)', \alpha \rangle = \mathbb{E} & \left[\int_0^T \alpha_s \left(P_s - Z_s^{\theta^*, u} - \int_s^T G^*(t-s) u_t dt - 2\lambda^* u_s \right. \right. \\ & \left. \left. + 2\phi \int_s^T Q_t^u dt + 2\varrho Q_T^u - P_T \right) ds \right], \end{aligned} \quad (5.3)$$

for any $\alpha \in \mathcal{A}$.

Since we are maximizing the strictly concave functional $u \mapsto J^{\theta^*}(u)$ over \mathcal{A} , a necessary and sufficient condition for the optimality of $u^{\theta^*} \in \mathcal{A}$ is given by

$$\langle J^{\theta^*}(u^{\theta^*})', \alpha \rangle = 0 \quad \text{for all } \alpha \in \mathcal{A}; \quad (5.4)$$

see e.g., [18]. Now we are ready to prove Proposition 5.1.

Proof of Proposition 5.1. From (2.6) and Fubini's theorem we get that

$$\begin{aligned} & \int_0^T (Z_t^{\theta^*, u} u_t - Z_t^{\theta^*, u^{\theta^*}} u_t^{\theta^*}) dt \\ &= \int_0^T (u_t - u_t^{\theta^*}) (Z_t^{\theta^*, u} - Z_t^{\theta^*, u^{\theta^*}}) dt \\ & \quad + \int_0^T (u_t - u_t^{\theta^*}) Z_t^{\theta^*, u^{\theta^*}} dt + \int_0^T u_t^{\theta^*} (Z_t^{\theta^*, u} - Z_t^{\theta^*, u^{\theta^*}}) dt. \\ &= \int_0^T (u_t - u_t^{\theta^*}) (Z_t^{\theta^*, u} - Z_t^{\theta^*, u^{\theta^*}}) dt + \int_0^T (u_t - u_t^{\theta^*}) Z_t^{u^{\theta^*, \theta^*}} dt \\ & \quad + \int_0^T (u_s - u_s^{\theta^*}) \int_s^T G^*(t-s) u_t^{\theta^*} dt ds. \end{aligned} \quad (5.5)$$

From (2.6) and application of Cauchy-Schwarz inequality twice we get

$$\begin{aligned} & \mathbb{E} \left[\left| \int_0^T (u_t - u_t^{\theta^*}) (Z_t^{\theta^*, u} - Z_t^{\theta^*, u^{\theta^*}}) dt \right| \right] \\ & \leq \left(\mathbb{E} \left[\int_0^T (Z_t^{\theta^*, u} - Z_t^{\theta^*, u^{\theta^*}})^2 dt \right] \right)^{1/2} \left(\mathbb{E} \left[\int_0^T (u_t - u_t^{\theta^*})^2 dt \right] \right)^{1/2} \\ & \leq \|u^{\theta^*} - u\|_{\mathcal{H}^2} \left(\mathbb{E} \left[\int_0^T \left(\int_0^t G^*(t-s) (u_s - u_s^{\theta^*}) ds \right)^2 dt \right] \right)^{1/2} \\ & \leq C (\|u^{\theta^*} - u\|_{\mathcal{H}^2})^2 \left(\int_0^T (G^*(t-s))^2 ds \right)^{1/2} \\ & \leq C \|u^{\theta^*} - u\|_{\mathcal{H}^2}^2, \end{aligned} \quad (5.6)$$

where we used Definition 2.14 in the last inequality.

Using (2.4) and Jensen's inequality we get that

$$\begin{aligned}\mathbb{E}[(Q_T^u - Q_T^{u^{\theta^*}})^2] &\leq C\mathbb{E}\left[\int_0^T (u_t^{\theta^*} - u_t)^2 dt\right] \\ &\leq C\|u^{\theta^*} - u\|_{\mathcal{H}^2}^2.\end{aligned}\tag{5.7}$$

Similarly we have

$$\mathbb{E}\left[\int_0^T (Q_t^u - Q_t^{u^{\theta^*}})^2 dt\right] \leq \|u^{\theta^*} - u\|_{\mathcal{H}^2}^2.\tag{5.8}$$

From (2.7) and (5.5) we therefore get for any $u \in \mathcal{A}$ that

$$\begin{aligned}&J^{\theta^*}(u) - J^{\theta^*}(u^{\theta^*}) \\ &= \mathbb{E}\left[\int_0^T P_t(u_t - u_t^{\theta^*})dt - \int_0^T (Z_t^{\theta^*,u}u_t - Z_t^{\theta^*,u^{\theta^*}}u_t^{\theta^*})dt \right. \\ &\quad - \lambda \int_0^T (u_t^2 - (u_t^{\theta^*})^2)dt - \phi \int_0^T ((Q_t^u)^2 - (Q_t^{u^{\theta^*}})^2)dt + (Q_T^u - Q_T^{u^{\theta^*}})P_T \\ &\quad \left. - \varrho((Q_T^u)^2 - (Q_T^{u^{\theta^*}})^2)\right] \\ &= \mathbb{E}\left[- \int_0^T (Z_t^{\theta^*,u} - Z_t^{\theta^*,u^{\theta^*}})(u_t - u_t^{\theta^*})dt - \lambda \int_0^T (u_t^{\theta^*} - u_t)^2 dt \right. \\ &\quad \left. - \phi \int_0^T (Q_t^u - Q_t^{u^{\theta^*}})^2 dt - \varrho(Q_T^u - Q_T^{u^{\theta^*}})^2\right] \\ &\quad + \mathbb{E}\left[\int_0^T P_t(u_t - u_t^{\theta^*})dt - \int_0^T (u_t - u_t^{\theta^*})Z_t^{\theta^*,u^{\theta^*}} dt \right. \\ &\quad - \int_0^T (u_s - u_s^{\theta^*}) \int_s^T G^*(t-s)u_t^{\theta^*} dt ds - 2\lambda \int_0^T u_t^{\theta^*}(u_t - u_t^{\theta^*})dt \\ &\quad + (Q_T^{u^{\theta^*}} - Q_T^u)P_T + 2\phi \int_0^T Q_t^{u^{\theta^*}}(Q_t^u - Q_t^{u^{\theta^*}})dt \\ &\quad \left. + 2\varrho Q_T^{u^{\theta^*}}(Q_T^u - Q_T^{u^{\theta^*}})\right].\end{aligned}\tag{5.9}$$

Since $u_t - u_t^{\theta^*} \in \mathcal{A}$ it follows from (2.4), (5.3) and (5.4) that

$$\begin{aligned} & J^{\theta^*}(u) - J^{\theta^*}(u^{\theta^*}) \\ &= \mathbb{E} \left[- \int_0^T (Z_t^{\theta^*, u} - Z_t^{\theta^*, u^{\theta^*}})(u_t - u_t^{\theta^*}) dt - \lambda \int_0^T (u_t^{\theta^*} - u_t)^2 dt \right. \\ & \quad \left. - \phi \int_0^T (Q_t^u - Q_t^{u^{\theta^*}})^2 dt - \varrho(Q_T^u - Q_T^{u^{\theta^*}})^2 \right]. \end{aligned} \quad (5.10)$$

Together with (5.6) and (5.7) we get that

$$|J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u)| \leq C \|u^{\theta^*} - u\|_{\mathcal{H}^2}^2,$$

which completes the proof. \square

The following proposition proves the stability of u^{θ^*} with respect to the parameter θ^* . The proof of Proposition 5.3 is postponed to Section 8.

Proposition 5.3. *For each $\theta \in \Xi_\varepsilon$ let u^θ be defined as in (2.8). Then, there exists a constant $C > 0$ such that*

$$\|u^\theta - u^{\theta'}\|_{\mathcal{H}^2} \leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}), \quad \text{for all } \theta, \theta' \in \Xi_\varepsilon.$$

Now we are ready to prove Theorem 2.16.

Proof of Theorem 2.16. The proof of Theorem 2.16 follows immediately from Propositions 5.1 and 5.3. \square

6 Analysis of regularised least-squares estimators

This section quantifies the convergence rate of $(\theta^N)_{N \in \mathbb{N}}$ in (2.17) in high probability, and hence proves Theorem 2.10. To simplify the notation, for each $f \in H^1([0, T], \mathbb{R})$, we denote by \dot{f} the derivative of f .

The following lemma proves the injectivity of \mathbf{u}^e in (2.13), and further characterises the range of the adjoint operator $(\mathbf{u}^e)^*$.

Lemma 6.1. *Suppose Assumption 2.2 holds. Then $\mathbf{u}^e : L^2([0, T], \mathbb{R}) \rightarrow L^2([0, T], \mathbb{R})$ is injective and compact. Moreover, the range of the adjoint operator $(\mathbf{u}^e)^*$ is given by $\mathcal{R}((\mathbf{u}^e)^*) = \{f \in H^1([0, T], \mathbb{R}) \mid f(T) = 0\}$.*

Proof. The compactness of \mathbf{u}^e follows from the fact that it is a Hilbert-Schmidt integral operator. Next we show that \mathbf{u}^e is injective. Let $f \in L^2([0, T], \mathbb{R})$ such

that $(\mathbf{u}^e f)(t) = 0$ for a.e. $t \in [0, T]$. By (2.13) and the Leibniz integral rule, for a.e. $t \in [0, T]$,

$$\begin{aligned} 0 &= \frac{d}{dt} \left(\int_0^t u^e(t-s)f(s)ds \right) = u^e(0)f(t) + \int_0^t \dot{u}^e(t-s)f(s)ds \\ &= ((u^e(0)\text{id} + \dot{\mathbf{u}}^e)f)(t), \end{aligned}$$

where $\dot{\mathbf{u}}^e$ is the Volterra operator induced by the square-integrable kernel $[0, T]^2 \ni (t, s) \mapsto \dot{u}^e(t-s)\mathbf{1}_{\{s < t\}} \in \mathbb{R}$. This along with [27, Chapter 9, Corollary 3.16] and $u^e(0) \neq 0$ yields $f = 0$. This proves the injectivity of \mathbf{u}^e .

It remains to characterise the range of $(\mathbf{u}^e)^*$. By (2.13) and Fubini's theorem, the adjoint of \mathbf{u}^e satisfies for all $f \in L^2([0, T], \mathbb{R})$ and a.e. $t \in [0, T]$, $((\mathbf{u}^e)^* f)(t) = \int_t^T u^e(s-t)f(s)ds$, and hence $((\mathbf{u}^e)^* f)(T) = 0$. By Assumption 2.2 and the Leibniz integral rule, for all $f \in L^2([0, T], \mathbb{R})$ and for a.e. $t \in [0, T]$,

$$\begin{aligned} \frac{d}{dt} ((\mathbf{u}^e)^* f)(t) &= \frac{d}{dt} \left(\int_t^T u^e(s-t)f(s)ds \right) = -u^e(0)f(t) - \int_t^T \dot{u}^e(s-t)f(s)ds \\ &= -((u^e(0)\text{id} + \dot{\mathbf{u}}^e)^* f)(t), \end{aligned} \tag{6.1}$$

where $(\dot{\mathbf{u}}^e)^*$ is the adjoint of $\dot{\mathbf{u}}^e$. This along with the integrability of f implies that $\frac{d}{dt} ((\mathbf{u}^e)^* f) \in L^2([0, T], \mathbb{R})$ and hence $\mathcal{R}((\mathbf{u}^e)^*) \subset \{f \in H^1([0, T], \mathbb{R}) \mid f(T) = 0\}$. On the other hand, let $f \in H^1([0, T], \mathbb{R})$ be such that $f(T) = 0$. Then by (6.1), [27, Chapter 9, Corollary 3.16] and $u^e(0) \neq 0$, there exists $g \in L^2([0, T], \mathbb{R})$ such that $\frac{d}{dt} ((\mathbf{u}^e)^* g) = f$, which along with $f(T) = ((\mathbf{u}^e)^* g)(T) = 0$ implies that $f \equiv (\mathbf{u}^e)^* g$. This shows $\{f \in H^1([0, T], \mathbb{R}) \mid f(T) = 0\} \subset \mathcal{R}((\mathbf{u}^e)^*)$ and finishes the proof. \square

To analyse the convergence rate of (2.17), we interpret (2.16) as a Tikhonov regularisation for (2.12). We first adapt the theoretical framework of Tikhonov regularisation for (deterministic) linear inverse problems to the present setting. Let $(X, \|\cdot\|_X)$ be a Hilbert space, and let $K : X \rightarrow X$ be an injective compact linear operator with a non-closed range $\mathcal{R}(K)$. Let $x_0, y_0 \in X$ be such that $y_0 = Kx_0$, let $y^\delta \in X$ be a noisy observation of y_0 , and consider approximations of y_0 via the regularised Ritz approach. Specifically, let $(V_m)_{m \in \mathbb{N}} \subset X$ be a sequence of finite-dimensional subspaces such that $\bigcup_{m=1}^\infty V_m = X$. For each $\alpha > 0$ and $m \in \mathbb{N}$, let $x_{\alpha, m}^\delta$ be the unique minimiser of the following Tikhonov functional:

$$x_{\alpha, m}^\delta = \arg \min_{x \in V_m} (\|Kx - y^\delta\|_Y^2 + \alpha \|x\|_X^2) = (K_m^* K_m + \alpha \text{id})^{-1} K_m^* y^\delta, \tag{6.2}$$

where $K_m = KP_m$ and P_m is the orthogonal projection of X onto V_m .

The following lemma quantifies $\|x_{\alpha, m}^\delta - x_0\|_X$ in terms of α, m and the error $\|y^\delta - y_0\|_X$. The proof essentially combines the results in [28, 33], and is presented below for the reader's convenience.

Lemma 6.2. For each $\alpha > 0$ and $m \in \mathbb{N}$, let $x_{\alpha,m}^\delta \in X$ be defined by (6.2), and for each $R > 0$, let $\mathcal{D}(R) = \inf\{\|x_0 - K^*v\|_X \mid v \in X, \|v\|_X \leq R\}$. Then for all $\alpha, R > 0$ and $m \in \mathbb{N}$,

$$\|x_{\alpha,m}^\delta - x_0\|_X \leq \frac{\|y_0 - y^\delta\|_X}{2\sqrt{\alpha}} + \gamma_m \sqrt{1 + \frac{\gamma_m^2}{\alpha}} \left(R + \frac{1}{2\sqrt{\alpha}} \mathcal{D}(R) \right) + \mathcal{D}(R) + \frac{\sqrt{\alpha}}{2} R,$$

where $\gamma_m = \|(\text{id} - P_m)K^*\|_{\text{op}}$.

Proof. For each $m \in \mathbb{N}$ and $\alpha > 0$, let $x_{\alpha,m} = (K_m^*K_m + \alpha \text{id})^{-1}K_m^*y_0$ and $x_\alpha = (K^*K + \alpha \text{id})^{-1}K^*y_0$. Then for all $m \in \mathbb{N}$ and $\alpha > 0$,

$$\begin{aligned} \|x_{\alpha,m}^\delta - x_0\|_X &\leq \|x_{\alpha,m}^\delta - x_{\alpha,m}\|_X + \|x_{\alpha,m} - x_\alpha\|_X + \|x_\alpha - x_0\|_X \\ &\leq \frac{\|y_0 - y^\delta\|_X}{2\sqrt{\alpha}} + \sqrt{1 + \frac{\gamma_m^2}{\alpha}} \|(\text{id} - P_m)x_\alpha\|_X + \mathcal{D}(R) + \frac{\sqrt{\alpha}}{2} R, \end{aligned} \quad (6.3)$$

where for the second inequality in (6.3) is derived as follows: the first term used the spectral inequality $\|(A^*A + \alpha \text{id})^{-1}A^*\|_{\text{op}} \leq \sup_{\lambda \geq 0} \frac{\sqrt{\lambda}}{\lambda + \alpha} \leq \frac{1}{2\sqrt{\alpha}}$ for any compact operator A (e.g., [19, p. 45, equation (2.48)]), the second term used [28, Lemma 4.2.8], and the third term used [33, Lemma 1].

Finally, by $K^*(KK^* + \alpha \text{id}) = (K^*K + \alpha \text{id})K^*$, $x_\alpha = (K^*K + \alpha \text{id})^{-1}K^*y_0 = K^*(KK^* + \alpha \text{id})^{-1}y_0$. Thus by the definition of γ_m ,

$$\|(\text{id} - P_m)x_\alpha\|_X = \|(\text{id} - P_m)K^*(KK^* + \alpha \text{id})^{-1}Kx_0\|_X \leq \gamma_m \|(KK^* + \alpha \text{id})^{-1}Kx_0\|_X.$$

For each $R > 0$, let $(v_n^R)_{n \in \mathbb{N}} \subset X$ such that $\lim_{n \rightarrow \infty} \|x_0 - K^*v_n^R\|_X = \mathcal{D}(R)$ and $\|v_n^R\|_X \leq R$ for all $n \in \mathbb{N}$. Then for all $n \in \mathbb{N}$, by spectral inequalities (see [19, p. 45]),

$$\begin{aligned} \|(KK^* + \alpha \text{id})^{-1}Kx_0\|_X &= \|(KK^* + \alpha \text{id})^{-1}K(K^*v_n^R + x_0 - K^*v_n^R)\|_X \\ &\leq \|v_n^R\|_X + \frac{1}{2\sqrt{\alpha}} \|x_0 - K^*v_n^R\|_X \leq R + \frac{1}{2\sqrt{\alpha}} \|x_0 - K^*v_n^R\|_X, \end{aligned}$$

from which by passing $n \rightarrow \infty$ yields $\|(KK^* + \alpha \text{id})^{-1}Kx_0\|_X \leq R + \frac{1}{2\sqrt{\alpha}} \mathcal{D}(R)$. Combining the above estimate with (6.3) leads to the desired result. \square

We now apply Lemma 6.2 in order to analyse (2.17) under a general regularity assumption of the kernel G^* . To this end, let $(V_N)_{N \in \mathbb{N}} \subset L^2([0, T], \mathbb{R})$ be a given family of finite-dimensional subspaces such that $\overline{\bigcup_{N=1}^\infty V_N} = L^2([0, T], \mathbb{R})$, and consider a slight generalisation of (2.17), where G^N , $N \in \mathbb{N}$, is defined as the minimiser of (2.16) over V_N . Then, by the first-order condition to (2.16), for all $N \in \mathbb{N}$,

$$G^N = ((\mathbf{u}_N^e)^* \mathbf{u}_N^e + \tau_N \text{id})^{-1} (\mathbf{u}_N^e)^* \left(-\frac{1}{N} \sum_{m=1}^N (S^m - A^m + \lambda^N u^e) \right), \quad (6.4)$$

with $\mathbf{u}_N^e := \mathbf{u}^e P_N$, where P_N is the projection of $L^2([0, T], \mathbb{R})$ onto V_N . For each $N \in \mathbb{N}$, let

$$\gamma_N = \|(\text{id} - P_N)(\mathbf{u}^e)^*\|_{\text{op}}. \quad (6.5)$$

Note that if $(V_N)_{N \in \mathbb{N}}$ is the space of piecewise constant functions as in (2.15), then there exists $C > 0$, depending on \mathbf{u}^e , such that $\gamma_N \leq C|\pi_N|$ for all $N \in \mathbb{N}$. Indeed, for all $N \in \mathbb{N}$ and $f \in L^2([0, T], \mathbb{R})$,

$$\|(\text{id} - P_N)(\mathbf{u}^e)^* f\|_{L^2([0, T])} \leq |\pi_N| \left\| \frac{d}{dt}((\mathbf{u}^e)^* f) \right\|_{L^2([0, T])} \leq C|\pi_N| \|f\|_{L^2([0, T])}, \quad (6.6)$$

where the first and second inequalities used [47, Theorem 6.1] and (6.1), respectively. We further introduce the following function $\mathcal{D} : (0, \infty) \rightarrow [0, \infty)$ to measure the regularity of G^* : for all $R > 0$,

$$\mathcal{D}(R) = \inf\{\|G^* - (\mathbf{u}^e)^* v\|_{L^2([0, T])} \mid v \in L^2([0, T], \mathbb{R}), \|v\|_{L^2([0, T])} \leq R\}. \quad (6.7)$$

Note that $G^* \in \mathcal{R}((\mathbf{u}^e)^*)$ if and only if $\mathcal{D}(R) = 0$ for all large $R > 0$. As most commonly used kernels are not in $\mathcal{R}((\mathbf{u}^e)^*)$ (see Remark 2.7 and Lemma 6.1), the subsequent analysis focuses on the case with $\mathcal{D}(R) > 0$ for all $R > 0$.

The following theorem presents a general version of Theorem 2.10, and quantifies the accuracy of $(\theta^N)_{N \in \mathbb{N}}$ under a power-type decay rate of the function \mathcal{D} .

Theorem 6.3. *Suppose Assumptions 2.2 and 2.8 hold. Assume further that there exists $\beta \in (0, 1]$ such that $\limsup_{R \rightarrow \infty} \mathcal{D}(R)R^\beta < \infty$. Let $C \geq 1$. Then for all $\eta \in (0, 1)$, by setting $(\tau_N)_{N \in \mathbb{N}} \subset (0, \infty)$ and $(V_N)_{N \in \mathbb{N}} \subset L^2([0, T], \mathbb{R})$ such that*

$$\frac{1}{C} \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{2(\beta+1)}{2\beta+1}} \leq \tau_N \leq C \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{2(\beta+1)}{2\beta+1}}, \quad \gamma_N \leq C\tau_N^{1/2},$$

it holds with probability at least $1 - \eta$ that, for all $N \in \mathbb{N} \cap [2, \infty)$,

$$|\lambda^N - \lambda^*| \leq C' \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right), \quad \|G^N - G^*\|_{L^2([0, T])} \leq C' \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^{\frac{\beta}{2\beta+1}}.$$

for some constant $C' > 0$ independent of η and N .

Proof. Throughout this proof, let $\eta \in (0, 1)$ be fixed, and let C' be a generic constant, which is independent of N and η , and may take a different value at each occurrence.

For each $N \in \mathbb{N}$, by (2.12) and (2.14),

$$\lambda^N - \lambda^* = -\frac{1}{Nu^e(0)} \sum_{m=1}^N (M_0^m - \lambda^* u^e(0)) - \lambda^* = -\frac{1}{Nu^e(0)} \sum_{m=1}^N M_0^m. \quad (6.8)$$

Let $y_0 = -(\mathbb{E}[S^1 - A^1] + \lambda^* u^e)$ and for each $N \in \mathbb{N}$, let $y_N = -\frac{1}{N} \sum_{m=1}^N (S^m - A^m + \lambda^N u^e)$, and let $\delta_N = \|y_N - y_0\|_{L^2([0, T])}$. Then by (2.12) and $\mathbb{E}[M_t^m] = 0$ for all $t \in [0, T]$

and $m \in \mathbb{N}$,

$$\begin{aligned} y_0 - y_N &= \frac{1}{N} \sum_{m=1}^N (S^m - A^m) - \mathbb{E}[S^1 - A^1] + (\lambda^N - \lambda^*)u^e \\ &= \frac{1}{N} \sum_{m=1}^N M^m + (\lambda^N - \lambda^*)u^e. \end{aligned} \quad (6.9)$$

Observe that $\mathbf{u}^e G^* = y_0$ and the choice of $(\tau_N, V_N)_{N \in \mathbb{N}}$ ensures that $\gamma_N \leq C' \sqrt{\tau_N}$ for all $N \in \mathbb{N}$. Moreover, $\limsup_{R \rightarrow \infty} \mathcal{D}(R)R^\beta < \infty$ implies that there exists $C > 0$ such that $\mathcal{D}(R) \leq CR^{-\beta}$ for all sufficiently large $R > 0$. Then by Lemma 6.2 (cf. (6.2) and (6.4)) and the decay condition of $\mathcal{D}(R)$, for all $N \in \mathbb{N}$ and $R > 0$,

$$\|G^N - G^*\|_{L^2([0,T])} \leq \frac{\delta_N}{2\sqrt{\tau_N}} + C' \sqrt{\tau_N} \left(R + \frac{1}{R^\beta \sqrt{\tau_N}} \right) + \frac{1}{R^\beta} + \frac{R\sqrt{\tau_N}}{2}.$$

from which by setting $R = \tau_N^{-1/2(\beta+1)}$, it holds for all $N \in \mathbb{N}$,

$$\|G^N - G^*\|_{L^2([0,T])} \leq \frac{\delta_N}{2\sqrt{\tau_N}} + C' \tau_N^{\frac{\beta}{2(\beta+1)}}. \quad (6.10)$$

We now estimate $(\delta_N)_{N \in \mathbb{N}}$ in high probability. For each $N \in \mathbb{N}$, by applying Assumption 2.8 with $\tilde{\eta} = \eta/N^2$, with probability at least $1 - \eta/N^2$,

$$\left| \frac{1}{N} \sum_{m=1}^N M_0^m \right| + \left\| \frac{1}{N} \sum_{m=1}^N M^m \right\|_{L^2([0,T])} \leq C' (\log(2\eta^{-1}) + \log N) N^{-\frac{1}{2}},$$

which along with $\sum_{N=2}^{\infty} \frac{1}{N^2} < 1$ implies that with probability at least $1 - \eta$,

$$\left| \frac{1}{N} \sum_{m=1}^N M_0^m \right| + \left\| \frac{1}{N} \sum_{m=1}^N M^m \right\|_{L^2([0,T])} \leq C' \frac{\log(\eta^{-1}) + \log N}{\sqrt{N}}, \quad \text{for all } N \in \mathbb{N} \cap [2, \infty). \quad (6.11)$$

Consider the event where (6.11) holds. Then, by (6.8) and (6.9),

$$|\lambda^N - \lambda^*| + \|y_N - y_0\|_{L^2([0,T])} \leq C' \frac{\log(\eta^{-1}) + \log N}{\sqrt{N}}, \quad \text{for all } N \in \mathbb{N} \cap [2, \infty). \quad (6.12)$$

Substituting the bound of $\|y_N - y_0\|_{L^2([0,T])}$ and the choice of τ_N into (6.10) yields the estimate of $\|G^N - G^*\|_{L^2([0,T])}$. \square

Based on Theorem 6.3, it suffices to establish the precise decay rate of $\mathcal{D} : (0, \infty) \rightarrow [0, \infty)$, in order to conclude Theorem 2.10 (recall the bounds of $(\gamma_N)_{N \in \mathbb{N}}$ following from (6.5) and (6.6)).

Theorem 6.4. *Suppose Assumption 2.2 holds.*

- (1) *If Assumption 2.6(1) holds, then $\limsup_{R \rightarrow \infty} \mathcal{D}(R)R < \infty$. Consequently, Theorem 6.3 holds with $\beta = 1$.*
- (2) *If Assumption 2.6(2) holds, then $\limsup_{R \rightarrow \infty} \mathcal{D}(R)R^{\frac{1-2\alpha}{1+2\alpha}} < \infty$. Consequently, Theorem 6.3 holds with $\beta = \frac{1-2\alpha}{1+2\alpha}$.*

Proof. Throughout this proof $C, C', \tilde{C} > 0$ are generic constants that may take a different value at each occurrence, and are independent of R . The proof relies on constructing specific sequence $(v^R)_{R>0}$ in $L^2([0, T], \mathbb{R})$ such that $\|v^R\|_{L^2([0, T])} \leq CR$ for all large R , and quantifying the decay rate of $\|G^* - (\mathbf{u}^e)^*v^R\|_{L^2([0, T])}$ as $R \rightarrow \infty$.

To prove Item (1), let $R_0 = T^{-1/2} > 0$. Then for each $R > R_0$, define $G_R \in H^1([0, T], \mathbb{R})$ such that $G_R(t) = G^*(t)$ for all $t \in [0, T - R^{-2}]$ and $G_R(t) = (T - t)R^2G^*(T - R^{-2})$ for all $t \in [T - R^{-2}, T]$. Lemma 6.1 implies that $(G_R)_{R>R_0} \subset \mathcal{R}((\mathbf{u}^e)^*)$. For each $R > R_0$, let $v^R \in L^2([0, T], \mathbb{R})$ be such that $G_R = (\mathbf{u}^e)^*v^R$, which along with (6.1) implies that $\dot{G}_R = -(u^e(0)\text{id} + \dot{\mathbf{u}}^e)^*v^R$. As $(u^e(0)\text{id} + \dot{\mathbf{u}}^e)^*$ has a bounded inverse on $L^2([0, T], \mathbb{R})$ (see [27, Corollary 9.3.16, p 238]), $\|v^R\|_{L^2([0, T])} \leq C\|\dot{G}_R\|_{L^2([0, T])}$ for all $R > R_0$. Observe that $\dot{G}_R(t) = \dot{G}^*(t)$ for all $t \in [0, T - R^{-2}]$ and $\dot{G}_R(t) = -R^2G^*(T - R^{-2})$ for all $t \in [T - R^{-2}, T]$. This along with $G^* \in H^1([0, T], \mathbb{R})$ implies that $\|\dot{G}_R\|_{L^2([0, T])} \leq CR$ for all $R > R_0$, and hence $\|v^R\|_{L^2([0, T])} \leq CR$ for all $R > R_0$. To estimate $\|G^* - (\mathbf{u}^e)^*v^R\|_{L^2([0, T])}$, note that by the definition of G_R and the fact that $G^* \in H^1([0, T], \mathbb{R})$ with $G^*(T) \neq 0$,

$$\begin{aligned} \|G^* - G_R\|_{L^2([0, T])}^2 &= \int_{T-R^{-2}}^T |G^*(t) - (T - t)R^2G^*(T - R^{-2})|^2 dt \\ &\leq \tilde{C} \left(R^{-2} + R^4 \int_{T-R^{-2}}^T (T - t)^2 dt \right) \leq \tilde{C}R^{-2}. \end{aligned}$$

Recalling (6.7), this shows that there exists $C, \tilde{C} > 0$ such that

$$\mathcal{D}(CR) \leq \|G^* - G_R\|_{L^2([0, T])} \leq \tilde{C}R^{-1}, \quad \text{for all } R > R_0.$$

Rescaling the inequality yields Item (1).

Next, we prove Item (2). As $|\dot{G}^*(t)| \leq C_0t^{-\alpha-1}$ for all $t < t_0$, we have for all $t < t_0$,

$$|G^*(t)| = \left| G^*(t_0) - \int_t^{t_0} \dot{G}^*(s) ds \right| \leq |G^*(t_0)| + \frac{C}{\alpha} (t^{-\alpha} - t_0^{-\alpha}) \leq Ct^{-\alpha}. \quad (6.13)$$

Let $R_0 > 1$ be such that $R_0^{-2} < R_0^{-2/(2\alpha+1)} < \min(t_0, T - t_0)$. For each $R > R_0$, define $G_R \in H^1([0, T], \mathbb{R})$ such that

$$G_R(t) = \begin{cases} G^*(R^{-2/(2\alpha+1)}), & t \in \mathcal{I}_1 := (0, R^{-2/(2\alpha+1)}), \\ G^*(t), & t \in \mathcal{I}_2 := [R^{-2/(2\alpha+1)}, T - R^{-2}], \\ (T - t)R^2G^*(T - R^{-2}), & t \in \mathcal{I}_3 := [T - R^{-2}, T]. \end{cases}$$

Lemma 6.1 implies that $(G_R)_{R>R_0} \subset \mathcal{R}((\mathbf{u}^e)^*)$. For each $R > 0$, let $v^R \in L^2([0, T], \mathbb{R})$ be such that $G_R = (\mathbf{u}^e)^* v^R$. Then by similar arguments as above, $\|v^R\|_{L^2([0, T])} \leq C \|\dot{G}_R\|_{L^2([0, T])}$ for all $R > R_0$, where for each $R > R_0$,

$$\dot{G}_R(t) = 0, \quad t \in \mathcal{I}_1; \quad \dot{G}_R(t) = \dot{G}^*(t), \quad t \in \mathcal{I}_2; \quad \dot{G}_R(t) = -R^2 G^*(T - R^{-2}), \quad t \in \mathcal{I}_3.$$

Then by $G^* \in H^1([t_0, T], \mathbb{R})$, $|G^*(t)| \leq C$ for all $t \in [T - R_0^{-2}, T]$. Hence by (6.13), for all $R > R_0$,

$$\begin{aligned} \|\dot{G}_R\|_{L^2([0, T])}^2 &\leq \int_{R^{-\frac{2}{2\alpha+1}}}^{t_0} (\dot{G}^*(t))^2 dt + \int_{t_0}^{T-R^{-2}} (\dot{G}^*(t))^2 dt + R^4 \int_{T-R^{-2}}^T (G^*(T - R^{-2}))^2 dt \\ &\leq C \left(R^2 - t_0^{-2\alpha-1} + \|\dot{G}^*\|_{L^2([t_0, T])}^2 + R^2 \right) \leq CR^2. \end{aligned}$$

This implies that $\|v^R\|_{L^2([0, T])} \leq CR$ for all $R > R_0$. To estimate $\|G^* - (\mathbf{u}^e)^* v^R\|_{L^2([0, T])}$, for all $R > R_0$,

$$\begin{aligned} \|G^* - G_R\|_{L^2([0, T])}^2 &\leq C \left(\|G^*\|_{L^2(\mathcal{I}_1)}^2 + \|G^*\|_{L^2(\mathcal{I}_3)}^2 + \|G_R\|_{L^2(\mathcal{I}_1)}^2 + \|G_R\|_{L^2(\mathcal{I}_3)}^2 \right) \\ &\leq C \left(\int_0^{R^{-\frac{2}{2\alpha+1}}} (G^*(t))^2 dt + R^{-2} + \int_0^{R^{-\frac{2}{2\alpha+1}}} (G^*(R^{-\frac{2}{2\alpha+1}}))^2 dt + R^4 \int_{T-R^{-2}}^T (T-t)^2 dt \right) \\ &\leq CR^{-\frac{2(1-2\alpha)}{2\alpha+1}}, \end{aligned}$$

where the last inequality used (6.13) and $R^{-\frac{2(1-2\alpha)}{2\alpha+1}} \geq R^{-2}$ due to $\alpha \in (0, 1/2)$. This shows that there exists $C, \tilde{C} > 0$ such that $\mathcal{D}(CR) \leq \|G^* - G_R\|_{L^2([0, T])} \leq \tilde{C}R^{-\frac{1-2\alpha}{2\alpha+1}}$ for all sufficiently large $R > 0$. Rescaling the inequality yields Item (2). \square

Now we are ready to prove Theorem 2.10.

Proof of Theorem 2.10. By combining the result of Theorem 6.3 with the decay rate of $\mathcal{D} : (0, \infty) \rightarrow [0, \infty)$ in Theorem 6.4 and the bounds on $(\gamma_N)_{N \in \mathbb{N}}$ following from (6.5) and (6.6) we get the result. \square

The following proposition shows that the decay rates of \mathcal{D} in Theorem 6.4 are optimal, i.e., they are the maximal power-type decay rates under Assumption 2.6. Specifically, it considers the power law kernel $G^*(t) = t^{-\alpha}$, which satisfies Assumption 2.6(1) if $\alpha = 0$, and Assumption 2.6(2) if $\alpha \in (0, 1/2)$.

Proposition 6.5. *Let $\alpha \in [0, 1/2)$, let $G^* \in L^2([0, T], \mathbb{R})$ be such that $G^*(t) = t^{-\alpha}$ for $t \in (0, T]$, and let $u^e \in H^1([0, T], \mathbb{R})$ be such that $u^e(t) = 1$ for all $t \in [0, T]$. Then $\limsup_{R \rightarrow \infty} \mathcal{D}(R)R^{\frac{1-2\alpha}{1+2\alpha}} < \infty$ and for all $\varepsilon > \frac{1-2\alpha}{1+2\alpha}$, $\limsup_{R \rightarrow \infty} \mathcal{D}(R)R^\varepsilon = \infty$.*

Proof. In the sequel, we focus on the case with $\alpha \in (0, 1/2)$, as the result for $\alpha = 0$ has been proved in [33]. By [33, Remark 1], the decay rate of \mathcal{D} can be characterised in terms of the singular system of \mathbf{u}^e . More precisely, let $\sigma_1 \geq \sigma_2 \geq \dots > 0$ be the ordered singular values of \mathbf{u}^e with $\lim_{n \rightarrow \infty} \sigma_n = 0$, and $(\mathbf{u}_n)_{n \in \mathbb{N}}$ be the orthonormal eigensystems of $(\mathbf{u}^e)^* \mathbf{u}^e$. Then for any $G^* \in L^2([0, T], \mathbb{R})$ with

$$\kappa_0 := \sup \left\{ \kappa > 0 \mid \sum_{n=1}^{\infty} \frac{1}{\sigma_n^{2\kappa}} \langle G^*, \mathbf{u}_n \rangle_{L^2([0, T])}^2 < \infty \right\} \in (0, 1), \quad (6.14)$$

and for any $\kappa \in (0, 1)$, $\limsup_{R \rightarrow \infty} \mathcal{D}(R) R^{\frac{\kappa}{1-\kappa}} < \infty$ if and only if $0 < \kappa \leq \kappa_0$. In the sequel, we compute the value κ_0 for $G^*(t) = t^{-\alpha}$ and $u^e(t) = 1$, $t \in [0, T]$.

The fact that $u^e \equiv 1$ implies that $(\mathbf{u}^e f)(t) = \int_0^t f(s) ds$ for all $f \in L^2([0, T], \mathbb{R})$ and $t \in [0, T]$. Let $\alpha \in (0, 1/2)$ and $G^*(t) = t^{-\alpha}$ for all $t \in (0, T]$. A direct computation shows that for all $n \in \mathbb{N}$, $\sigma_n = \frac{T}{\pi(n-\frac{1}{2})}$, $\mathbf{u}_n(t) = \sqrt{\frac{2}{T}} \cos((n - \frac{1}{2})\pi \frac{t}{T})$ for all $t \in [0, T]$, and

$$\langle G^*, \mathbf{u}_n \rangle_{L^2([0, T])} = \sqrt{\frac{2}{T}} \int_0^T t^{-\alpha} \cos\left(\frac{n-\frac{1}{2}}{T} \pi t\right) dt = \sqrt{\frac{2}{T}} \left(\frac{n-\frac{1}{2}}{T}\right)^{\alpha-1} \int_0^{n-\frac{1}{2}} t^{-\alpha} \cos(\pi t) dt. \quad (6.15)$$

To estimate $\langle G^*, \mathbf{u}_n \rangle_{L^2([0, T])}$ for large n , we prove that $\lim_{R \rightarrow \infty} \int_0^R t^{-\alpha} \cos(\pi t) dt \in (0, \infty)$. Indeed, for each $0 < R_1 < R_2 < \infty$, consider the closed contour $\gamma = (\gamma_1, \gamma_2, \gamma_3, \gamma_4) \subset \mathbb{C}$, where $\gamma_1(\theta) = R_1 e^{i(\frac{\pi}{2}-\theta)}$ for $\theta \in [0, \frac{\pi}{2}]$, $\gamma_2(r) = r$ for $r \in [R_1, R_2]$, $\gamma_3(\theta) = R_2 e^{i\theta}$ for $\theta \in [0, \frac{\pi}{2}]$, and $\gamma_4(r) = i(R_1 + R_2 - r)$ for $r \in [R_1, R_2]$. As $\alpha > 0$, the function $\mathbb{C} \setminus \{0\} \ni z \mapsto z^{-\alpha} e^{i\pi z} \in \mathbb{C}$ is analytic. Hence by the Cauchy-Goursat theorem and the definition of contour integral, for each $0 < R_1 < R_2 < \infty$,

$$0 = \int_{\gamma} z^{-\alpha} e^{i\pi z} dz = I_{R_1} + \int_{R_1}^{R_2} r^{-\alpha} e^{i\pi r} dr + I_{R_2} + \int_{R_2}^{R_1} (ir)^{-\alpha} e^{i\pi(ir)} i dr, \quad (6.16)$$

where I_{R_1} and I_{R_2} are the integrals along the curves γ_1 and γ_3 , respectively:

$$I_{R_1} = \int_0^{\frac{\pi}{2}} \gamma_1(\theta)^{-\alpha} e^{i\pi \gamma_1(\theta)} R_1 e^{i(\frac{\pi}{2}-\theta)} (-i) d\theta, \quad I_{R_2} = \int_0^{\frac{\pi}{2}} \gamma_3(\theta)^{-\alpha} e^{i\pi \gamma_3(\theta)} R_2 e^{i\theta} i d\theta.$$

By the definition of γ_1 , for all $\theta \in (0, \pi/2)$ and $R_1 > 0$, $|e^{i\pi \gamma_1(\theta)}| = e^{-\pi R_1 \sin(\frac{\pi}{2}-\theta)} \leq 1$. As $|\int_0^{\frac{\pi}{2}} f(\theta) d\theta| \leq \int_0^{\frac{\pi}{2}} |f(\theta)| d\theta$ for any integrable $f : [0, \frac{\pi}{2}] \rightarrow \mathbb{C}$,

$$\lim_{R_1 \rightarrow 0^+} |I_{R_1}| \leq \lim_{R_1 \rightarrow 0^+} \int_0^{\frac{\pi}{2}} R_1^{1-\alpha} d\theta \leq \lim_{R_1 \rightarrow 0^+} \frac{\pi}{2} R_1^{1-\alpha} = 0,$$

where the last identity used $\alpha \in (0, 1/2)$. Similarly, by the definition of γ_3 , for all $\theta \in (0, \pi/2)$ and $R_2 > 0$, $|e^{i\pi \gamma_3(\theta)}| = e^{-\pi R_2 \sin(\theta)}$ and $|I_{R_2}| \leq R_2^{1-\alpha} \int_0^{\frac{\pi}{2}} e^{-\pi R_2 \sin(\theta)} d\theta$. As

$\sin \theta \geq \frac{2}{\pi} \theta$ for all $\theta \in (0, \pi/2)$,

$$\lim_{R_2 \rightarrow \infty} |I_{R_2}| \leq \lim_{R_2 \rightarrow \infty} \left(R_2^{1-\alpha} \int_0^{\frac{\pi}{2}} e^{-2R_2\theta} d\theta \right) = \lim_{R_2 \rightarrow \infty} \left(\frac{1}{2} R_2^{-\alpha} (1 - e^{-R_2\pi}) \right) = 0.$$

Thus, letting $R_1 \rightarrow 0^+$ and $R_2 \rightarrow \infty$ in (6.16) yields

$$\lim_{R_2 \rightarrow \infty} \int_0^{R_2} r^{-\alpha} e^{i\pi r} dr = i^{1-\alpha} \lim_{R_2 \rightarrow \infty} \int_0^{R_2} r^{-\alpha} e^{-\pi r} dr.$$

The upper bound $\lim_{R \rightarrow \infty} \int_0^R t^{-\alpha} \cos(\pi t) dt < \infty$ follows from $\lim_{R_2 \rightarrow \infty} \int_0^{R_2} r^{-\alpha} e^{-\pi r} dr < \infty$ and $\int_0^R t^{-\alpha} \cos(\pi t) dt$ is the real part of $\int_0^R t^{-\alpha} e^{i\pi t} dt$ for all $R > 0$. For the lower bound, consider the sequence $(a_n)_{n \in \mathbb{N}}$ with $a_n := \int_0^{2n} t^{-\alpha} \cos(\pi t) dt$ for all $n \in \mathbb{N}$. As $t \mapsto t^{-\alpha}$ is decreasing on $(0, \infty)$, $a_{n+1} \geq a_n > 0$ for all $n \in \mathbb{N}$, and hence $\lim_{R \rightarrow \infty} \int_0^R t^{-\alpha} \cos(\pi t) dt = \lim_{n \rightarrow \infty} \int_0^{2n} t^{-\alpha} \cos(\pi t) dt \geq a_1 > 0$.

Therefore, by (6.15), for each $\alpha \in (0, 1/2)$, there exists $C > 0$, depending on α and T , such that $\frac{1}{C} n^{2\alpha-2} \leq \langle G^*, \mathbf{u}_n \rangle_{L^2([0,T])}^2 \leq C n^{2\alpha-2}$ for all $n \in \mathbb{N}$, and for all $\kappa > 0$,

$$\frac{1}{C} \sum_{n=1}^{\infty} \frac{1}{n^{2-2\kappa-2\alpha}} \leq \sum_{n=1}^{\infty} \frac{1}{\sigma_n^{2\kappa}} \langle G^*, \mathbf{u}_n \rangle_{L^2([0,T])}^2 \leq C \sum_{n=1}^{\infty} \frac{1}{n^{2-2\kappa-2\alpha}}. \quad (6.17)$$

Hence $\sum_{n=1}^{\infty} \frac{1}{\sigma_n^{2\kappa}} \langle G^*, \mathbf{u}_n \rangle_{L^2([0,T])}^2 < \infty$ if and only if $\kappa \in (0, \frac{1}{2} - \alpha)$. This implies that $\kappa_0 = \frac{1}{2} - \alpha \in (0, 1)$ in (6.14), and subsequently proves that for all $\kappa \in (\frac{1}{2} - \alpha, 1)$, $\limsup_{R \rightarrow \infty} \mathcal{D}(R) R^{\frac{\kappa}{1-\kappa}} = \infty$. The desired statement follows from the fact that $\{\frac{\kappa}{1-\kappa} \mid \kappa \in (\frac{1}{2} - \alpha, 1)\} = (\frac{1-2\alpha}{1+2\alpha}, \infty)$. \square

7 Proof of Theorem 2.19

Throughout this section, we denote by $C > 0$ a generic constant, which is independent of $\eta \in (0, 1)$ and $N \in \mathbb{N}$, and may take a different value at each occurrence.

The following lemma proves that for sufficiently large $\mathbf{m}_0^e \in \mathbb{N}$, the estimators $(\theta^k)_{k \geq 0}$ from Algorithm 1 are in the admissible parameter set Ξ_ε from Definition 2.14.

Lemma 7.1. *Assume that condition (2.26) holds. Let $(\theta^k)_{k \geq 0}$ be generated from Algorithm 1. Then there exists $C_0 > 0$ such that for all $\eta \in (0, 1)$, and $\mathbf{m}_0^e \geq C_0(\log(\eta^{-1})^2 + 1)$, the following hold with probability at least $1 - \eta$.*

- (1) $\theta^k = (\lambda^k, G^k) \in \Xi_\varepsilon$ for all $k \in \mathbb{N} \cup \{0\}$,
- (2) There exists $C > 0$, independent of η , such that for all $k \in \mathbb{N} \cup \{0\}$,

$$|\lambda^k - \lambda^*| + \|G^k - G^*\|_{L^2([0,T])} \leq C \left(\frac{\log(\eta^{-1}) + \log(\mathbf{m}_0^e + k)}{\sqrt{\mathbf{m}_0^e + k}} \right)^\kappa.$$

Proof. Observe that for all $G, f \in L^2([0, T], \mathbb{R})$, by the Cauchy-Schwarz inequality,

$$\begin{aligned}
\int_0^T \int_0^T G(|t-s|) f(s) f(t) ds dt &= 2 \int_0^T \int_0^T G(t-s) \mathbb{1}_{\{s \leq t\}} f(s) f(t) ds dt \\
&\leq 2 \|f\|_{L^2([0, T])} \int_0^T \left(\int_0^T (G(t-s))^2 \mathbb{1}_{\{s \leq t\}} ds \right)^{1/2} f(t) dt \\
&\leq 2 \|f\|_{L^2([0, T])}^2 \left(\int_0^T \int_0^T (G(t-s))^2 \mathbb{1}_{\{s \leq t\}} ds dt \right)^{1/2} \\
&\leq 2\sqrt{T} \|f\|_{L^2([0, T])}^2 \|G\|_{L^2([0, T])}.
\end{aligned}$$

Consequently, by (2.3) and Definition 2.14 there exists $\varepsilon' > 0$ such that $\mathbb{B}(\theta^*, \varepsilon') := \{(\lambda, G) \in \mathbb{R} \times L^2([0, T], \mathbb{R}) \mid |\lambda - \lambda^*| + \|G - G^*\|_{L^2([0, T])} \leq \varepsilon'\} \subset \Xi_\varepsilon$.

Let $\tilde{C} > 0$ be the constant in the condition (2.26), and assume without loss of generality that $\varepsilon'/\tilde{C} < 1$. Observe that for all $x, y > 0$, $\frac{\partial}{\partial y} \left(\frac{x+\log y}{\sqrt{y}} \right) = \frac{2-\log y-x}{2y^{3/2}}$, and hence $y \mapsto \frac{x+\log y}{\sqrt{y}}$ decreases on $[e^2, \infty)$ for all $x > 0$. Let $C_0 \geq e^2$ be a constant such that $\frac{x+\log(C_0(x^2+1))}{\sqrt{C_0(x^2+1)}} \leq \left(\frac{\varepsilon'}{\tilde{C}}\right)^{1/\kappa}$ for all $x > 0$. Then for any $\eta \in (0, 1)$ and $N \geq C_0(\log(\eta^{-1})^2 + 1)$,

$$\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \leq \left(\frac{\log(\eta^{-1}) + \log N}{\sqrt{N}} \right)^\kappa \leq \frac{\varepsilon'}{\tilde{C}},$$

where the first inequality used $\kappa \in (0, 1)$. Now let $\mathbf{m}_0^e \geq C_0(\log(\eta^{-1})^2 + 1)$. By the definitions of $(\theta^k)_{k \geq 0}$ in (2.24) and (2.25), for each $k \in \mathbb{N} \cup \{0\}$, the estimator θ^k is computed using $k + \mathbf{m}_0^e$ samples, which along with (2.26) shows that $\theta^k \in \mathbb{B}(\theta^*, \varepsilon')$ for all $k \geq 0$ in Item (1). The convergence rate of $(\theta^k)_{k \geq 0}$ in Item (2) follows from (2.26) with $N = k + \mathbf{m}_0^e$. \square

We are now ready to prove Theorem 2.19.

Proof of Theorem 2.19. Throughout this proof, let $\mathcal{E} = \{1, \dots, \mathbf{m}_0^e\} \cup \{\mathbf{m}_0^e + \sum_{j=1}^k \mathbf{n}(j) + k \mid k \in \mathbb{N}\}$ be the indices of exploration episodes, let $\mathfrak{L}(k) = \mathbf{m}_0^e + \sum_{j=1}^k \mathbf{n}(j) + k$, $k \in \mathbb{N} \cup \{0\}$ be the index of the last episode in the k -th cycle (cf. Algorithm 1), and let $\mathbf{c}(m)$, $m \in \mathbb{N}$, be the corresponding cycle for the m -th episode, i.e., $\mathbf{c}(m) = \min\{k \in \mathbb{N} \cup \{0\} \mid \mathfrak{L}(k) \geq m\}$. Let $\mathbf{n}(k) = \lfloor k^\delta \rfloor$, $k \in \mathbb{N}$, with some $\delta \in (0, 1)$ to be determined later. Observe that for all $m > \mathbf{m}_0^e$, $\mathfrak{L}(\mathbf{c}(m) - 1) < m \leq \mathfrak{L}(\mathbf{c}(m))$, which along with the inequality that $k^\delta - 1 \leq \lfloor k^\delta \rfloor \leq k^\delta$ for all $k \in \mathbb{N}$ implies that

$$\sum_{j=1}^{\mathbf{c}(m)-1} (j^\delta - 1) + \mathbf{c}(m) - 1 < m - \mathbf{m}_0^e \leq \sum_{j=1}^{\mathbf{c}(m)} j^\delta + \mathbf{c}(m).$$

Hence there exists $\bar{c}, \underline{c} > 0$, depending on δ , such that for all $m > \mathbf{m}_0^e$,

$$\underline{c}(\mathbf{c}(m) - 1)^{\delta+1} \leq m - \mathbf{m}_0^e \leq \bar{c}(\mathbf{c}(m))^{\delta+1} \quad (7.1)$$

In the following, we optimise the growth rate of $R(N)$ over δ .

Let $\eta \in (0, 1)$ be fixed, and let $\mathbf{m}_0^e = \lceil C(\log(\eta^{-1})^2 + 1) \rceil$ with $C \geq C_0$ in Lemma 7.1, and consider the event (which holds with probability at least $1 - \eta$) such that $\theta^k \in \Xi_\varepsilon$ for all $k \in \mathbb{N} \cup \{0\}$. For each $N \in \mathbb{N}$, by (2.10) and Algorithm 1,

$$\begin{aligned} R(N) &= \sum_{m \in [1, N] \cap \mathcal{E}} (J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^m)) + \sum_{m \in [1, N] \setminus \mathcal{E}} (J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^m)) \\ &= \sum_{m \in [1, N] \cap \mathcal{E}} (J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^e)) + \sum_{m \in [1, N] \setminus \mathcal{E}} (J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^{\theta^{\mathbf{c}(m)-1}})). \end{aligned} \quad (7.2)$$

If $N \leq \mathbf{m}_0^e$, the fact that $J^{\theta^*}(u^{\theta^*}), J^{\theta^*}(u^e) < \infty$ implies that $R(N) \leq \mathbf{m}_0^e |J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^e)| \leq C(\log(\eta^{-1})^2 + 1)$. Now consider $N > \mathbf{m}_0^e$. The number of exploration episodes up to the N -th episode is bounded by $\mathbf{m}_0^e + \mathbf{c}(N)$. Moreover, by Lemma 7.1 and Theorem 2.16, $|J^{\theta^*}(u^{\theta^*}) - J^{\theta^*}(u^{\theta^{\mathbf{c}(m)-1}})| \leq C(|\lambda^{\mathbf{c}(m)-1} - \lambda^*|^2 + \|G^{\mathbf{c}(m)-1} - G^*\|_{L^2([0, T])}^2)$ for some constant $C > 0$. Hence, by Lemma 7.1 and (7.2),

$$\begin{aligned} R(N) &\leq C(\mathbf{m}_0^e + \mathbf{c}(N)) + C \sum_{k=1}^{\mathbf{c}(N)} \mathbf{n}(k) \left(|\lambda^k - \lambda^*|^2 + \|G^k - G^*\|_{L^2([0, T])}^2 \right) \\ &\leq C(\log(\eta^{-1})^2 + 1 + \mathbf{c}(N)) + C \sum_{k=1}^{\mathbf{c}(N)} \mathbf{n}(k) \left(\frac{\log(\eta^{-1}) + \log(\mathbf{m}_0^e + k)}{\sqrt{\mathbf{m}_0^e + k}} \right)^{2\kappa} \\ &\leq C(\log(\eta^{-1})^2 + 1 + \mathbf{c}(N)) + C \sum_{k=1}^{\mathbf{c}(N)} k^{\delta-\kappa} (\log(\eta^{-1}) + \log N)^{2\kappa} \\ &\leq C(\log(\eta^{-1})^2 + 1 + \mathbf{c}(N)) + C \frac{\mathbf{c}(N)^{\delta-\kappa+1}}{\delta - \kappa + 1} (\log(\eta^{-1}) + \log N)^{2\kappa}. \end{aligned}$$

Then from (7.1) it follows that $\mathbf{c}(N) \leq C(N - \mathbf{m}_0^e)^{1/(1+\delta)}$, which implies that for all $N > \mathbf{m}_0^e$,

$$R(N) \leq C \left(\log(\eta^{-1})^2 + N^{\frac{1}{1+\delta}} + \frac{1}{\delta - \kappa + 1} N^{\frac{\delta-\kappa+1}{1+\delta}} (\log(\eta^{-1}) + \log N)^{2\kappa} \right).$$

Hence it is clear that the growth rate of $(R(N))_{N \in \mathbb{N}}$ is optimised at $\delta = \kappa$. This proves the desired estimate. \square

8 Proof of Proposition 5.3

The proof of Proposition 5.3 will follow by proving the stability of the coefficients a^θ and B^θ in (2.8) with respect to θ . Note that by (4.13) the stability of these coefficients

will depend on the stability of the operator Γ_t^{-1} in (4.9) and of $\{\Theta_t(s)\}_{\{t \in [0, s], s \in [0, T]\}}$ in (4.11) with respect to θ . We emphasise the dependence of the ingredients of a^θ and B^θ in θ in the following. Throughout this section we assume that $\phi > 0$ in (2.7), where the proof of Proposition 5.3 for the case where $\phi = 0$ follows the same lines but is much simpler as a^θ and B^θ considerably simplify.

We recall that for any operator \mathbf{G} from $L^2([0, T], \mathbb{R}^2)$ to itself we define the operator norm,

$$\|\mathbf{G}\|_{\text{op}} = \sup_{f \in L^2([0, T], \mathbb{R}^2), f \neq 0} \frac{\|\mathbf{G}f\|_{L^2([0, T])}}{\|f\|_{L^2([0, T])}}. \quad (8.1)$$

The following Lemma and Propositions are essential ingredients for the proof of Proposition 5.3. Recall that Ξ_ε was defined in Definition 2.14.

Lemma 8.1. *Let $(\Gamma_t^\theta)^{-1}$ be defined as in (4.9), then we have,*

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \|(\Gamma_t^\theta)^{-1}\|_{\text{op}} < \infty.$$

The proof of Lemma 8.1 is postponed to Section 9.

Proposition 8.2 (Stability of $(\Gamma_t^\theta)^{-1}$). *There exists a constant $C \geq 0$ such that*

$$\sup_{t \in [0, T]} \|(\Gamma_t^\theta)^{-1} - (\Gamma_t^{\theta'})^{-1}\|_{\text{op}} \leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}), \quad \text{for all } \theta, \theta' \in \Xi_\varepsilon.$$

The proof of Proposition 8.2 is postponed to Section 9.

For any $\theta \in \Xi_\varepsilon$ we define $\Theta^\theta = \{\Theta_t^\theta(s) : t \in [0, s], s \in [0, T]\}$ as in (4.11), that is,

$$\Theta_t^\theta(s) = - \left((\Gamma_t^\theta)^{-1} \mathbf{1}_t \mathbb{E} \left[A_{(\cdot)} - A_T \mid \mathcal{F}_t \right] e_1 \right) (s). \quad (8.2)$$

Proposition 8.3 (Stability of Θ). *There exists a constant $C > 0$ such that for all $t \in [0, T]$ and $\theta, \theta' \in \Xi_\varepsilon$ we have,*

$$\|\Theta_t^\theta(\cdot) - \Theta_t^{\theta'}(\cdot)\|_{L^2([0, T])} \leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}) \|\mathbb{E}[A_{(\cdot)} - A_T \mid \mathcal{F}_t]\|_{L^2([0, T])}.$$

Proof. Throughout this proof let $(\omega, t) \in \Omega \times [0, T]$ and define

$$f_t(\cdot) = \mathbf{1}_t(\cdot) \mathbb{E}[A_{(\cdot)} - A_T \mid \mathcal{F}_t](\omega).$$

From (8.2) it follows that

$$\Theta_t^\theta(s) - \Theta_t^{\theta'}(s) = - \left(\left((\Gamma_t^\theta)^{-1} - (\Gamma_t^{\theta'})^{-1} \right) \mathbf{1}_t \mathbb{E} \left[A_{(\cdot)} - A_T \mid \mathcal{F}_t \right] e_1 \right) (s). \quad (8.3)$$

From (8.3) and Theorem 8.2 we get

$$\begin{aligned} \left\| \Theta_t^\theta(\cdot) - \Theta_t^{\theta'}(\cdot) \right\|_{L^2([0, T])} &\leq C \left\| (\Gamma_t^\theta)^{-1} - (\Gamma_t^{\theta'})^{-1} \right\|_{\text{op}} \|f_t(\cdot) e_1\|_{L^2([0, T])} \\ &\leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}) \|f_t(\cdot)\|_{L^2([0, T])}, \end{aligned}$$

where $C > 0$ is a constant depending only on T, L and ε , but independent of $t, s \in [0, T]$, $\theta, \theta' \in \Xi_\varepsilon$ and $\omega \in \Omega$. \square

Proof of Proposition 5.3. Inspired by (4.13) we define for any $\theta \in \Xi_\varepsilon$,

$$\begin{aligned} a_t^\theta &= \frac{1}{2\lambda} \left(\mathbb{E}[A_t - A_T \mid \mathcal{F}_t] + 2\varrho q + \langle \Theta_t^\theta, K_t^\theta \rangle_{L^2} + \langle (\mathbf{\Gamma}_t^\theta)^{-1} K_t^\theta, \mathbf{1}_t(-2\varrho q, q)^\top \rangle_{L^2} \right), \\ B^\theta(t, s) &= \mathbf{1}_{\{s < t\}} \frac{1}{2\lambda} \left(\langle (\mathbf{\Gamma}_t^\theta)^{-1} K_t^\theta, \mathbf{1}_t(\tilde{G}^\theta(\cdot, s), -1)^\top \rangle_{L^2} - \tilde{G}^\theta(t, s) \right). \end{aligned} \quad (8.4)$$

From Definition 2.14 we have $\lambda, \lambda' \geq L^{-1}$ hence

$$\begin{aligned} |a_t^\theta - a_t^{\theta'}| &\leq \left| \frac{1}{2\lambda} \left(\langle \Theta_t^\theta, K_t^\theta \rangle_{L^2} + \langle (\mathbf{\Gamma}_t^\theta)^{-1} K_t^\theta, \mathbf{1}_t(-2\varrho q, q)^\top \rangle_{L^2} \right) \right. \\ &\quad \left. - \frac{1}{2\lambda'} \left(\langle \Theta_t^{\theta'}, K_t^{\theta'} \rangle_{L^2} + \langle (\mathbf{\Gamma}_t^{\theta'})^{-1} K_t^{\theta'}, \mathbf{1}_t(-2\varrho q, q)^\top \rangle_{L^2} \right) \right| \\ &\leq C_1 |\lambda - \lambda'| \left| \langle \Theta_t^\theta, K_t^\theta \rangle_{L^2} + \langle (\mathbf{\Gamma}_t^\theta)^{-1} K_t^\theta, \mathbf{1}_t(-2\varrho q, q)^\top \rangle_{L^2} \right| \\ &\quad + C_2 |\langle \Theta_t^{\theta'}, K_t^{\theta'} \rangle_{L^2} - \langle \Theta_t^\theta, K_t^\theta \rangle_{L^2}| \\ &\quad + C_3 |\langle (\mathbf{\Gamma}_t^{\theta'})^{-1} K_t^{\theta'} - (\mathbf{\Gamma}_t^\theta)^{-1} K_t^\theta, \mathbf{1}_t(-2\varrho q, q)^\top \rangle_{L^2}| \\ &=: \sum_{j=1}^3 I_j^{\theta, \theta'}(t). \end{aligned} \quad (8.5)$$

From (4.5), (4.14) and Definition 2.14 we get

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \|K_t^\theta\|_{L^2([0, T])} < \infty. \quad (8.6)$$

From (2.2), Lemma 8.1 and (8.2) we get,

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \mathbb{E} \left[\int_0^T (\Theta_t^\theta(s))^2 ds \right] < \infty. \quad (8.7)$$

From (8.5), (8.6), (8.7), Lemma 8.1 and Cauchy-Schwarz inequality it follows that

$$\sup_{t \in [0, T]} \mathbb{E}[(I_1^{\theta, \theta'}(t))^2] \leq C |\lambda - \lambda'|^2. \quad (8.8)$$

Note that,

$$|I_2^{\theta, \theta'}(t)| \leq C_1 |\langle \Theta_t^{\theta'}, K_t^{\theta'} - K_t^\theta \rangle_{L^2}| + C_2 |\langle \Theta_t^\theta - \Theta_t^{\theta'}, K_t^\theta \rangle_{L^2([0, T])}|. \quad (8.9)$$

From (4.5), (4.14) and Definition 2.14 we get

$$\|K_t^{\theta'} - K_t^\theta\|_{L^2} \leq C \|G - G'\|_{L^2([0, T])}. \quad (8.10)$$

From Proposition 8.3, (8.6), (8.7), (8.9) and (8.10) it follows that

$$\sup_{t \in [0, T]} \mathbb{E}[(I_2^{\theta, \theta'}(t))^2] \leq C \|G - G'\|_{L^2([0, T])}^2. \quad (8.11)$$

Following similar steps as in (8.9)–(8.11), only using Proposition 8.2 instead of Proposition 8.3 we get,

$$\sup_{t \in [0, T]} I_3^{\theta, \theta'}(t) \leq C \|G - G'\|_{L^2([0, T])}. \quad (8.12)$$

Plugging in (8.8), (8.11) and (8.12) into (8.5) we get,

$$\sup_{t \in [0, T]} \mathbb{E}[(a_t^\theta - a_t^{\theta'})^2] \leq C \left(|\lambda - \lambda'|^2 + \|G - G'\|_{L^2([0, T])}^2 \right), \quad \text{for all } \theta, \theta' \in \Xi_\varepsilon. \quad (8.13)$$

By repeating a similar argument leading to (8.13), using (8.6), (8.10), Lemma 8.1 and Proposition 8.2 on (8.4) we obtain for all $\theta, \theta' \in \Xi_\varepsilon$,

$$\sup_{t \in [0, T]} \|B^\theta(t, \cdot) - B^{\theta'}(t, \cdot)\|_{L^2([0, T])} \leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}), \quad (8.14)$$

and

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \|B^\theta(t, \cdot)\|_{L^2([0, T])} < \infty. \quad (8.15)$$

The following bound can be obtained from (4.12), (8.13), (8.15) and Gronwall's lemma by using standard arguments,

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \mathbb{E}[(u_t^\theta)^2] < \infty. \quad (8.16)$$

Since in the following we use a similar argument to derive the stability of u^θ , we omit the details in order to avoid unnecessary repetition.

Recall that the \mathcal{H}^2 -norm was defined in (5.1). From (4.12), (8.13), (8.14) and Cauchy-Schwarz inequality we therefore get

$$\begin{aligned} \mathbb{E}[(u_t^\theta - u_t^{\theta'})^2] &\leq C \left(\mathbb{E}[(a_t^\theta - a_t^{\theta'})^2] + \mathbb{E} \left[\left(\int_0^t (B^\theta(t, s) - B^{\theta'}(t, s)) u_s^\theta ds \right)^2 \right] \right. \\ &\quad \left. + \mathbb{E} \left[\left(\int_0^t B^\theta(t, s) (u_s^\theta - u_s^{\theta'}) ds \right)^2 \right] \right) \\ &\leq C_1 \left(|\lambda - \lambda'|^2 + \|G - G'\|_{L^2([0, T])}^2 \right) (1 + \|u^\theta\|_{\mathcal{H}^2}) \\ &\quad + C_2 \sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \|B^\theta(t, \cdot)\|_{L^2([0, T])} \int_0^t \mathbb{E}[(u_s^\theta - u_s^{\theta'})^2] ds. \end{aligned} \quad (8.17)$$

Together with (8.15) and (8.16) it follows that there exist $C_i > 0$, $i = 1, 2$ not depending on $\theta, \theta' \in \Xi_\varepsilon$ such that,

$$\begin{aligned} \sup_{r \in [0, t]} \mathbb{E}[(u_r^\theta - u_r^{\theta'})^2] &\leq C_1 \left(|\lambda - \lambda'|^2 + \|G - G'\|_{L^2([0, T])}^2 \right) \\ &\quad + C_2 \int_0^t \sup_{r \in [0, s]} \mathbb{E}[(u_s^\theta - u_s^{\theta'})^2] ds. \end{aligned} \quad (8.18)$$

Then from Gronwall's lemma it follows that

$$\sup_{r \in [0, T]} \mathbb{E}[(u_r^\theta - u_r^{\theta'})^2] \leq C \left(|\lambda - \lambda'|^2 + \|G - G'\|_{L^2([0, T])}^2 \right),$$

and we get the result. \square

9 Proofs of Lemma 8.1 and Proposition 8.2

As we mentioned at the beginning of Section 8, we assume that $\phi > 0$ in (2.7), where the case of $\phi = 0$ is much simpler and is left to reader. From (4.9) we note that for $\phi > 0$, $(\mathbf{\Gamma}_t^\theta)^{-1}$ is the inverse of the operator

$$\mathbf{\Gamma}_t^\theta = \begin{pmatrix} \mathbf{D}_t^\theta - 2\phi \mathbf{1}_t^* \mathbf{1}_t & -\mathbf{1}_t^* \\ -\mathbf{1}_t & -\frac{1}{2\phi} \text{id} \end{pmatrix}. \quad (9.1)$$

The proofs of Lemma 8.1 and Proposition 8.2 will use the following auxiliary lemmas.

Lemma 9.1. *There exists a constant $C > 0$ such that,*

$$\sup_{t \in [0, T]} \|\mathbf{\Gamma}_t^\theta - \mathbf{\Gamma}_t^{\theta'}\|_{\text{op}} \leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}), \quad \text{for all } \theta, \theta' \in \Xi_\varepsilon.$$

Proof. From (9.1) it follows that it is enough to prove that there exists a constant $C > 0$ such that,

$$\sup_{t \in [0, T]} \|\mathbf{D}_t^\theta - \mathbf{D}_t^{\theta'}\|_{\text{op}} \leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}), \quad \text{for all } \theta, \theta' \in \Xi_\varepsilon. \quad (9.2)$$

From (4.7) we get,

$$\mathbf{D}_t^\theta - \mathbf{D}_t^{\theta'} := 2(\lambda - \lambda') \text{id} + (\tilde{\mathbf{G}}_t^\theta - \tilde{\mathbf{G}}_t^{\theta'} + (\tilde{\mathbf{G}}_t^\theta)^* - (\tilde{\mathbf{G}}_t^{\theta'})^*) \mathbf{1}_t.$$

Note that for $\theta = (\lambda, G)$, using (4.5) and (4.6), the kernel of $\tilde{\mathbf{G}}_t^\theta$ is given by,

$$\tilde{G}_t^\theta(s, u) = (2\varrho + G(s - u)) \mathbf{1}_{\{u < s\}} \mathbf{1}_{\{u > t\}}. \quad (9.3)$$

Hence

$$\begin{aligned} (\mathbf{D}_t^\theta - \mathbf{D}_t^{\theta'}) f(s) &= 2(\lambda - \lambda') f(s) + \int_0^T (G(s - u) - G'(s - u)) \mathbf{1}_{\{s > u > t\}} f(u) du \\ &\quad + \int_0^T (G(u - s) - G'(u - s)) \mathbf{1}_{\{u > s > t\}} f(u) du. \end{aligned} \quad (9.4)$$

From (9.4) and (8.1) and an application of Cauchy-Schwarz inequality we get (9.2). \square

Lemma 9.2. *For any $\theta \in \Xi_\varepsilon$ the operator \mathbf{D}_t^θ is positive definite, self-adjoint, invertible and moreover we have*

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \langle f, \mathbf{D}_t^\theta f \rangle \geq (L^{-1} - 2\varepsilon) \|f\|_{L^2([0, T], \mathbb{R})}^2, \quad \text{for all } f \in L^2([0, T], \mathbb{R}). \quad (9.5)$$

Proof. We first note that from (4.7) it follows that \mathbf{D}_t^θ is a self-adjoint operator. We will prove (9.5), which under the condition of Definition 2.14 implies that \mathbf{D}_t^θ is positive definite, hence it is invertible.

Recall that \tilde{G}^θ was defined in (4.5) and that $\tilde{\mathbf{G}}_t^\theta$ is the operator induced by the kernel $\tilde{G}^\theta(s, u) \mathbf{1}_{\{u \geq t\}}$. Since $\lambda > 0$, the operator λid is positive definite and the operator $\mathbf{1}_t^* \mathbf{1}_t$ is nonnegative definite. It follows from (4.7) that in order to prove that \mathbf{D}_t^θ satisfies (9.5) we need to derive a lower bound on $(\tilde{\mathbf{G}}_t^\theta + (\tilde{\mathbf{G}}_t^\theta)^*)$. Let $f \in L^2([0, T], \mathbb{R})$. Repeating the same steps as in the proof of Lemma 4.1 in [1] we get

$$\begin{aligned} & \int_0^T \int_0^T (G_t^\theta(s, u) + (G_t^\theta)^*(s, u)) f(s) f(u) ds du \\ &= \int_0^T \int_0^T 2G(|s - u|) f_t(s) f_t(u) ds du \\ &\geq -2\varepsilon \|f\|_{L^2([0, T], \mathbb{R})}^2, \quad \text{for all } t \in [0, T], \theta \in \Xi_\varepsilon, \end{aligned} \quad (9.6)$$

where we used Definition 2.14 in the last inequality.

Moreover, we have

$$\begin{aligned} & \int_t^T \int_u^T f(s) f(u) ds du + \int_0^T \int_0^T \mathbf{1}_{\{t \leq s \leq u\}} f(u) f(s) du ds \\ &= \int_t^T \int_u^T f(s) f(t) ds dt + \int_t^T \int_t^u f(u) f(s) ds du \\ &= \int_t^T \int_t^T f(s) f(u) ds du = \left(\int_t^T f(s) ds \right)^2 \geq 0. \end{aligned} \quad (9.7)$$

From (9.3), (9.6) and (9.7) it follows that $\tilde{\mathbf{G}}_t^\theta$ satisfies,

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \langle f, (\tilde{\mathbf{G}}_t^\theta + (\tilde{\mathbf{G}}_t^\theta)^*) f \rangle \geq -2\varepsilon \|f\|_{L^2([0, T], \mathbb{R})}^2, \quad \text{for all } f \in L^2([0, T], \mathbb{R}). \quad (9.8)$$

Since $\mathbf{1}_t^* \mathbf{1}_t$ is nonnegative definite and by Definition 2.14, $\lambda \geq L^{-1} > 2\varepsilon$, (9.5) follows from (4.7) and (9.8). \square

Now we are ready to prove Lemma 8.1 and Proposition 8.2.

Proof of Lemma 8.1 . Note that from (4.9) it follows that it enough to show that

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \|(\mathbf{D}_t^\theta)^{-1}\|_{\text{op}} < \infty. \quad (9.9)$$

Considering (9.5), we choose $\delta \in (0, L^{-1} - 2\varepsilon)$ and define the operator

$$\mathbf{S}_t^\theta := (2\lambda - \delta)\text{id} + (\tilde{\mathbf{G}}_t^\theta + (\tilde{\mathbf{G}}_t^\theta)^*) + 2\phi \mathbf{1}_t^* \mathbf{1}_t, \quad 0 \leq t \leq T. \quad (9.10)$$

From Lemma 9.2 it follows that there exists $\bar{\delta} > 0$ such that

$$\sup_{\theta \in \Xi_\varepsilon} \sup_{t \in [0, T]} \langle f, \mathbf{S}_t^\theta f \rangle > \bar{\delta} \|f\|_{L^2([0, T], \mathbb{R})}^2, \quad \text{for all } f \in L^2([0, T], \mathbb{R}). \quad (9.11)$$

and in particular \mathbf{S}_t^θ is positive definite, invertible, self-adjoint and compact with respect to the space of bounded operators on $L^2([0, T])$, equipped with the operator norm given in (8.1). From Theorem 4.15 in [44] it follows that \mathbf{S}_t^θ admits a spectral decomposition in terms of a sequence of positive eigenvalues $(\mu_{t,n})_{n=1}^\infty$ and an orthonormal sequence of eigenvectors $(\varphi_{t,n}^\theta)_{n=1}^\infty$ in $L^2([0, T])$, such that

$$\mathbf{S}_t^\theta = \sum_k \mu_{t,k}^\theta \langle \varphi_{t,k}^\theta, \cdot \rangle_{L^2} \varphi_{t,k}^\theta.$$

By application of Cauchy Schwarz and the fact that \mathbf{S}_t^θ is self-adjoint we get for any $\theta \in \Xi_\varepsilon$,

$$\begin{aligned} \sup_{t \leq T} \sum_k (\mu_{t,k}^\theta)^2 &\leq C \left((2\lambda - \delta)^2 + \sup_{t \leq T} \int_0^T \left((\tilde{G}_t^\theta + (\tilde{G}_t^\theta)^*) + 2\phi \mathbf{1}_t^* \mathbf{1}_t \right)^2 (s, s) ds \right) \\ &< \infty, \end{aligned}$$

where the second inequality follows from Definition 2.14 and (4.5). From (4.7) and (9.10) it follows that we can rewrite $\mathbf{D}_t^\theta = \mathbf{S}_t^\theta + \delta \text{id}$ as follows,

$$\mathbf{D}_t^\theta = \sum_k (\delta + \mu_{t,k}^\theta) \langle \varphi_{t,k}^\theta, \cdot \rangle_{L^2} \varphi_{t,k}^\theta.$$

We can therefore represent $(\mathbf{D}_t^\theta)^{-1}$ as follows,

$$(\mathbf{D}_t^\theta)^{-1} = \sum_k \frac{1}{(\delta + \mu_{t,k}^\theta)} \langle \varphi_{t,k}^\theta, \cdot \rangle_{L^2} \varphi_{t,k}^\theta.$$

Since $\delta > 0$ and $\mu_{t,k}^\theta \geq 0$, for all $t \in [0, T]$, $\theta \in \Xi_\varepsilon$ and $k = 1, 2, \dots$, we get that for any $f \in L^2([0, T], \mathbb{R})$,

$$\|(\mathbf{D}_t^\theta)^{-1} f\|_{L^2} \leq \frac{1}{\delta} \|f\|_{L^2}, \quad \text{for all } 0 \leq t \leq T, \theta \in \Xi_\varepsilon.$$

Together with (8.1) we get (9.9) and this completes the proof. \square

Proof of Proposition 8.2. We observe that

$$(\mathbf{\Gamma}_t^\theta)^{-1} - (\mathbf{\Gamma}_t^{\theta'})^{-1} = (\mathbf{\Gamma}_t^\theta)^{-1} (\mathbf{\Gamma}_t^\theta - \mathbf{\Gamma}_t^{\theta'}) (\mathbf{\Gamma}_t^{\theta'})^{-1}. \quad (9.12)$$

By taking the operator norm on both sides of (9.12) and using Lemma 8.1 and then Lemma 9.1 we get for all $\theta, \theta' \in \Xi_\varepsilon$ and $t \in [0, T]$,

$$\left\| (\mathbf{\Gamma}_t^\theta)^{-1} - (\mathbf{\Gamma}_t^{\theta'})^{-1} \right\|_{\text{op}} \leq C \|\mathbf{\Gamma}_t^\theta - \mathbf{\Gamma}_t^{\theta'}\|_{\text{op}} \leq C (|\lambda - \lambda'| + \|G - G'\|_{L^2([0, T])}).$$

□

A Regression-based algorithm for signal estimation

Observe that the trading strategy (2.8) (i.e., the process a in (4.13)) involves the conditional process $(t, s) \mapsto \mathbb{E}[A_s \mid \mathcal{F}_t]$ of the signal process A on $\Delta_T = \{(t, s) \mid 0 \leq t \leq s \leq T\}$. As the conditional distribution of A is in general not known analytically, the section proposes a regression-based algorithm to estimate the process $(t, s) \mapsto \mathbb{E}[A_s \mid \mathcal{F}_t]$ based on observed signal trajectories. Since the agent's trading strategy does not affect the signals, we assume for simplicity that the signal estimation has been carried out separately from learning the price impact coefficients (λ^*, G^*) .

Throughout this section, we assume that there exist independent copies $(A^m)_{m \in \mathbb{N}}$ of A , and impose the regularity condition of the signal A as specified in Assumption A.1. To simplify the notation, we write $\|X\|_{L^p(\Omega)}$, $p > 1$, for the L^p -norm of a random variable $X : \Omega \rightarrow \mathbb{R}$, and write $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot \mid \mathcal{F}_t]$ for each $t \in [0, T]$.

Assumption A.1. *There exists $I : \Omega \times [0, T] \rightarrow \mathbb{R}$ such that $A_t = \int_0^t I_s ds$ for all $t \in [0, T]$, and I is Markov with respect to the filtration $(\mathcal{F}_t)_{t \in [0, T]}$. Moreover, there exists $\vartheta > 2$ and $L \geq 0$ such that $\sup_{t \in [0, T]} \|I_t\|_{L^\vartheta(\Omega)} < \infty$, and for all $0 \leq t \leq s \leq r \leq T$ and $x, y \in \mathbb{R}$, $\|I_t - I_s\|_{L^2(\Omega)} \leq L|t - s|^{1/2}$ and $|\mathbb{E}^{t,x}[I_r] - \mathbb{E}^{s,y}[I_r]| \leq L(|x - y| + |t - s|^{1/2}(1 + |x| + |y|))$, where $\mathbb{E}^{t,x}[I_r] := \mathbb{E}[I_r \mid I_t = x]$.*

Remark A.2. *Assumption A.1 allows for non-Markovian signals A (with respect to $(\mathcal{F}_t)_{t \in [0, T]}$), but requires the time derivative of A to be Markov with a sufficiently regular transition kernel. This assumption includes as special cases the signal processes in [36, 40], where I is an Ornstein–Uhlenbeck process. More generally, Assumption A.1 holds if I solves a jump-diffusion stochastic differential equation with sufficiently regular (e.g., Lipschitz continuous [16, Theorem 4.1.1]) coefficients.*

Least-squares Monte Carlo for signal estimation. By Assumption A.1, for each $m \in \mathbb{N}$ and $t \in [0, T]$, $A_t^m = \int_0^t I_s^m ds$, where $(I^m)_{m \in \mathbb{N}}$ are independent copies of the Markov process I . In the sequel, we approximate $(t, s) \mapsto \mathbb{E}_t[A_s] = A_t + \int_t^s \mathbb{E}_t[I_r] dr$

by constructing a regression-based estimator of $\Delta_T \ni (t, s) \mapsto \mathbb{E}[I_s \mid I_t] \in \mathbb{R}$ based on $(I^m)_{m \in \mathbb{N}}$. We first discretise $(t, s) \mapsto \mathbb{E}_t[I_s]$ in time. More precisely, for each $N \in \mathbb{N}$, consider the grid $\pi_N = \{t_i\}_{i=0}^N \in \mathcal{P}_{[0,T]}$ with $t_i = iT/N$ for all $0 \leq i \leq N$, and define the following approximation of $\mathbb{E}_t[A_s]$: for all $(t, s) \in \Delta_T$,

$$\mathbb{E}_t[A_s] \approx A_t + \int_t^s \mathcal{I}_{t,r}^N dr, \quad \text{with } \mathcal{I}_{t,r}^N := \sum_{i,j=0}^{N-1} \mathbb{E}_{t_i}[I_{t_j}] \mathbb{1}_{[t_i, t_{i+1}) \times [t_j, t_{j+1})}(t, r), \quad (t, r) \in \Delta_T. \quad (\text{A.1})$$

The conditional expectations in (A.1) are then projected on prescribed basis functions via least-squares Monte Carlo methods (see e.g., [31]). To this end, let \mathcal{V} be a finite-dimensional vector space of functions $\psi : \mathbb{R} \rightarrow \mathbb{R}$, and for each $R \geq 0$, let $\mathcal{T}_R : \mathbb{R} \rightarrow \mathbb{R}$ be the truncation function such that $\mathcal{T}_R(x) = \max(-R, \min(x, R))$ for all $x \in \mathbb{R}$. Then for each $N, M \in \mathbb{N}$ and $R \geq 0$, we define the following truncated least-squares estimate of $\mathbb{E}_{t_i}[I_{t_j}]$: for all $0 \leq i \leq j \leq N$,

$$\psi_{i,j}(\cdot) = \mathcal{T}_R \tilde{\psi}_{i,j}(\cdot), \quad \text{with } \tilde{\psi}_{i,j} \in \arg \min_{\psi \in \mathcal{V}} \frac{1}{M} \sum_{m=1}^M |\psi(I_{t_i}^m) - \mathcal{T}_R(I_{t_j}^m)|^2. \quad (\text{A.2})$$

As \mathcal{V} is a vector space, $\tilde{\psi}_{i,j}$ in (A.2) can be computed by solving a least-squares problem over the weights of some fixed basis functions, whose solution may not be unique (see e.g., page 162 of [31]). Note that there can be more than one solution to (A.2). We then consider the approximation $\mathbb{E}_t[A_s] \approx A_t + \int_t^s \mathcal{I}_{t,r}^{N,M,R} dr$, with

$$\mathcal{I}_{t,r}^{N,M,R} = \sum_{i,j=0}^{N-1} \psi_{i,j}(I_{t_i}) \mathbb{1}_{[t_i, t_{i+1}) \times [t_j, t_{j+1})}(t, r), \quad (t, r) \in \Delta_T. \quad (\text{A.3})$$

The accuracy of (A.3) depends on the expressivity and complexity of the vector space \mathcal{V} (see Proposition A.5). By the Lipschitz continuity of the map $x \mapsto \mathbb{E}^{t_i, x}[I_{t_j}]$ (Assumption A.1), we set the vector space \mathcal{V} as the space of piecewise constant functions defined on a spatial grid. This allows for optimally balancing the expressivity and complexity of \mathcal{V} and obtaining a precise error estimate of (A.3) (see Theorem A.6). More precisely, for each $R \geq 0$ and $K \in \mathbb{N}$, let $\mathbf{p}_K := \{-\frac{R}{2} = x_0 < x_1 < \dots < x_K = \frac{R}{2}\}$ be a uniform grid of $[-\frac{R}{2}, \frac{R}{2}]$ such that $x_{i+1} - x_i = \frac{R}{K}$ for all $0 \leq i \leq K-1$, and let \mathcal{V}_K be the space of real-valued functions that are piecewise constant on the grid \mathbf{p}_K and zero outside $[-\frac{R}{2}, \frac{R}{2}]$. It is clear that \mathcal{V}_K is of the dimension K . Then for each $N, M, K \in \mathbb{N}$ and $R \geq 0$, consider the approximation $\mathbb{E}_t[A_s] \approx A_t + \int_t^s \mathcal{I}_{t,r}^{N,M,K,R} dr$, where the process $\Delta_T \ni (t, r) \mapsto \mathcal{I}_{t,r}^{N,M,K,R} \in \mathbb{R}$ is defined by (cf. (A.3)):

$$\mathcal{I}_{t,r}^{N,M,K,R} = \sum_{i,j=0}^{N-1} \psi_{i,j}(I_{t_i}) \mathbb{1}_{[t_i, t_{i+1}) \times [t_j, t_{j+1})}(t, r), \quad (t, r) \in \Delta_T, \quad (\text{A.4})$$

with $\psi_{i,j}$ being the truncated least-squares estimate (A.2) over the space $\mathcal{V} = \mathcal{V}_K$:

$$\psi_{i,j}(\cdot) = \mathcal{T}_R \tilde{\psi}_{i,j}(\cdot), \quad \text{with } \tilde{\psi}_{i,j} \in \arg \min_{\psi \in \mathcal{V}_K} \frac{1}{M} \sum_{m=1}^M |\psi(I_{t_i}^m) - \mathcal{T}_R(I_{t_j}^m)|^2.$$

Convergence rates of (A.3) and (A.4). To quantify the accuracy of (A.3) and (A.4), we start with the following technical lemma.

Lemma A.3. *Suppose that Assumption A.1 holds. Then there exists $C \geq 0$ such that for all $0 \leq t \leq s \leq r \leq T$, $\|\mathbb{E}_t[I_r] - \mathbb{E}_t[I_s]\|_{L^2(\Omega)} \leq C|r - s|^{1/2}$ and $\|\mathbb{E}_t[I_r] - \mathbb{E}_s[I_r]\|_{L^2(\Omega)} \leq C|t - s|^{1/2}$.*

Proof. Let $0 \leq t \leq s \leq r \leq T$, and $C \geq 0$ be a generic constant independent of t, s and r . Jensen's inequality and Assumption A.1 imply that $\mathbb{E}[|\mathbb{E}_t[I_r] - \mathbb{E}_t[I_s]|^2] \leq \mathbb{E}[|I_r - I_s|^2] \leq L^2|r - s|$. Moreover, by the Markov property of I and Assumption A.1,

$$|\mathbb{E}_t[I_r] - \mathbb{E}_s[I_r]| = |\mathbb{E}^{t, I_t}[I_r] - \mathbb{E}^{s, I_s}[I_r]| \leq L(|I_t - I_s| + |t - s|^{1/2}(1 + |I_t| + |I_s|)).$$

Taking L^2 -norm on both sides and apply Young's inequality yield

$$\mathbb{E}[|\mathbb{E}_t[I_r] - \mathbb{E}_s[I_r]|^2] \leq 2L^2(\mathbb{E}[|I_t - I_s|^2] + |t - s|\mathbb{E}[(1 + |I_t| + |I_s|)^2]) \leq C|t - s|,$$

where the last inequality used $\sup_{t \in [0, T]} \mathbb{E}[|I_t|^2] < \infty$. \square

The following proposition quantifies the time discretisation error of (A.1) based on Lemma A.3.

Proposition A.4. *Suppose that Assumption A.1 holds. Then there exists $C \geq 0$ such that*

$$\sup_{t \in [0, T]} \left\| \sup_{s \in [t, T]} \left| \mathbb{E}_t[A_s] - \left(A_t + \int_t^s \mathcal{I}_{t,r}^N dr \right) \right| \right\|_{L^2(\Omega)} \leq CN^{-\frac{1}{2}}, \quad \text{for all } N \in \mathbb{N}.$$

Proof. Throughout this proof, let $t \in [0, T)$ and $N \in \mathbb{N}$ be fixed and C be a generic constant independent of t and N . Let $t_i \in \pi_N$ such that $t \in [t_i, t_{i+1})$. By (A.1), for all $s \in [t, T]$,

$$\begin{aligned} \left| \int_t^s \mathbb{E}_t[I_r] dr - \int_t^s \mathcal{I}_{t,r}^N dr \right| &\leq \left| \int_t^s (\mathbb{E}_t[I_r] - \mathbb{E}_{t_i}[I_r] + \mathbb{E}_{t_i}[I_r] - \mathcal{I}_{t,r}^N) dr \right| \\ &\leq \int_t^s |\mathbb{E}_t[I_r] - \mathbb{E}_{t_i}[I_r]| dr + \int_t^s \left| \mathbb{E}_{t_i}[I_r] - \sum_{j=i}^{N-1} \mathbb{E}_{t_i}[I_{t_j}] \mathbb{1}_{[t_j, t_{j+1})}(r) \right| dr \\ &\leq \int_t^T |\mathbb{E}_t[I_r] - \mathbb{E}_{t_i}[I_r]| dr + \sum_{j=i}^{N-1} \int_{t_j}^{t_{j+1}} |\mathbb{E}_{t_i}[I_r] - \mathbb{E}_{t_i}[I_{t_j}]| dr. \end{aligned}$$

By taking the supremum over $s \in [t, T]$ and the L^2 -norm on both sides of the above estimate, and applying Lemma A.3,

$$\begin{aligned}
& \left\| \sup_{s \in [t, T]} \left| \mathbb{E}_t[A_s] - \left(A_t + \int_t^s \mathcal{I}_{t,r}^N dr \right) \right| \right\|_{L^2(\Omega)} \\
& \leq \int_t^T \|\mathbb{E}_t[I_r] - \mathbb{E}_{t_i}[I_r]\|_{L^2(\Omega)} dr + \sum_{j=i}^{N-1} \int_{t_j}^{t_{j+1}} \|\mathbb{E}_{t_i}[I_r] - \mathbb{E}_{t_i}[I_{t_j}]\|_{L^2(\Omega)} dr \\
& \leq C \int_t^T |t - t_i|^{1/2} dt + \sum_{j=i}^{N-1} \int_{t_j}^{t_{j+1}} |r - t_j|^{1/2} dr \leq C N^{-1/2},
\end{aligned}$$

due to the fact that $t \in [t_i, t_{i+1})$. Taking the supremum over $t \in [0, T]$ yields the desired estimate. \square

The following proposition quantifies the accuracy of $\left(A_t + \int_t^s \mathcal{I}_{t,r}^{N,M,R} dr \right)_{(t,s) \in \Delta_T}$ in terms of the number of time discretisation N , the sample size M , the truncation level R and the complexity of the function space \mathcal{V} .

Proposition A.5. *Suppose that Assumption A.1 holds. Then there exists $C \geq 0$ such that for all $N, M \in \mathbb{N}$, $R \geq 0$ and vector spaces \mathcal{V} of functions $\psi : \mathbb{R} \rightarrow \mathbb{R}$,*

$$\begin{aligned}
& \sup_{t \in [0, T]} \left\| \sup_{s \in [t, T]} \left| \mathbb{E}_t[A_s] - \left(A_t + \int_t^s \mathcal{I}_{t,r}^{N,M,R} dr \right) \right| \right\|_{L^2(\Omega)} \\
& \leq C \left(\frac{1}{\sqrt{N}} + R \sqrt{\frac{(\ln M + 1)n_{\mathcal{V}}}{M}} + \sup_{t \in [0, T]} \mathbb{E}[|I_t|^2 \mathbf{1}_{|I_t| \geq R}] + \max_{0 \leq i \leq j \leq N} \inf_{\psi \in \mathcal{V}} \|\mathbb{E}_{t_i}[I_{t_j}] - \psi(I_{t_i})\|_{L^2(\Omega)} \right),
\end{aligned}$$

where $n_{\mathcal{V}}$ is the vector space dimension of \mathcal{V} .

Proof. Throughout this proof, let $N, M, R \in \mathbb{N}$, $t \in [0, T)$ and \mathcal{V} be fixed, and let C be a generic constant independent of the above quantities. By Jensen's inequality,

$$\begin{aligned}
\|\mathbb{E}_{t_i}[\mathcal{T}_R(I_{t_j})] - \mathbb{E}_{t_i}[I_{t_j}]\|_{L^2(\Omega)}^2 & \leq \|\mathbb{E}_{t_i}[(|I_{t_j}| + R)\mathbf{1}_{|I_{t_j}| > R}]\|_{L^2(\Omega)}^2 \\
& \leq 4\mathbb{E}[|I_{t_j}|^2 \mathbf{1}_{|I_{t_j}| > R}].
\end{aligned} \tag{A.5}$$

Then by observing that $|\mathcal{T}_R(I_{t_j})| \leq R$ and applying [31, Theorem 11.3], there exists $C \geq 0$ such that for all $0 \leq i \leq j \leq N$,

$$\begin{aligned}
& \|\mathbb{E}_{t_i}[I_{t_j}] - \psi_{i,j}(I_{t_i})\|_{L^2(\Omega)}^2 \\
& \leq 2\|\mathbb{E}_{t_i}[I_{t_j}] - \mathbb{E}_{t_i}[\mathcal{T}_R(I_{t_j})]\|_{L^2(\Omega)}^2 + 2\|\mathbb{E}_{t_i}[\mathcal{T}_R(I_{t_j})] - \psi_{i,j}(I_{t_i})\|_{L^2(\Omega)}^2 \\
& \leq 8\mathbb{E}[|I_{t_j}|^2 \mathbf{1}_{|I_{t_j}| > R}] + C \left(R^2 \frac{(\ln M + 1)n_{\mathcal{V}}}{M} + \inf_{\psi \in \mathcal{V}} \|\mathbb{E}_{t_i}[\mathcal{T}_R(I_{t_j})] - \psi(I_{t_i})\|_{L^2(\Omega)}^2 \right) \tag{A.6} \\
& \leq C \left(R^2 \frac{(\ln M + 1)n_{\mathcal{V}}}{M} + \sup_{t \in [0, T]} \mathbb{E}[|I_t|^2 \mathbf{1}_{|I_t| \geq R}] + \inf_{\psi \in \mathcal{V}} \|\mathbb{E}_{t_i}[I_{t_j}] - \psi(I_{t_i})\|_{L^2(\Omega)}^2 \right),
\end{aligned}$$

where $n_{\mathcal{V}}$ is the vector space dimension of \mathcal{V} . Hence, let $i \in \mathbb{N}$ be such that $t \in [t_i, t_{i+1})$, by (A.1), (A.3) and (A.6),

$$\begin{aligned}
& \left\| \sup_{s \in [t, T]} \left| \int_t^s \mathcal{I}_{t,r}^N dr - \int_t^s \mathcal{I}_{t,r}^{N,M,R} dr \right| \right\|_{L^2(\Omega)} \leq \int_t^T \|\mathcal{I}_{t,r}^N - \mathcal{I}_{t,r}^{N,M,R}\|_{L^2(\Omega)} dr \\
& \leq \int_{t_i}^T \left\| \sum_{j=i}^{N-1} \mathbb{E}_{t_i}[I_{t_j}] \mathbb{1}_{[t_j, t_{j+1})}(r) - \sum_{j=i}^{N-1} \psi_{i,j}(I_{t_i}) \mathbb{1}_{[t_j, t_{j+1})}(r) \right\|_{L^2(\Omega)} dr \\
& \leq \sum_{j=i}^{N-1} \int_{t_j}^{t_{j+1}} \|\mathbb{E}_{t_i}[I_{t_j}] - \psi_{i,j}(I_{t_i})\|_{L^2(\Omega)} dr \\
& \leq C \left(R \left(\frac{(\ln M + 1)n_{\mathcal{V}}}{M} \right)^{\frac{1}{2}} + \sup_{t \in [0, T]} \mathbb{E}[|I_t|^2 \mathbb{1}_{|I_t| \geq R}] + \max_{0 \leq i \leq j \leq N} \inf_{\psi \in \mathcal{V}} \|\mathbb{E}_{t_i}[I_{t_j}] - \psi(I_{t_i})\|_{L^2(\Omega)} \right),
\end{aligned}$$

which along with Proposition A.4 leads to the desired estimate. \square

The following theorem simplifies the error bound in Proposition A.5 with \mathcal{V} being the space of piecewise constant functions. It specifies the precise dependence of the hyperparameters N, K, R on the sample size M , and establishes the convergence rate of $\left(A_t + \int_t^s \mathcal{I}_{t,r}^{N,M,K,R} dr \right)_{(t,s) \in \Delta_T}$ as M tends to infinity.

Theorem A.6. *Suppose that Assumption A.1 holds. Then there exists $C \geq 0$ such that for all $N, M, K \in \mathbb{N}$ and $R \geq 0$,*

$$\begin{aligned}
& \sup_{t \in [0, T]} \left\| \sup_{s \in [t, T]} \left| \mathbb{E}_t[A_s] - \left(A_t + \int_t^s \mathcal{I}_{t,r}^{N,M,K,R} dr \right) \right| \right\|_{L^2(\Omega)} \\
& \leq C \left(\frac{1}{\sqrt{N}} + R \sqrt{\frac{(\ln M + 1)K}{M}} + \frac{R}{K} + \frac{1}{R^{\frac{\vartheta-2}{2}}} \right).
\end{aligned} \tag{A.7}$$

Consequently, if one sets $(N_M)_{M \in \mathbb{N}}, (K_M)_{M \in \mathbb{N}} \subset \mathbb{N}$, and $(R_M)_{M \in \mathbb{N}} \subset (0, \infty)$ such that

$$N_M \sim \left(\frac{M}{\ln M + 1} \right)^{\frac{2}{3}}, \quad K_M \sim \left(\frac{M}{\ln M + 1} \right)^{\frac{1}{3}}, \quad R_M \sim \left(\frac{M}{\ln M + 1} \right)^{\frac{2}{3\vartheta}},^2$$

then for all $M \in \mathbb{N}$,

$$\sup_{t \in [0, T]} \left\| \sup_{s \in [t, T]} \left| \mathbb{E}_t[A_s] - \left(A_t + \int_t^s \mathcal{I}_{t,r}^{N_M, M, K_M, R_M} dr \right) \right| \right\|_{L^2(\Omega)} \leq C \left(\frac{\ln M + 1}{M} \right)^{\frac{\vartheta-2}{3\vartheta}}. \tag{A.8}$$

²For any sequences $(a_M)_{M \in \mathbb{N}}, (b_M)_{M \in \mathbb{N}} \subset (0, \infty)$, we write $a_M \sim b_M$ if $0 < \liminf_{M \rightarrow \infty} \frac{a_M}{b_M} \leq \limsup_{M \rightarrow \infty} \frac{a_M}{b_M} < \infty$.

Proof. Throughout this proof, let $N, M, K \in \mathbb{N}$, $R \geq 0$ and $t \in [0, T)$ be fixed, and let C be a generic constant independent of the above quantities. For each $0 \leq i \leq j \leq N$, let $\varphi_{i,j} : \mathbb{R} \rightarrow \mathbb{R}$ be such that $\varphi_{i,j}(x) = \mathbb{E}^{t_i, x}[I_{t_j}]$ for all $x \in \mathbb{R}$. By Assumption A.1, $|\varphi_{i,j}(x) - \varphi_{i,j}(y)| \leq L|x - y|$ for all $x, y \in \mathbb{R}$. For any given $0 \leq i \leq j \leq N$, by considering $\beta_\ell = \varphi_{i,j}(x_\ell)$ for all $\ell = 0, \dots, K-1$,

$$\begin{aligned}
\inf_{\psi \in \mathcal{V}_K} \|\mathbb{E}_{t_i}[I_{t_j}] - \psi(I_{t_i})\|_{L^2(\Omega)} &\leq \left\| \mathbb{E}_{t_i}[I_{t_j}] - \sum_{\ell=0}^{K-1} \tilde{\beta}_\ell \mathbf{1}_{[x_\ell, x_{\ell+1})}(I_{t_i}) \right\|_{L^2(\Omega)} \\
&\leq \left\| \varphi_{i,j}(I_{t_i}) - \varphi_{i,j}(I_{t_i}) \mathbf{1}_{|x| \leq \frac{R}{2}}(I_{t_i}) \right\|_{L^2(\Omega)} \\
&\quad + \left\| \varphi_{i,j}(I_{t_i}) \mathbf{1}_{|x| \leq \frac{R}{2}}(I_{t_i}) - \sum_{\ell=0}^{K-1} \varphi_{i,j}(x_\ell) \mathbf{1}_{[x_\ell, x_{\ell+1})}(I_{t_i}) \right\|_{L^2(\Omega)} \\
&\leq \|\varphi_{i,j}(I_{t_i}) \mathbf{1}_{|I_{t_i}| > \frac{R}{2}}\|_{L^2(\Omega)} + \left\| \sum_{\ell=0}^{K-1} |\varphi_{i,j}(I_{t_i}) - \varphi_{i,j}(x_\ell)| \mathbf{1}_{[x_\ell, x_{\ell+1})}(I_{t_i}) \right\|_{L^2(\Omega)} \\
&\leq \|I_{t_j} \mathbf{1}_{|I_{t_i}| > \frac{R}{2}}\|_{L^2(\Omega)} + \frac{LR}{K},
\end{aligned}$$

where the first term in the last inequality used $\varphi_{i,j}(I_{t_i}) \mathbf{1}_{|I_{t_i}| > R/2} = \mathbb{E}_{t_i}[I_{t_j} \mathbf{1}_{|I_{t_i}| > R/2}]$ and Jensen's inequality, and the second term used the L -Lipschitz continuity of $\varphi_{i,j}$. By Hölder's inequality (with $p = \frac{\vartheta}{2} > 1$) and Markov's inequality,

$$\begin{aligned}
\left\| I_{t_j} \mathbf{1}_{|I_{t_i}| > \frac{R}{2}} \right\|_{L^2(\Omega)} &\leq \|I_{t_j}\|_{L^\vartheta(\Omega)} \mathbb{P}(|I_{t_i}| > \frac{R}{2})^{\frac{\vartheta-2}{2\vartheta}} \leq \|I_{t_j}\|_{L^\vartheta(\Omega)} \left(\frac{2^\vartheta \|I_{t_i}\|_{L^\vartheta(\Omega)}^\vartheta}{R^\vartheta} \right)^{\frac{\alpha-2}{2\alpha}} \\
&\leq 2^{\frac{\vartheta-2}{2}} R^{-\frac{\vartheta-2}{2}} \sup_{t \in [0, T]} \|I_t\|_{L^\vartheta(\Omega)}^{\frac{\vartheta}{2}}.
\end{aligned} \tag{A.9}$$

This along with Proposition A.5 and $n_{\mathcal{V}_K} = K$ leads to the desired estimate (A.7).

Finally, let $\mathfrak{C}_\vartheta = 2^{\frac{\vartheta-2}{2}} \sup_{t \in [0, T]} \|I_t\|_{L^\vartheta(\Omega)}^{\frac{\vartheta}{2}} < \infty$. Observe that if

$$K_M \sim \sqrt{N}, \quad R_M \sim \frac{\mathfrak{C}_\vartheta^{2/\vartheta}}{\left(\sqrt{\frac{(\ln M + 1)\sqrt{N}}{M}} + \frac{1}{\sqrt{N}} \right)^{2/\vartheta}}, \quad N_M \sim \left(\frac{M}{\ln M + 1} \right)^{2/3}, \tag{A.10}$$

then the following error estimate holds: for all $M \in \mathbb{N}$,

$$\begin{aligned}
& \frac{1}{\sqrt{N_M}} + R_M \sqrt{\frac{(\ln M + 1)K_M}{M}} + \frac{R_M}{K_M} + \frac{\mathfrak{C}_\vartheta}{R_M^{\frac{\vartheta-2}{2}}} \\
& \leq C \left(R_M \left(\sqrt{\frac{(\ln M + 1)\sqrt{N_M}}{M}} + \frac{1}{\sqrt{N_M}} \right) + \frac{\mathfrak{C}_\vartheta}{R_M^{\frac{\vartheta-2}{2}}} \right) \\
& \leq C \mathfrak{C}_\vartheta^{\frac{2}{\vartheta}} \left(\sqrt{\frac{(\ln M + 1)\sqrt{N_M}}{M}} + \frac{1}{\sqrt{N_M}} \right)^{1-\frac{2}{\vartheta}} \leq C \mathfrak{C}_\vartheta^{\frac{2}{\vartheta}} \left(\frac{\ln M + 1}{M} \right)^{\frac{\vartheta-2}{3\vartheta}}.
\end{aligned}$$

The desired estimate (A.8) follows from the fact that $1 \leq (2^{\frac{\vartheta-2}{2}})^{\frac{2}{\vartheta}} \leq 2$ for all $\vartheta > 2$, and the observation that the choices of hyperparameters N, K, R in the statement satisfies the criterion (A.10). \square

Acknowledgments

We are very grateful to the Associate Editor and to the anonymous referees for careful reading of the manuscript and for a number of useful comments and suggestions that significantly improved this paper.

Funding. Not applicable.

Availability of data and material. Not applicable.

Compliance with ethical standards. The authors have no conflicts of interest to declare that are relevant to the content of this article.

Code availability. Not applicable.

References

- [1] E. Abi Jaber and E. Neuman. Optimal liquidation with signals: the general propagator case. *arXiv:2211.00447*, 2022.
- [2] R. Almgren. Real time trading signals, 2018. URL <https://www.youtube.com/watch?v=s4IdoWUhrDA>.
- [3] R. Almgren and N. Chriss. Value under liquidation. *Risk*, 12:61–63, 1999.
- [4] R. Almgren and N. Chriss. Optimal execution of portfolio transactions. *Journal of Risk*, 3(2):5–39, 2000.

- [5] M. Basei, X. Guo, A. Hu, and Y. Zhang. Logarithmic regret for episodic continuous-time linear-quadratic reinforcement learning over a finite-time horizon. *Journal of Machine Learning Research*, 23(178):1–34, 2022.
- [6] C. Belak, J. Muhle-Karbe, and K. Ou. Liquidation in target zone models. *Market Microstructure and Liquidity*, 2019. URL <https://doi.org/10.1142/S2382626619500102>.
- [7] C. Bellani, D. Brigo, A. Done, and E. Neuman. Static vs adaptive strategies for optimal execution with signals. *arXiv:1811.11265*, 2018.
- [8] D. Benatia, M. Carrasco, and J-P. Florens. Functional linear regression with functional response. *Journal of Econometrics*, 201(2):269–291, 2017. ISSN 0304-4076. doi: <https://doi.org/10.1016/j.jeconom.2017.08.008>. URL <https://www.sciencedirect.com/science/article/pii/S0304407617301586>. THEORETICAL AND FINANCIAL ECONOMETRICS: ESSAYS IN HONOR OF C. GOURIEROUX.
- [9] G. Blanchard and N. Mücke. Optimal rates for regularization of statistical inverse learning problems. *Foundations of Computational Mathematics*, 18(4):971–1013, 2018.
- [10] J.-P. Bouchaud, Y. Gefen, M. Potters, and M. Wyart. Fluctuations and response in financial markets: the subtle nature of ‘random’ price changes. *Quantitative finance*, 4(2):176–190, 2004.
- [11] J-P. Bouchaud, J. Bonart, J. Donier, and M. Gould. *Trades, Quotes and Prices: Financial Markets Under the Microscope*. Cambridge University Press, 2018. doi: 10.1017/9781316659335.
- [12] Á. Cartea and S. Jaimungal. Incorporating order-flow into optimal execution. *Mathematics and Financial Economics*, 10(3):339–364, 2016. ISSN 1862-9660. doi: 10.1007/s11579-016-0162-z. URL <http://dx.doi.org/10.1007/s11579-016-0162-z>.
- [13] Á. Cartea, S. Jaimungal, and J. Penalva. *Algorithmic and High-Frequency Trading (Mathematics, Finance and Risk)*. Cambridge University Press, 1 edition, October 2015. ISBN 1107091144. URL <http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/1107091144>.
- [14] Y. Chen, U. Hort, and H.H. Tran. Portfolio liquidation under transient price impact - theoretical solution and implementation with 100 nasdaq stocks. *arXiv preprint arXiv:1912.06426*, 2019.

- [15] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20(4):633–679, 2020.
- [16] L. Delong. *Backward stochastic differential equations with jumps and their actuarial and financial applications*. Springer, 2013.
- [17] W. Stockinger E. Neuman and Y. Zhang. An offline learning approach to propagator models. *arXiv:2309.02994*, 2023.
- [18] I. Ekeland and R. Témam. *Convex Analysis and Variational Problems*. Society for Industrial and Applied Mathematics, 1999. doi: 10.1137/1.9781611971088. URL <http://epubs.siam.org/doi/abs/10.1137/1.9781611971088>.
- [19] H. W. Engl, M. Hanke, and A. Neubauer. *Regularization of inverse problems*, volume 375. Springer Science & Business Media, 1996.
- [20] M. Forde, L. Sánchez-Betancourt, and B. Smith. Optimal trade execution for gaussian signals with power-law resilience. *Quantitative Finance*, 22(3):585–596, 03 2022. doi: 10.1080/14697688.2021.1950919. URL <https://doi.org/10.1080/14697688.2021.1950919>.
- [21] X. Gao and X. Y. Zhou. Logarithmic regret bounds for continuous-time average-reward markov decision processes. *arXiv preprint arXiv:2205.11168*, 2022.
- [22] X. Gao and X. Y. Zhou. Square-root regret bounds for continuous-time episodic markov decision processes. *arXiv preprint arXiv:2210.00832*, 2022.
- [23] J. Gatheral. No-dynamic-arbitrage and market impact. *Quantitative finance*, 10(7):749–759, 2010.
- [24] J. Gatheral and A. Schied. Dynamical models of market impact and algorithms for order execution. In Jean-Pierre Fouque and Joseph Langsam, editors, *Handbook on Systemic Risk*, pages 579–602. Cambridge University Press, 2013.
- [25] J. Gatheral, A. Schied, and A. Slynko. Transient linear price impact and Fredholm integral equations. *Math. Finance*, 22:445–474, 2012.
- [26] S. Gökay, A. Roch, and H.M. Soner. Liquidity models in continuous and discrete time. In Giulia di Nunno and Bern Øksendal, editors, *Advanced Mathematical Methods for Finance*, pages 333–366. Springer-Verlag, 2011.
- [27] G. Gripenberg, S.O. Londen, and O. Staffans. *Volterra integral and functional equations*. Number 34. Cambridge University Press, 1990.
- [28] C. W. Groetsch. The theory of Tikhonov regularization for Fredholm equations. 104p, *Boston Pitman Publication*, 1984.

- [29] O. Guéant. *The Financial Mathematics of Market Liquidity*. New York: Chapman and Hall/CRC, 2016.
- [30] X. Guo, A. Hu, and Y. Zhang. Reinforcement learning for linear-convex models with jumps via stability analysis of feedback controls. *SIAM Journal on Control and Optimization*, 61(2):755–787, 2023.
- [31] L. Györfi, M. Kohler, A. Krzyzak, H. Walk, et al. *A distribution-free theory of nonparametric regression*, volume 1. Springer, 2002.
- [32] J. He, D. Zhou, and Q. Gu. Logarithmic regret for reinforcement learning with linear function approximation. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4171–4180. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/he21c.html>.
- [33] D. Hofmann, D. Düvelmeyer, and K. Krumbiegel. Approximate source conditions in Tikhonov regularization-new analytical results and some numerical studies. *Mathematical Modelling and Analysis*, 11(1):41–56, 2006.
- [34] T. Lattimore and C. Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [35] C. A. Lehalle and O. Mounjid. Limit Order Strategic Placement with Adverse Selection Risk and the Role of Latency, October 2016. URL <http://arxiv.org/abs/1610.00261>.
- [36] C. A. Lehalle and E. Neuman. Incorporating signals into optimal trading. *Finance and Stochastics*, 23(2):275–311, 2019. doi: 10.1007/s00780-019-00382-7. URL <https://doi.org/10.1007/s00780-019-00382-7>.
- [37] C. A. Lehalle, S. Laruelle, R. Burgot, S. Pelin, and M. Lasnier. *Market Microstructure in Practice*. World Scientific publishing, 2013. URL <http://www.worldscientific.com/worldscibooks/10.1142/8967>.
- [38] A. Lipton, U. Pesavento, and M. G. Sotiropoulos. Trade arrival dynamics and quote imbalance in a limit order book, December 2013. URL <http://arxiv.org/abs/1312.0514>.
- [39] E. Neuman and A. Schied. Optimal portfolio liquidation in target zone models and catalytic superprocesses. *Finance and Stochastics*, 20:495–509, 2016.
- [40] E. Neuman and M. Voss. Optimal signal-adaptive trading with temporary and transient price impact. *SIAM Journal on Financial Mathematics*, 13(2):551–575, 2022.

- [41] E. Neuman and M. Voss. Trading with the crowd. *Mathematical Finance*, 33(3):548–617, 2023. doi: <https://doi.org/10.1111/mafi.12390>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/mafi.12390>.
- [42] A. A. Obizhaeva and J. Wang. Optimal trading strategy and supply/demand dynamics. *Journal of Financial Markets*, 16(1):1 – 32, 2013. ISSN 1386-4181. doi: <http://dx.doi.org/10.1016/j.finmar.2012.09.001>. URL <http://www.sciencedirect.com/science/article/pii/S1386418112000328>.
- [43] I. Osband and B. Van Roy. Model-based reinforcement learning and the eluder dimension. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. URL https://proceedings.neurips.cc/paper_files/paper/2014/file/1141938ba2c2b13f5505d7c424ebae5f-Paper.pdf.
- [44] D. Porter and D. S. G. Stirling. *Integral equations: A Practical Treatment, from Spectral Theory to Applications*. Cambridge University Press, 1990.
- [45] A. Kukanov R. Cont and S. Stoikov. The price impact of order book events. *Journal of Financial Econometrics*, 12(1):47–88, 2014.
- [46] Abhishake Rastogi, Gilles Blanchard, and Peter Mathé. Convergence analysis of Tikhonov regularization for non-linear statistical inverse problems. *Electronic Journal of Statistics*, 14(2):2798 – 2841, 2020. doi: 10.1214/20-EJS1735. URL <https://doi.org/10.1214/20-EJS1735>.
- [47] L. Schumaker. *Spline functions: basic theory*. Cambridge University Press, 2007.
- [48] L. Szpruch, T. Treetanthiploet, and Y. Zhang. Exploration-exploitation trade-off for continuous-time episodic reinforcement learning with linear-convex models. *To appear in Annals of Applied Probability*, 2024.
- [49] Lukasz Szpruch, Tanut Treetanthiploet, and Yufei Zhang. Optimal scheduling of entropy regularizer for continuous-time linear-quadratic reinforcement learning. *SIAM Journal on Control and Optimization*, 62(1):135–166, 2024.
- [50] D. E. Taranto, G. Bormetti, J-P. Bouchaud, F. Lillo, and B. Toth. Linear models for the impact of order flow on prices I. Propagators: Transient vs. History Dependent Impact, February 2016. URL <http://arxiv.org/abs/1602.02735>.
- [51] B. Tóth, Z. Eisler, and J. P. Bouchaud. The short-term price impact of trades is universal. *Market Microstructure and Liquidity*, 03(02):1850002, 2017. doi: 10.1142/S2382626618500028. URL <https://doi.org/10.1142/S2382626618500028>.

- [52] M. Vodret, I. Mastromatteo, B. Tóth, and M. Benzaquen. Do fundamentals shape the price response? a critical assessment of linear impact models. *Quantitative Finance*, 22(12):2139–2150, 2022. doi: 10.1080/14697688.2022.2114376. URL <https://doi.org/10.1080/14697688.2022.2114376>.
- [53] J. Yong and X.Y. Zhou. *Stochastic controls: Hamiltonian systems and HJB equations*, volume 43. Springer Science & Business Media, 1999.