

# Estimation of an Order Book Dependent Hawkes Process for Large Datasets\*

Luca Mucciante<sup>†</sup>      Alessio Sancetta<sup>‡</sup>

July 19, 2023

## Abstract

A point process for event arrivals in high frequency trading is presented. The intensity is the product of a Hawkes process and high dimensional functions of covariates derived from the order book. Conditions for stationarity of the process are stated. An algorithm is presented to estimate the model even in the presence of billions of data points, possibly mapping covariates into a high dimensional space. The large sample size can be common for high frequency data applications using multiple liquid instruments. Convergence of the algorithm is shown, consistency results under weak conditions is established, and a test statistic to assess out of sample performance of different model specifications is suggested. The methodology is applied to the study of four stocks that trade on the New York Stock Exchange (NYSE). The out of sample testing procedure suggests that capturing the nonlinearity of the order book information adds value to the self exciting nature of high frequency trading events.

**Key Words:** Counting process, forecast evaluation, high frequency trading, high dimensional estimation, one-hot encoding, trade arrival.

**JEL Codes:** C13; C32; C55.

## 1 Introduction

This paper presents an intensity model for event arrivals in high frequency trading. The intensity depends on order book information. The model is a Hawkes process where the intensity does not only depend on the time from an event arrival but also on the order book.

---

\*We are very grateful to the Editor Dacheng Xiu and the Referees for their comments that have led to substantial improvements both in content and presentation. We are also grateful to Francesco Cordonni (Royal Holloway), Giuliano De Rossi (Goldman Sachs) and Yuri Taranenko (ADIA) for useful conversations related to the topic of this paper.

<sup>†</sup>Department of Economics, Royal Holloway University of London, Egham TW20 0EX, UK. Email: luca-host@gmail.com.

<sup>‡</sup>Corresponding author. Department of Economics, Royal Holloway University of London, Egham TW20 0EX, UK. Email: asancetta@gmail.com.

The model is specifically designed for estimation in the presence of high dimensional covariates. Moreover, we present an estimation procedure that can deal with large datasets in a relatively simple way, using quadratic programming. High frequency data that include the order book may contain more than a million records in a day for a liquid instrument. When we consider information from other instruments, we may easily have millions of updates per day.

The problem of estimating the intensity of trading events conditional on order book and trade data relies to a single realization of trade and order book orders rather than a cross-section, as in the case of hazard models. The case when the number of conditioning variables is large has been studied recently (Sancetta, 2018, Mucciante and Sancetta, 2022). However, these references ignore the self exciting nature of event arrivals which is well documented by a number of authors (Bacry et al., 2015, Filimonov and Sornette, 2015 for reviews). The high dimensional model in Sancetta (2018) accounts for self excitement in the intensity in a simulation study, but does not provide a proof of the stationarity and ergodicity of the process.

Recently, the statistical properties of Hawkes processes that incorporate some information from the order book have been studied by a number of authors (inter alia, Fosset et al., 2020, Morariu-Patrichi and Pakkanen, 2022, Mounjid et al., 2019, Wu et al., 2019). These references define a probabilistic model for some order book information. For example they can be used to model the arrival of limit, market and cancel orders conditioning on the queue size of the order book. These models are general, but also rather complex when the dimension of the conditioning events grows. The state variables tend to be restricted to a finite set, which for practical reasons need to be low dimensional. Hence, applications usually focus on say one order book variable taking a finite number of discrete values. Hence, they are not suited for estimation, conditioning on a large information set. Moreover, the approach does not lend itself to the test of functional restriction on the impact that the order book variables can have on the intensity.

The focus of the current paper differs from the above in a number of ways. First, we allow for the covariates that capture the effect of the order book to take values in a subset of the real line rather than taking values in a finite state space. Second, the number of covariates can be large in the order of hundreds if not thousands. This is important, as the number of levels in the order book can be large. For examples, when considering the first ten levels of the order book, we have ten values for prices and ten for quantities for the bid and the ask, respectively. When we include information from additional instruments and add covariates that capture the time dynamics, it is easy to see how the information set can grow fast. Third, we chose a parametrization that can model nonlinearities mapping covariates into a higher dimensional space. This is particularly suited for high dimensional estimation and testing restrictions.

In summary, the main focus of the present paper is on modelling and consistent estimation, allowing for possibly complex nonlinear impact of order book variables and high dimensional information sets. The goal is not to model the order book, but to use information from the

order book as possible predictors when modelling the intensity of high frequency event arrivals. adds to the literature in a complementary way.

We now introduce the model. Let  $N := (N(t))_{t \geq 0}$  be the number of trade arrivals adapted to a filtration  $\mathcal{F} = (\mathcal{F}_t)_{t \geq 0}$ . The time of the  $j^{\text{th}}$  event arrival is denoted by  $T_j$ ,  $j \geq 1$  with  $T_0 := 0$ . The counting process admits an  $\mathcal{F}_t$ -adapted stochastic intensity  $\lambda_0(t)$  such that

$$\lambda_0(t) = h_0(t) g_0(t), \quad (1)$$

where  $h_0$  is the predictable process

$$h_0(t) = c_0 + \int_{(-\infty, t)} \left( \sum_{l=1}^L d_{0,l} e^{-a_{0,l}(t-s)} \right) dN(s) = c_0 + \sum_{j \geq 0: T_j < t} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}(t-T_j)} \quad (2)$$

with  $a_{0,l}, c_{0,l}, d_{0,l} \geq 0$ , and

$$g_0(t) = X(t)' b_0, \quad (3)$$

where  $b_0$  is a positive  $K \times 1$  vector,  $X := (X(t))_{t \geq 0}$  is a positive  $K \times 1$  dimensional left-continuous process; throughout, the prime symbol  $'$  stands for transposition. The positivity condition on  $b_0$  and  $X$  ensures that (1) is positive. The intensity can be interpreted as  $\lambda_0(t) = \lim_{s \downarrow 0} \Pr(N(t+s) - N(t) > 0 | \mathcal{F}_t) / s$  so that  $M(t) = N(t) - \int_{-\infty}^t \lambda_0(s) ds$  is an  $\mathcal{F}_t$ -martingale. In this paper  $K$  is assumed to be large possibly in the order of thousands. Throughout, for identification reasons,  $d_{0,1} := 1$ .

When (3) is constant, (1) reduces to a Hawkes process, where the kernel is the sum of exponential functions. It is well known that Hawkes processes provide a clear interpretation in terms of exogenous information arrival, captured by  $c_0$ , and endogenous market activity resulting from the second term in (2) (Hawkes and Oakes, 1974). However this interpretation misses any information on the trading environment. The latter is captured by the order book and other state variables as summarized in (3). We interpret (3) as the market environment. Section A.1 in the Appendix provides additional remarks on the model interpretation. Our interest is on the impact of the state variables on the system. This impact can be non-linear and we wish to capture these nonlinearities in a way that is easy to interpret. To do so, in applications,  $X$  may map a vector valued covariate process  $Z := (Z(t))_{t \geq 0}$  into a high dimensional space, where  $Z$  represents information from the order book. The empirical application of this paper focuses on a special type of mapping that is called one-hot encoding and is popular in machine learning due to its ease of interpretation (Alaya et al., 2019). Other examples include series expansions such as Bernstein polynomials, estimation in reproducing kernel Hilbert spaces etc. One hot-encoding essentially discretizes the variables and maps them into dummy variables. The cost of increasing the dimensionality of the problem comes at the advantage of having variables that are easy to interpret and to which linear constraints such

as monotonicity can be imposed in a natural way. For the model in (1) we state conditions for stationarity, ergodicity and strong mixing.

The model needs to be estimated under the constraints that we have briefly outlined: very large datasets, high dimensional parameters, and parameter’s restrictions. We present an estimation procedure that alternate between estimation of  $h_0$  for fixed  $g$  and uses surrogate loss function for estimation of  $g_0$  for fixed  $h$ . Mutatis mutandis, our approach is close in spirit to estimation by coordinate descent, which is standard in many problems (Friedman et al., 2007) and has sound theoretical properties (Beck and Tetrushvili, 2013). We then show that we can consistently estimate the high dimensional parameter  $b_0$  as long as  $T^{-1} \ln K \rightarrow 0$ , where  $[0, T]$  is the interval over which we collect the sample data.

The empirical application is based on Level 3 data for four stocks and the ETF on the S&P500 as auxiliary instrument. All data are traded on the NYSE. The data has been accessed through the Lobster dataset (Huang and Polak, 2011). Level 3 data allows us to correctly synchronize trades and order book. In the application, we estimate the model, and are interested in establishing the shape of the non-linear impact of variables derived from the order book. These include order book volume imbalance and spread among others. We also include information from other instruments.

We observe that our non-negativity constraint on  $b_0$  in (3) leads to a form of “shrinkage”. Most of the coefficients in the estimation are zero due to the non-negativity constraint. This is what we mean by “shrinkage”, and it does not require to explicitly employ a penalty on the coefficients. This phenomenon has been formally studied in Mucciante and Sancetta (2022). We shall briefly comment on this in Section 2.2.

## 1.1 Further Remarks on the Literature

The use of intensity models in high frequency econometrics was pioneered by Engle and Russell (1998). A review of different methodologies can be found in Bauwens and Hautsch (2009). As a result of data availability, there has been increasing interest in order book data on top of transaction data. High frequency trading strategies rely on order book features. For example, the literature has found that order book volume imbalances and other order book variables have an impact on price movements and trade arrivals at very short term horizons (Hall and Hautsch, 2007, Cont et al., 2014, Sancetta, 2018). MacKenzie (2017) reports anonymous interviews with ex algorithmic traders of the market maker Automated Trading Desk. These interviews confirm the importance of order book information for price movements. The empirical application confirms the importance of this within the context of our model.

The estimation of the intensity for hazard models conditioning on a large number of covariates has been considered by Gaïffas and Guilloux (2012). The main difference is that that the problem of estimating the intensity of trading events conditional on order book and

trades information relies on a single realization of trade and order book orders rather than a cross-section.

A model similar to (1) with (2) as baseline intensity has been considered in Sancetta (2018). However that approach cannot easily scale to a number  $m$  of covariate updates in the order of millions or more and is not as easily interpretable as the current method. When  $h_0$  is constant, Sancetta and Mucciante (2022) propose a consistent estimation method when the dimension  $K$  of  $b_0$  is large relatively to the number of jump events, with no need to impose a penalty. However, they require the variables to be linearly independent. This is not necessarily the case in many large scale problems. Moreover, they do not allow for estimation of a baseline intensity  $h_0$ . The method in that paper is complementary to the current one and provides theoretical insights into the shrinkage aspect of the procedure.

There is a rich literature that deals with estimation using large datasets, for example using distributed computing and averaging estimators across different smaller subsamples (Zhang et al., 2015). However, the structure of our problem is such that we can use all the data in one single estimation, at least as far as estimation of  $g_0$  in (3) is concerned, and it benefits from good asymptotic properties. Estimation of  $h_0$  is a low dimensional problem. Given  $g_0$ , estimation of  $h_0$  is equivalent to estimation of a Hawkes process. Despite being low dimensional, it may still pose challenges as first recognized by Ogata and Akaike (1982), and recently documented by Filimonov and Sornette (2015). Amongst other reasons, alternative procedures to likelihood estimation of Hawkes processes have been proposed in the literature (Da Fonseca and Zaatour, 2014, Kirchner, 2017, Cartea et al., 2021). We do not attempt to address these problems here. However, using our quadratic estimating function for  $h_0$ , estimates appeared relatively stable.

## 1.2 Outline of the Paper

Section 2 states the regularity conditions for estimation of the model as well as conditions for stationarity and ergodicity (Theorem 1). Section 3 discusses the estimation challenges and presents an algorithm to address these. Section 3.2 contains a simulation to assess the validity of the estimation algorithm and its finite sample performance. The focus is mostly on its convergence, but as a byproduct, we also show that it possesses the variable screening property under appropriate conditions. Asymptotic results concerning consistency when  $b_0$  is high dimensional are presented in Section 4. We extend the validity of the test procedure discussed in Sancetta (2018) to compare the fit of two competing intensities to the case of unbounded intensities. The model in (1) has unbounded intensity. The empirical application to four liquid stocks traded on the NYSE can be found in Section 5. There, we show that accounting for order book information and its dynamics adds value. Moreover, we find that the resulting impact of covariates derived from the order book is nonlinear and can be remarkably regular in some cases. Proofs and additional details can be found in the Appendix.

## 2 The Model

Here, we introduce the regularity conditions used for estimation of the model. We also consider a subset of these conditions to show that a slight generalization of our model is a strongly mixing stationary and ergodic process.

### 2.1 Regularity Conditions

We use  $|\cdot|_1$  to denote the  $\ell_1$  norm of a vector and  $|\cdot|_\infty$  for the uniform norm  $|g_0|_\infty = \sup_{t \geq 0} |g_0(t)|$ . In order to deduce stationarity of the process, we need some form of weak exogeneity of  $X(t)$  conditioning on  $(N_s)_{s < t}$ . We shall show that this is not restrictive. To introduce the exogeneity condition, we need some additional notation. Let  $\mathbb{M}$  be the space of Radon measures on  $\mathbb{R}$ . Define  $S_t$  to be the operator on  $\mathbb{M}$  that shifts mass  $t$  units to the left:  $S_t \nu(\cdot) = \nu(t + \cdot)$  for  $\nu \in \mathbb{M}$ . A map  $f$  from  $\mathbb{M}$  to the reals is causal if  $f(\nu) = f(\tilde{\nu})$  whenever  $\nu = \tilde{\nu}$  on  $(-\infty, 0)$ ,  $\nu, \tilde{\nu} \in \mathbb{M}$ . We write  $\tilde{\nu} \succeq \nu$  if  $\tilde{\nu}(C) \geq \nu(C)$  for any  $C \subset \mathbb{R}$ . Then, we call  $f : \mathbb{M} \rightarrow \mathbb{R}$  nondecreasing if  $f(\tilde{\nu}) \geq f(\nu)$  for any  $\tilde{\nu} \succeq \nu$  such that  $\nu, \tilde{\nu} \in \mathbb{M}$ . With slight abuse of notation, let  $N(C) = \int_C dN(t)$  for any  $C \subset \mathbb{R}$ . Hence, here we view  $N$  as an element in  $\mathbb{M}$ : the random measure associated to a point process. Let  $(W(t))_{t \geq 0}$  be a predictable stationary and ergodic process with values in  $\mathbb{R}^l$ , for some positive integer  $l$ , such that  $W(t)$  is independent of  $N(t + C)$  for any  $C \subseteq (-\infty, 0)$ . We refer to this as the independence condition.

**Condition 1** (*Weak Exogeneity*) *There is a strictly positive  $f_0 : \mathbb{R}^l \times \mathbb{M} \rightarrow \mathbb{R}$  such that for any  $x \in \mathbb{R}^l$ ,  $f_0(x, \cdot)$  is nondecreasing and causal and such that  $g_0(t) = f_0(W(t), S_t N)$ , where the latent process  $W$  is stationary and ergodic, and satisfies the independence condition.*

In the above, we do not need  $g_0(t) = X(t)' b_0$ . We now specialize the conditions for the purpose of estimation.

**Condition 2** (*True Model*) *The intensity of the process  $N$  is not identically zero, is as in (1),  $X := (X(t))_{t \geq 0}$  is left continuous and takes values in  $[0, 1]^K$  for every  $t \geq 0$ ,  $b_0$  has non-negative entries and is such that  $|b_0|_1 \leq B$ , where  $B \sum_{l=1}^L \frac{d_l}{a_l} < 1$  and  $d_{0,1} := 1$ .*

**Condition 3** (*Set  $\mathcal{G}$* )  $\mathcal{G} = \{g = X'b : |b|_1 \leq B, b_k \geq 0, k = 1, 2, \dots, K\}$ ,  $g_0 \in \mathcal{G}$ .

**Condition 4** (*Set  $\mathcal{H}$* ) *We have that  $\mathcal{H} := \{h_\psi : \psi \in \Psi\}$  where*

$$h_\psi(t) = c + \int_{(-\infty, t)} \left( \sum_{l=1}^L d_l e^{-a_l(t-s)} \right) dN(s),$$

$\psi = \{c, d, a\}$ ,  $\Psi = \mathcal{C} \cup \mathcal{D} \cup \mathcal{A}$  where  $\mathcal{C} \subset (0, \infty)$ ,  $\mathcal{D} \subset [0, \infty)^L$  and  $\mathcal{A} \subset (0, \infty)^L$  are compact sets that contain  $c_0$ ,  $d_{0,l}$  and  $a_{0,l}$  as in (2).

We refer to the above as the Regularity Conditions. Note that the true model is in the sets used for estimation. Here,  $B$  is a free parameter that needs to be tuned. For reasons to be discussed in Section 3.1, its choice is not crucial. In Theorem 1 in Section 2.3, we show that  $B \sum_{l=1}^L d_{0,l}/a_{0,l} < 1$  together with the weak exogeneity implies that the counting process is stationary.

## 2.2 Remarks on Regularity Conditions

We remark on regularity conditions.

**Condition 1.** We show how Condition 1 applies within our framework, i.e. when Condition 2 hold. To this end, we suppose that the  $k^{\text{th}}$  element in  $X(t)$  can be written as  $X_k(t) = f_{0,k}(W(t), S_t N)$  where  $f_{0,k} : \mathbb{R}^l \times \mathbb{M} \rightarrow \mathbb{R}$  such that for any  $x \in \mathbb{R}^l$ ,  $f_{0,k}(x, \cdot)$  is nondecreasing and causal,  $k = 1, 2, \dots, K$ . Then,  $g_0(t) = f_0(W(t), S_t N) = \sum_{k=1}^K b_k f_{0,k}(W(t), S_t N)$  satisfies Condition 1 because  $b_k \geq 0$ ,  $k = 1, 2, \dots, K$ . The exact functional form of  $f_{0,k}$  is not important for our purposes. However, the following is a simple example of the monotonicity and independence condition:  $f_{0,k}(W(t), S_t N) = \max\{\min\{1, W(t) + \nu(S_t N)\}, 0\}$  where  $\nu(S_t N) := \int_{(t-v, t)} dN(s)$  for some finite constant  $v > 0$ , and  $S_t N$  and  $W(t)$  are independent. Here,  $\nu(S_t N)$  counts the most recent number of events. Recall that we are assuming covariates in  $[0, 1]$ , hence  $X_k(t) := f_{0,k}(W(t), S_t N)$  is bounded accordingly.

Condition 1 says that we need to be able to decompose  $g_0(t)$  into a part that depends positively (non-negatively) on past event arrivals (monotonicity) and a component independent of these. Hence, Condition 1 rules out the case where the impact of  $X_k(t) b_{0,k}$  on the current intensity is positive, but the impact of past events on  $X_k(t) b_{0,k}$  is negative. This does not mean that quantities in the order book cannot reduce the intensity, as this is just a parametrization problem which is discussed in the remarks to Condition 2. Next we give an example.

For simplicity, with no loss of generality, let  $K = 1$ , so that  $X(t)' b_0 = X_1(t) b_{0,1}$ . Let  $X_1(t) = 1 - F(\text{Dur}(t))$  where  $\text{Dur}(t)$  is the last trade duration at time  $t$  and  $F(\cdot)$  is the distribution function of this duration. Then, for  $b_{0,1} > 0$  we have that a longer duration implies a smaller  $X_1(t) b_{0,1}$ . This makes sense, as a longer duration implies a smaller intensity. Moreover, it is reasonable that the next duration is likely to be smaller when past event arrivals are large. It follows that  $X_1(t) b_{0,1}$  is increasing in  $(N(s))_{s < t}$  so that Condition 1 is satisfied. This argument can be extended for  $K > 1$  to other variables like spread and order book volume imbalances in exactly the same way.

In the extreme case where  $X_1(t)$  is independent of past event arrivals  $(N(s))_{s < t}$ , we have that  $X_1(t)$  is strictly exogenous. Then,  $X_1(t) b_{0,1}$  can have an impact on  $N(t)$  and future event arrivals through the intensity  $\lambda(t)$ , but is not affected by past event arrivals  $(N(s))_{s < t}$ . Given that we do not assume independence of  $(N(s))_{s < t}$ , it is clear why we regard Condition

1 to be a weak exogeneity condition.

**Condition 2.** The model is parametrized in a way that is simple and intuitive to analyze. In practice we will have some raw covariates, say  $Z := (Z(t))_{t \geq 0}$  with values in  $\mathbb{R}^K$  and map them into  $[0, 1]^K$ . This is always possible because the extended real line is isomorphic to the unit interval. Such map can change the interpretation of the covariates. A notable example of such transformation is  $F(Z_k(t))$  where  $F(x) = \Pr(Z_k(t) \leq x)$  which can be approximated by the empirical distribution when the variables are stationary and ergodic. If the covariates take values in a known compact interval inside the real line, we can just use a linear transformation.

The raw covariates  $Z$  may actually take values in  $\mathbb{R}^L$  where  $L < K$ . This is common in many applications where we map the data into a higher dimensional space in order to capture nonlinearities. Notable examples are reproducing kernel Hilbert spaces, one-hot encoding and Bernstein polynomials. In the application of this paper we focus on one-hot encoding because of their simplicity; a precise definition of one-hot encoding will be given in due course. Ideally, we would parametric the model in a parsimonious way. However, ex ante we may not have a good understanding of how the variables impact the intensity. Hence, the approach can be seen as either nonparametric or the initial stage in the analysis of a high frequency dataset.

Given that we can map the raw variables into a higher dimensional space, the non-negativity restriction is a parametrization assumption. In fact, for the one-hot encoding method used in the empirical section, we have that variables can have a negative impact as soon as the linear coefficients are decreasing (see Figure 2 in Section 5.4). Similar comments pertains to Bernstein polynomials (e.g. Sancetta, 2018, Section 3.6.6, Mucciante and Sancetta, 2022, Section 2.2). For the sake of clarity, we now consider two examples.

Consider the intensity  $\lambda(Z(t)) = a_0 + a_1 Z_1(t) - a_2 Z_2(t)$  that depends on the raw covariate  $Z(t) = [Z_1(t), Z_2(t)]'$  with values in  $[0, 1]^2$ , where  $a_i \geq 0$ ,  $i = 0, 1, 2$ ,  $a_0 - a_2 \geq 0$ . The parameters' restriction ensures that this intensity is always nonnegative. Then,  $\lambda(Z(t))$  can be written as (3) where  $X_1(t) = 1$ ,  $X_2(t) = Z_1(t)$ ,  $X_3(t) = 1 - Z_2(t)$ , and  $b_{0,1} = a_0 - a_2$ ,  $b_{0,2} = a_1$ ,  $b_{0,3} = a_2$  where  $b_0 \geq 0$ . Hence, in our framework we are able to control the direction of the impact by the linear transformation  $x \mapsto 1 - x$ . From a computational point of view, this is equivalent to changing the sign of the covariate and imposing an additional inequality constraint. This example also makes clear that a purely linear model for (3) is only possible if the raw covariates  $Z$  are bounded and we impose constraints in the estimation. On the other hand, the use of unbounded variables can lead to a negative intensity.

As a second example we consider the case where each covariate is mapped into a higher dimensional space by one-hot encoding. For ease of exposition, consider a model with only one explanatory variable  $Z = (Z(t))_{t \geq 0}$ , where  $Z$  is a predictable ergodic stochastic process such that  $Z(t)$  takes values in  $\mathcal{Z}$ , a subset of the real line, for all  $t \geq 0$ . Let  $\mathcal{P} = \bigcup_{k=1}^K \mathcal{P}_k$  be a partition of  $\mathcal{Z}$ . One-hot encoding maps the univariate raw covariate  $Z(t)$  into a  $K$



dimensional covariate  $X(t)$ , where the  $k^{\text{th}}$  element in  $X(t)$  is  $X^{(k)}(t) = 1_{\{Z(t) \in \mathcal{P}_k\}}$ . Then,  $X$  is said to be a one-hot encoding of  $Z$ . It follows that  $g_0(t) = \sum_{k=1}^K b_{0,k} 1_{\{Z(t) \in \mathcal{P}_k\}}$  so that the coefficient  $b_{0,k}$  is the value of  $g_0(t)$  when  $Z(t) \in \mathcal{P}_k$ . As a concrete illustration, consider trade durations, where time is measured in seconds. Again, for simplicity of notation, suppose that the model only uses trade durations as explanatory variables so that there is only one raw covariate. In this case  $\mathcal{Z} = [0, \infty)$  because durations are nonnegative. For a very liquid instrument, we could set  $\mathcal{P}_k := \left[ \frac{(k-1)}{10}, \frac{k}{10} \right)$ ,  $k = 1, 2, \dots, 10$ , and  $\mathcal{P}_{11} := [1, \infty)$ . This means that the univariate raw covariate  $Z(t)$  is mapped into an 11 dimensional covariate  $X(t)$ , so that  $K = 11$ . In this case if the last duration at time  $t$  is 50 milliseconds,  $X^{(1)}(t) = 1$  and  $X^{(k)}(t) = 0$  for  $k \neq 1$ . This is because 50 milliseconds is 0.05 of a second. Hence, the duration of 50 milliseconds falls within the first bin  $\mathcal{P}_1 = [0, 0.1)$ . On the other hand if the duration at time  $t$  is 2.5 seconds, we have  $X^{(11)}(t) = 1$  and  $X^{(k)}(t) = 0$  for  $k \neq 11$ . Hence, if the impact of the duration  $Z(t)$  is monotonically decreasing, we shall have that  $b_{0,k} \geq b_{0,k+1}$  as in fact is the case for the estimated coefficients in Figure 2. In our empirical application we use slightly different bins for the durations and other covariate and the details are reported in Section 5.2, in the Appendix.

The extension to more than one raw covariate is achieved by applying the one-hot encoding to each variable separately. Suppose that  $Z(t)$  is a  $K_Z$  dimensional vector where  $K_Z > 1$ . The  $j^{\text{th}}$  element in  $Z(t)$  is mapped into a  $K_j$  dimensional covariate via one-hot encoding. Then, the resulting number of covariates  $X$  is  $K = \sum_{j=1}^{K_Z} K_j$ .

**Condition 3.** The bound  $|b|_1 \leq B$  is used for technical reasons to ensure that we achieve consistency under additional regularity conditions (e.g. Sancetta, 2016, 2018). This is equivalent to estimation via Lasso because by duality a penalty on the  $l_1$  norm of the linear coefficients is equivalent to a constraint on the  $l_1$  norm of these coefficients. However, given the nonnegativity constraint, the actual value of  $B$  becomes less relevant. The nonnegativity constraint tends to produce sparse solutions and has some set identification properties under certain assumptions (Mucciante and Sancetta, 2022). Using (3), the  $L_2$  norm of  $\lambda_0$  in (1) can be written as  $b'Ab$  where  $A = \frac{1}{T} \int_0^T (h_0^2(t) X(t) X(t)') dt$  is a  $K \times K$  matrix. Then, if for any  $b \geq 0$  (elementwise)  $b'Ab \geq \nu |b|_1^2$  for some  $\nu > 0$  we have that a bound on  $b'Ab$  controls  $|b|_1^2$ . If this is the case, intuitively, the sign constraint allows us to control the  $l_1$  norm as in Lasso, but through the control of the square error loss. This condition holds in numerous circumstances (Mucciante and Sancetta, 2022, Section 2.3). To see this note that all the entries in  $A$  are nonnegative. If we suppose that they are strictly positive, the inequality  $b'Ab \geq \nu |b|_1^2$  is trivially satisfied for all  $b \geq 0$ . Mucciante and Sancetta (2022, Section 2.3) discusses how the inequality can hold when  $A$  is not necessarily strictly positive.

**Condition 4.** A sum of exponential kernels is a simple but flexible parametrization and was originally discussed in Hawkes (1971). We choose it for the sake of definiteness. This allows us to avoid abstract technical conditions. Any parametric kernel that decays fast enough and is smooth in the parameters can be used. In this case, consistency can be proved following the steps in the proofs. Finally, the restriction on the parameter space to a compact set ensures non-negativity of the intensity and consistent estimation.

### 2.3 Stationarity and Ergodicity of the Point Process

For consistent estimation of the process, we shall use stationarity and ergodicity. Such statistical properties are of interest in their own merit and do not need the full extent of the Regularity Conditions. Hence, we state the following.

**Condition 5** *There is a filtration  $\mathcal{F} = (\mathcal{F}_t)$  such that  $\lambda_0(t) = h_0(t)g_0(t)$  is an  $\mathcal{F}_t$ -intensity for  $N(t)$  where  $h_0$  is as in (2) and  $g_0$  is a uniformly bounded nonnegative stochastic process satisfying Condition 1 and such that  $|g_0|_\infty \sum_{l=1}^L \frac{d_{0,l}}{a_{0,l}} < 1$ .*

When  $g_0$  is constant, as in the case of a standard Hawkes process, Condition 5 is the usual one for stationarity and ergodicity of the Hawkes process (Brémaud and Massoulié, 1996). The uniform upper bound on  $g_0$  ensures that we can still apply those results together with the help of Condition 1. When  $g_0$  is not constant, Condition 5 restricts  $g_0$  to be weakly exogenous in the sense of Condition 1.

Under the above condition, the point process is strictly stationary, ergodic with strongly mixing coefficients decaying exponentially fast. Recall that a point process is stationary if  $N(t+C)$  and  $N(C)$  have same distribution for any  $t \in \mathbb{R}$  and  $C \in \mathbb{R}$ . We say that the stationary point process  $N$  is ergodic if  $T^{-1}N(T) \rightarrow \mathbb{E}\lambda_0(0)$  in probability;  $\lambda_0(0)$  is  $\lambda_0(t)$  for  $t = 0$ . A stationary point process  $N$  is strongly mixing if

$$\alpha(t) = \sup \{ |\Pr(A \cap B) - \Pr(A)\Pr(B)| : A \in \mathcal{E}_{-\infty}^0, B \in \mathcal{E}_t^\infty \}$$

goes to zero as  $t \rightarrow \infty$ , where  $\mathcal{E}_r^t$  is the sigma algebra generated by the cylinder sets on the interval  $(r, t]$  (e.g. Cheysson and Lang, 2022). We have the following.

**Theorem 1** *Suppose that Condition 1 hold.*

1. *Then, there is a unique stationary distribution of  $N$  with finite average intensity  $\mathbb{E} \int_0^1 dN(t)$  and dynamics as in (1) and the process is ergodic with exponentially decaying strongly mixing coefficients.*
2. *Suppose that for  $t \leq 0$ ,  $N$  is restricted to the set  $\mathcal{A} := \{N(t) = 0 : t \leq 0\}$  so that  $h_0(t) := c_0 + \int_{(0,t)} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}(t-s)} dN(s)$ . Then, there is a stationary point process*

$\tilde{N}$  with same dynamics as in (1), and a stochastic time  $\tau$  such that  $\Pr(\tau < \infty) = 1$  and  $\tilde{N}(t) = N(t)$  for all  $t \geq \tau$ .

The processes  $\tilde{N}$  and  $N$  couple in finite time irrespective of the initial condition. This is important because in practice we only have data for  $t > 0$ . This is equivalent to restricting  $N$  to  $\mathcal{A}$ .

## 2.4 Reduced Form Model for Buy and Sell Events

Multivariate extensions of intensity models to buy and sell events have been considered in the literature, but to keep focus we limit the discussion to univariate intensity (Bauwens and Hautsch, 2009, and the references in Section 1.1 for such extensions). Nevertheless, we show how we can separately estimate a reduced form model for buy and sell event arrivals. Let  $g^{buy}$  and  $g^{sell}$  be as in (3) but for buy and sell events separately. Both can depend on order book variables as well as past durations and other quantities. Suppose the following feedback loop effect through  $g^{buy}$  and  $g^{sell}$ ,

$$\begin{aligned}\lambda^{buy}(t) &= h^{buy}(t) \left( g^{buy}(t) + \rho^{buy} g^{sell}(t) \right) \\ \lambda^{sell}(t) &= h^{sell}(t) \left( g^{sell}(t) + \rho^{sell} g^{buy}(t) \right),\end{aligned}\tag{4}$$

where  $\rho^{buy}$  and  $\rho^{sell}$  are constants in  $[0, 1)$ . Because of (3), this system has the reduced form

$$\begin{aligned}\lambda^{buy}(t) &= h^{buy}(t) X(t)' \left( b^{buy} + \rho^{buy} b^{sell} \right) \left( 1 - \rho^{buy} \rho^{sell} \right)^{-1} \\ \lambda^{sell}(t) &= h^{sell}(t) X(t)' \left( b^{sell} + \rho^{sell} b^{buy} \right) \left( 1 - \rho^{buy} \rho^{sell} \right)^{-1}.\end{aligned}$$

In consequence, separate estimation of the buy and sell intensities is equivalent to estimation of the above reduced form model as opposed to (4). For large samples, the loss in efficiency for carrying out a separate estimation is secondary.

A natural extension of the multivariate Hawkes process considered in the literature (Hawkes, 1971) would be  $\lambda^{buy}(t) = h^{buy}(t) g^{buy}(t) + \rho^{buy} h^{sell}(t) g^{sell}(t)$  and similarly for the sell intensity. This differs from (4). In (4) we are assuming that the intensity for buy events is independent of  $h^{sell}(t)$  when we condition on  $g^{sell}(t)$  and  $g^{buy}(t)$ , and similarly for sell events.

## 3 Estimation

Throughout, to keep the notation more compact we may write  $\int_0^T f d\mu = \int_0^T f(t) dt$  for any measurable function, where  $\mu$  is the Lebesgue measure. The loglikelihood for  $\lambda = hg$  as in (1),

is

$$L_T(h, g) := \int_0^T \ln(hg) dN - \int_0^T hgd\mu. \quad (5)$$

Estimation of  $g = X'b$  requires a positivity constraint on  $b$ . This, coupled with the high dimension  $K$  of  $b$ , makes the problem unfeasible. Moreover, when the sample size is large, it is not possible to hold the data in memory. To see this, suppose that  $X(t) = X(t_j)$  for  $t \in (t_j, t_{j+1}]$  where the times  $t_j$  are update times for any of the covariates including the counting process,  $j = 1, 2, \dots, m$ . Throughout,  $m$  is the total number of event updates, including the jumps of the process  $N$ . The second term in the loglikelihood (5) is explicitly written as  $\sum_{j=1}^m \left( \int_{t_j}^{t_{j+1}} h(t) dt \right) X(t_j)' b$ . This requires to hold in memory a matrix  $m \times K$ . It is only feasible for small scale high frequency applications. Even for moderate scale problems, when  $K$  is in the order of hundreds or thousands, and  $m$  is a few million (e.g. a few days for some liquid instruments), estimation of  $g$  by loglikelihood with a positivity constraint is impractical. We may have only a few thousand events generated by  $N$  in a day, but the order book updates much faster and we wish to track the impact of order book derived covariates.

It is simple to see that around the true value the negative loglikelihood is quadratic. In fact, an alternative estimation strategy for intensity processes is based on least squares (Gaïffas and Guillaou, 2012, Mucciante and Sancetta, 2022, and references therein). As a first step, we suggest to minimize the quadratic loss

$$Q_T(h, g) := -\frac{2}{T} \int_0^T hgdN + \frac{1}{T} \int_0^T (hg)^2 d\mu. \quad (6)$$

When,  $h$  is known, minimization is just a quadratic programming problem, hence much easier to solve than (5). The rationale for minimizing (6) is that its expectation is equal to

$$\mathbb{E}Q_T(h, g) := -\frac{2}{T} \mathbb{E} \int_0^T hg\lambda_0 d\mu + \frac{1}{T} \mathbb{E} \int_0^T (hg)^2 d\mu$$

and it is minimized by  $h = h_0$  and  $g = g_0$  because  $\lambda_0 = h_0g_0$ . When  $h$  is unknown, the second term in (6) poses similar computational challenges as the second term in (5) for large datasets. To solve this, we add a second step in the optimization procedure. We suggest to fix  $g$  and estimate  $h$  and vice versa, as commonly done in coordinate descent algorithms. When it comes to optimization w.r.t.  $g$ , given an  $h$ , we use the loss function

$$R_T(g; h) := -\frac{2}{T} \int_0^T \frac{g}{h} dN + \frac{1}{T} \int_0^T g^2 d\mu \quad (7)$$

in place of (6). This is a much simpler problem, as there is a summary statistic for the second term as in standard regression problems. In summary, we alternate between estimation of  $h$  minimizing (6) with  $g$  fixed and estimation of  $g$  minimizing (7) with  $h$  fixed. The starting

value for  $g$  is set to a constant so that at the first iteration we estimate a standard Hawkes process using (6). Algorithm 1 summarizes the procedure. There, the only free parameter is  $B$ , as defined in Condition 2. We discuss its choice in Section 3.1.

The theory of coordinate descent algorithms justifies alternating between minimization w.r.t.  $g$  and  $h$  (Beck and Tetruashvili, 2013). We shall show that asymptotically, the minimizers of (6) and (7) w.r.t.  $g \in \mathcal{G}$  are the same if  $h$  is close to  $h_0$ .

Estimation via the loss function (7) is very efficient as it is a standard quadratic programming problem. Let  $T_j$  be the  $j^{\text{th}}$  jump time of  $N$ ,  $j = 1, 2, \dots, n$  with  $n = N(T)$  and for simplicity suppose that  $t_m = T$ . For numerical estimation we rewrite (7) explicitly as

$$R_T(b) := -\frac{2}{T}b'\Phi'\Gamma + \frac{1}{T}b'\Phi'\Sigma\Phi b \quad (8)$$

where the  $j^{\text{th}}$  row of  $\Phi$  is  $X(t_j)$ ,  $\Gamma$  is a vector with  $j^{\text{th}}$  entry  $1/h(T_l)$  if  $t_j = T_l$  for some  $l$  (i.e., if  $t_j$  is a jump time of  $N$ ) and zero otherwise;  $\Sigma$  is a diagonal matrix with  $(j, j)^{\text{th}}$  entry  $(t_j - t_{j-1})$ . Note that when a new observation is collected, we only need to update  $\Phi'\Gamma$  and  $\Phi'\Sigma\Phi$ , which are low dimensional matrices ( $K \times 1$  and  $K \times K$ , respectively) relatively to  $m$ . Moreover, the dimension of  $\Phi'\Sigma\Phi$  does not depend on  $m$  (very large) and on  $h$ .

We may incur additional difficulties when the  $K$ -dimensional process  $X$  has linearly dependent entries. This is the case for our application based on one-hot encoding as well as for other methods such as estimation in reproducing kernel Hilbert spaces or additive multivariate Bernstein polynomials. However, the nonnegativity constraints mitigates this problem if we can find a  $\nu > 0$  such that  $b'(\Phi'\Sigma\Phi/T)b \geq \nu b'b$  for all  $b \geq 0$ , where the inequality on  $b$  is meant elementwise (see the discussion to Condition 3 in Section 2.2). In the applications we have considered the problem remained convex and a fast solution could be obtained. For more difficult problems, we may rely on greedy algorithms (e.g. Sancetta, 2016, 2018, and references therein).

When the number of updates is very large, computation of the second integral in (6) can be slow even for fixed  $g$ . Monte Carlo techniques can be used in this case (Cartea et al., 2021).

---

**Algorithm 1** Intensity Estimation

---

Start with  $g^{(0)} = \gamma$ .

For each  $v = 1, 2, \dots$ , find the minimizer of  $Q_T(h, g^{(v-1)})$  (as in (6)) w.r.t.  $h$  and denote it by  $h^{(v)}$ . When  $v = 1$  we shall also minimize w.r.t.  $\gamma > 0$ .

Minimize  $R_T(X'b; h)$  (as in (7)) w.r.t.  $b \in [0, \infty)^K$  s.t.  $\sum_k b_k \leq B$ .

Define the minimizer by  $b^{(v)}$  so that  $g^{(v)} = X'b^{(v)}$  is the estimator for  $g_0$ .

Stop when  $h^{(v)}g^{(v)}$  converges.

---

**Reparametrization for improved estimation.** We found improvements in the estimation of (1) using the reparametrization

$$\lambda_0(t) = \left( c_0 d_0 + \int_{(-\infty, t)} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}(t-s)} dN(s) \right) \frac{g_0(t)}{\mathbb{E}g_0}. \quad (9)$$

In (6), we would estimate  $c$  and  $d$  using  $g/\mathbb{E}g$  as given, instead of  $g$ . We would then estimate  $g$  in (7) using  $h(t) = \left( c + \int_{(-\infty, t)} \sum_{l=1}^L \frac{d_l}{d_1} e^{-a_{0,l}(t-s)} dN(s) \right)$  as given; for identification, this forces  $d_l/d_1 = 1$  when  $l = 1$ . This means that at every iteration  $\mathbb{E}g$  is replaced with the most recent estimator for  $\mathbb{E}g$ , for example, in the empirical section, at the  $v^{\text{th}}$  iteration we use  $\mathbb{E}g \simeq \frac{1}{m} \sum_{j=1}^m X(t_j)' b^{(v-1)}$ ; the update times  $t_j$  were defined below (5).

### 3.1 Choice of $B$

The non-negativity constraint leads to the shrinkage property studied in Mucciante and Sancetta (2022), among others. Therefore, the choice of  $B$  (Condition 3) is not very important.

Nevertheless, following Sancetta (2018) the value  $B$  can be chosen to minimize  $-L_T(\hat{h}_B, \hat{g}_B) + K_B$  w.r.t.  $B$ , where  $L_T$  is the loglikelihood in (5) and  $\hat{h}_B, \hat{g}_B$  are the estimators from Algorithm 1 for a given  $B$ . Here,  $K_B$  is the number of nonzero coefficient in the estimator for  $b_0$ . We note that even when  $B = \infty$ , the number  $K_B$  is smaller than  $K$ , which is the size of  $b_0$ . This is due to the aforementioned shrinkage property.

To identify a reasonable order of magnitude for  $B$ , we can use the fact that  $\mathbb{E}h_0 g_0 \leq \mathbb{E}h_0 |g_0|_\infty = \left[ c_0 / \left( 1 - \sum_{l=1}^L \frac{d_{0,l}}{a_{0,l}} \right) \right] B$ . Moreover,  $\mathbb{E}h_0 g_0$  is approximately  $N(T)/T$  so that we would expect  $B$  to be greater than  $\left[ c_0 / \left( 1 - \sum_{l=1}^L \frac{d_{0,l}}{a_{0,l}} \right) \right]^{-1} N(T)/T$ . In practice, we replace  $c_0$ ,  $d_{0,l}$  and  $a_{0,l}$  with the estimated parameters from a Hawkes process with  $g_0$  constant. We can also approximate  $B$  to be of the same order of magnitude as  $d_{0,1}$  when we use the parametrization in (9). Assuming that on average the covariates have expectation 1/2, then we would have  $B = 2d_{0,1}$ . The simulation results show that such large value is reasonable and results are not sensitive to it.

A related approach is to directly impose a constraint on the magnitude of each coefficient. Assume for the moment that  $b = \beta \mathbf{1}_K$  for some constant  $\beta > 0$ , and solve (8) under this constraint. The solution is  $\beta = \mathbf{1}'_K \Phi' \Gamma / (\mathbf{1}'_K \Phi' \Sigma \Phi \mathbf{1}_K)$ . We can then solve a quadratic programming problem under the constraint that  $b_k \in [0, \beta]$ , which clearly implies  $B = K\beta$ . We chose such approach in the empirical section and found our results not to be sensitive to the above alternatives.

### 3.2 Simulations: Number of Iterations for the Algorithm to Converge

We found that the algorithm would converge within two or three iterations. We showcase this with a simulation. We let  $X(t) = X(T_{i-1})$  when  $t \in (T_{i-1}, T_i]$  and let  $X(T_{i-1})$  be uniformly distributed in  $[0, 1]^K$  for  $i \geq 0$ . Moreover, we set  $b_0$  such that the first three entries are equal to  $2/3$  and the remaining are equal to zero. This choice implies that  $\mathbb{E}g = 1$ . Moreover, we consider  $L = 1$  and  $(c_0, d_0, a_0) = (1, 1, 2)$ , dropping the subscript  $l$ . We use Algorithm 1 for the estimation. We use the following parameters restrictions:  $B = K$  and  $a \in [10^{-9}, 10^4]$ ,  $c \in [10^{-9}, 10]$ ,  $d \in [10^{-9}, 10^3]$ . We cannot report all estimated values of  $b$ . To help assess how close estimated parameters are to the true ones, we compute

$$\text{Error}(\alpha) := \sum_{k=1}^K 1_{\{|b_k - b_{0,k}| > \alpha(2/3)\}} \quad (10)$$

where  $2/3$  is the value of the non-zero entries in  $b_0$  and  $\alpha$  is a small number. We compute the above quantities from a single realization of 100,000 jumps from the process. As we increase the number of iterations, the convergence is remarkably fast. To reduce the dependence on a the single sample realization, we also compute an average, at the fourth iteration, of the end results and their standard errors over 50 realization. The average numbers over the simulations confirm that the claim is not just the result of chance. Table 1 reports the details. From those results we deduce that the error in the estimation is relatively low (see columns  $\alpha = 0.1, 0.05$  in Table 1). We can deduce that two/three iterations are enough to converge. Moreover, the results in the row marked as ‘‘avg’’ of Table 1 show that the procedure has good finite sample properties. In particular, for this simulation design, we verify the variable screening properties of the methodology. Mutatis mutandis, simulations in Mucciante and Sancetta (2022) shed additional light in this respect, in a richer variety of simulation designs when  $h_0 = 1$ .

The results show that the parameters converge fast with the number of iterations. The convergence is towards the true parameters. We also report sensitivity results around different choices of  $B$ . In particular, we consider

$$B = \text{Mult.} \times [\hat{c} / (1 - \hat{a}^{-1})]^{-1} N(T) / T \quad (11)$$

where  $\hat{a}$  and  $\hat{c}$  are the estimators for  $a_0$  and  $c_0$  at the first iteration in Algorithm 1 and  $\text{Mult.} \in \{0.5, 1, 10\}$ . Note that when  $\text{Mult.} = 0.5$ , the resulting  $B$  is smaller than the recommendation from Section 3.1 and we expect a sub-optimal performance. We also consider  $B = K$  which satisfies the requirement that  $B > \mathbb{E}g_0$  as discussed in Section 3.1. Under various criteria, we find out that the choice of  $B$  is not critical as long as the recommendations from Section 3.1 are followed. The results in Table 2 corroborate this statement.

**Table 1:** Example of Algorithm Performance. Estimation is for  $b \in [0, B]^K$  with  $B = K$ . The true values are  $b_{0,k} = 2/3$  for  $k = 1, 2, 3$  and zero otherwise, and  $(c_0, d_0, a_0) = (1, 1, 2)$ . Their estimated values are reported in columns  $b_1, b_2, b_3, c, d, a$  for one single simulation over increasing number of iterations (iter) in Algorithm 1. The average (avg.) and standard error (s.e.) over 50 simulations are given in the last two rows for iter 4. Columns Error ( $\alpha$ ) compute the statistic (10) for different values of  $\alpha$ . Small values imply higher precision of the estimator.

$K$	iter	$b_1$	$b_2$	$b_3$	$c$	$d$	$a$	Error ( $\alpha$ ) in (10)		
								$\alpha$	0.1	0.05
3	1	0.675	0.640	0.684	0.701	1.130	2.171	0	0	3
	2	0.677	0.637	0.684	1.014	0.988	2.008	0	0	3
	3	0.677	0.637	0.684	1.014	0.988	2.008	0	0	3
	4	0.677	0.637	0.684	1.014	0.988	2.008	0	0	3
	avg	0.668	0.667	0.665	0.998	0.999	2.000	0	0	1.260
	s.e.	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.112
10	1	0.648	0.614	0.658	0.701	1.130	2.171	0	1	8
	2	0.677	0.637	0.684	1.007	0.997	2.013	0	0	3
	3	0.677	0.637	0.684	1.014	0.988	2.008	0	0	3
	4	0.677	0.637	0.684	1.014	0.988	2.008	0	0	3
	avg	0.664	0.663	0.661	0.998	1.000	2.001	0.000	0.000	2.300
	s.e.	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.188
100	1	0.635	0.600	0.645	0.701	1.130	2.171	0	1	10
	2	0.674	0.634	0.681	1.005	1.002	2.019	0	0	3
	3	0.675	0.635	0.682	1.014	0.989	2.009	0	0	3
	4	0.675	0.635	0.682	1.014	0.989	2.009	0	0	3
	avg	0.659	0.658	0.655	0.997	1.002	2.003	0.000	0.020	3.420
	s.e.	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.003	0.351



**Table 2:** Sensitivity of Results to  $B$ . Results for different metrics are reported as we vary the value of  $B$  using (11). When a value for Mult. is not reported, it means that  $B = K$ . The average (avg.) and standard error (s.e.) over 50 simulations are reported. The  $l_1$  and  $l_2$  norms of the difference between the estimator and the true one are scaled by the norm of the true parameter. The columns FP, FN and P report the false positive, false negative and the positives: the number of non-zero estimated coefficients that should in fact be zero, the number of zero estimated coefficient that instead should be non zero, and the number of non-zero coefficients in the population. The estimated value of  $B$  is also reported.

$K$	Mult.		$l_1$	$l_2$	FP	FN	P	estimated $B$
3	0.5	avg	0.0281	0.0334	0	0	3	0.6496
		s.e.	0.0004	0.0050	0	0	0	0.0003
	1	avg	0.0150	0.0191	0	0	3	1.2991
		s.e.	0.0001	0.0026	0	0	0	0.0006
	10	avg	0.0150	0.0191	0	0	3	12.9915
		s.e.	0.0001	0.0026	0	0	0	0.0057
-	avg	0.0150	0.0191	0	0	3	3	
	s.e.	0.0001	0.0026	0	0	0	0	
10	0.5	avg	0.0540	0.0412	7	0	3	0.6496
		s.e.	0.0005	0.0054	0	0	0	0.0003
	1	avg	0.0215	0.0221	7	0	3	1.2991
		s.e.	0.0002	0.0029	0	0	0	0.0006
	10	avg	0.0215	0.0221	7	0	3	12.9915
		s.e.	0.0002	0.0029	0	0	0	0.0057
-	avg	0.0215	0.0221	7	0	3	10	
	s.e.	0.0002	0.0029	0	0	0	0	
100	0.5	avg	0.0585	0.0385	97	0	3	0.6496
		s.e.	0.0005	0.0051	0	0	0	0.0003
	1	avg	0.0324	0.0265	97	0	3	1.2991
		s.e.	0.0003	0.0033	0	0	0	0.0006
	10	avg	0.0324	0.0265	97	0	3	12.9915
		s.e.	0.0003	0.0033	0	0	0	0.0057
-	avg	0.0324	0.0265	97	0	3	100	
	s.e.	0.0003	0.0033	0	0	0	0	

## 4 Asymptotic Results

First we show that estimation of the model using the quadratic loss (6) gives consistent estimators even in the ultra high dimensional case.

**Theorem 2** Define  $(\tilde{h}, \tilde{g}) := \arg \inf_{h,g} Q_T(h, g)$ . Under the Regularity Conditions,

$$\frac{1}{T} \int_0^T (h_0 g_0 - \tilde{h} \tilde{g})^2 d\mu = O_P \left( B \sqrt{\frac{\ln(1+K)}{T}} \right).$$

Define  $\hat{g}^h = \arg \inf_{g \in \mathcal{G}} R_T(g; h)$  and  $g^h = \arg \inf_{g \in \mathcal{G}} \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} - g \right)^2 d\mu$ . We show that the minimizer  $\hat{g}^h$  converges to the best  $L_2$  approximation  $g^h$  uniformly in  $h$  under certain conditions. We can then conclude that in a neighbourhood of  $h_0$  the estimators of  $g_0$  using the loss functions (6) and (7) are equivalent. We make this clear with the following.

**Theorem 3** Under the Regularity Conditions,

1.  $\sup_{g \in \mathcal{G}} \left| R_T(g; h) - \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} - g \right)^2 d\mu \right| = \frac{1}{T} \int_0^T (\lambda_0/h)^2 d\mu + O_P \left( B \sqrt{\frac{\ln(1+K)}{T}} \right)$  uniformly in  $h \in \mathcal{H}$ ; moreover we have that  $\sup_{h \in \mathcal{H}} \frac{1}{T} \int_0^T (\lambda_0/h)^2 d\mu \leq \frac{1}{T} \int_0^T \underline{c}^{-2} \lambda_0^2 d\mu$  where  $\underline{c} = \min \{c \in \mathcal{C}\}$ ;
2. for any set  $\mathcal{H}' \subseteq \mathcal{H}$ , we have that  $\sup_{h \in \mathcal{H}'} \frac{1}{T} \int_0^T (\hat{g}^h - g^h)^2 = O_P \left( \max \left\{ B \sqrt{\frac{\ln(1+K)}{T}}, A_T \right\} \right)$ , where  $A_T = \sup_{h \in \mathcal{H}'} \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} - g^h \right)^2 d\mu$ ;
3.  $\frac{1}{T} \int_0^T (\hat{g}^{\hat{h}} - g_0)^2 = O_P \left( B \sqrt{\frac{\ln(1+K)}{T}} \right)$ , if  $\hat{h} \in \mathcal{H}$  is such that  $\frac{1}{T} \int_0^T (\hat{h} g_0 - h_0 g_0)^2 = O_P \left( B \sqrt{\frac{\ln(1+K)}{T}} \right)$ .

In Point 1, the second term on the right hand side (r.h.s.) does not depend on  $g$  and  $h$ , hence the objective function converges uniformly to the average square loss plus some quantity independent of  $g$ . Note also that  $\int_0^T \left( \frac{\lambda_0}{h} - g \right)^2 d\mu = \int_0^T \left( \frac{\lambda_0 - hg}{h} \right)^2 d\mu$ . If the approximation of the scaled intensity  $\lambda_0/h$  by  $g^h$  is good in the sense that the approximation error  $A_T$  in Point 2 is small, then the estimator  $\hat{g}^h$  converges to the best approximation uniformly in  $h \in \mathcal{H}'$ . For example,  $\mathcal{H}'$  could be a ball of small radius around  $h_0$ . This is made clearer in Point 3. If we can find  $\hat{h} \rightarrow h_0$  in  $L_2$  in probability, then also  $\hat{g}^{\hat{h}} \rightarrow g_0$  in  $L_2$  in probability. The contribution to the error that comes from the dimensionality  $K$  of the covariates is only logarithmic. Hence, the methodology is suited for potentially ultra high dimensional models.

To put the rates of convergence into perspective, note that in the regression setting with i.i.d. Gaussian errors, no estimator of an arbitrary convex combinations of variables can

achieve a rate faster than  $(\frac{\ln K}{n})^{1/4}$ , where here  $n$  is the sample size (Tsybakov, 2003, Theorem 2). In our case  $n$  can be viewed as the number of event arrivals  $N(T)$ . Hence, the result from Theorems 2 and 3 can be considered optimal unless we add additional conditions.

#### 4.1 Test Statistic to Compare Two Intensity Estimators

A common diagnostic test for point processes is to use the fact that  $\int_{T_{i-1}}^{T_i} \lambda_0(t) dt$ ,  $i = 1, 2, 3, \dots$  are i.i.d. exponential random variables (Bauwens and Hautsch, 2009, Section 4.3); recall that the times  $T_i$  are the times of the event arrival. Such method cannot be used to compare the relative fit of two competing intensities. For this reason, methods based on the loglikelihood ratio are better suited. Due to the fact that the data is possibly high dimensional, the loglikelihood ratio can have nonstandard distribution in high dimensions. To avoid this problems we consider a sample splitting procedure. Then, the result also applies to the intensity that is recursively estimated. All that we need is a measurable intensity that is bounded away from zero. Using a sample up to time 0, the estimator for two competing intensities  $\lambda^{(k)}$ ,  $k = 1, 2$ , is denoted by  $\hat{\lambda}^{(k)}$ . This estimator is evaluated on data on the following days. For model  $k$ , the log likelihood (see (5)) evaluated on a sample  $(0, T]$  using  $\hat{\lambda}^{(k)}$  is denoted by  $L_T^{(k)}$ . Let  $q_\alpha$  be the  $\alpha$  quantile of the standard normal distribution, e.g.,  $q_{95} \simeq 1.64$ . At the  $(1 - \frac{\alpha}{100})\%$  significance level for a one sided test, we reject model 2 in favor of model 1 if

$$\frac{L_T^{(1)} - L_T^{(2)}}{\sqrt{T\hat{\sigma}_T^2}} \geq q_\alpha \quad (12)$$

where

$$\hat{\sigma}_T^2 = \frac{1}{T} \int_0^T [\ln(\hat{\lambda}^{(1)}/\hat{\lambda}^{(2)})]^2 dN. \quad (13)$$

The predictable part of the log-likelihood  $L_T^{(k)}$  is

$$H_T^{(k)} := \int_0^T \ln(\hat{\lambda}^{(k)}) \lambda_0 d\mu - \int_0^T \hat{\lambda}^{(k)} d\mu.$$

The predictable part of the likelihood ratio is such that  $H_T^{(k)}$  is maximized when  $\hat{\lambda}^{(j)} = \lambda_0$ . The closer  $\hat{\lambda}^{(1)}$  is to  $\lambda_0$  relatively to  $\hat{\lambda}^{(2)}$  the larger is  $H_T^{(1)} - H_T^{(2)}$ . Closeness of  $\hat{\lambda}^{(k)}$  to  $\lambda_0$  in measured in the sense of Kulback-Liebler divergence, i.e. how  $H_T^{(0)} - H_T^{(k)} \geq 0$  is close to zero, where  $H_T^{(0)}$  is as in the above display replacing  $\hat{\lambda}^{(j)}$  with  $\lambda_0$ . Define the predictable part of the log-likelihood ratio as

$$\epsilon_T := H_T^{(1)} - H_T^{(2)}. \quad (14)$$

Under the null hypothesis, we suppose that  $\epsilon_T = o_p(\sqrt{T})$ . Loosely speaking, the intensities  $\hat{\lambda}^{(1)}$  and  $\hat{\lambda}^{(2)}$  give similar predictions asymptotically, i.e. the predictable part of the

log-likelihood ratio diverges at a rate slower than  $o_p(\sqrt{T})$ . This means that we can only distinguish intensities that are sufficiently far apart as  $T \rightarrow \infty$ . The rate of divergence of  $\epsilon_T$  is clearly of order  $T$  when the two intensities differ by some positive function that is bounded away from zero. The following can be used to justify (12). The next result generalizes Proposition 1 in Sancetta (2018) to unbounded intensities.

**Theorem 4** *Suppose that the Regularity Conditions hold, that  $\hat{\lambda}^{(k)} \in \mathcal{H} \times \mathcal{G}$ ,  $k = 1, 2$ ,  $\hat{\lambda}^{(1)} \neq \hat{\lambda}^{(2)}$  and such that no intensity is exactly zero. If  $\epsilon_T = o_p(\sqrt{T})$ , then,*

$$\frac{L_T^{(1)} - L_T^{(2)}}{\sqrt{T\hat{\sigma}_T^2}} \rightarrow Z \quad (15)$$

*in distribution where  $Z$  is a standard normal random variable.*

## 5 Empirical Application

We separately model trade arrivals on each side of the order book. It is also of interest to consider trade arrivals for large trade sizes. This is because for large trade sizes a market maker might face adverse selection. This means being filled on a passive/resting order when the market goes in the opposite direction of the filled order. For example, a market maker may have an order to sell at the top of book ask. A large buy trade is initiated and it depletes the whole first level on the asks. In consequence, the market maker sells at the ask, but the price jumps up, making the sale unprofitable, at least in the short term. In what follows we consider two type of events: 1. trades arrivals on one side of the book, irrespective of size, and 2. trades arrivals on one side of the book for sizes that are at least as large as the size posted on the top of book. In what follows, we shall refer to these as Any Trade Arrivals and Large Trade Arrivals.

In particular, we study four stocks traded on the NYSE: Amazon (AMZN), Cisco (CSCO), Disney (DIS) and Coca Cola (KO). We also use the ETF on the S&P500 (SPY) as auxiliary instrument. The stocks tickers are given inside the parenthesis. The sample period is 01/March/2019-30/April/2019 from 9:30am until 4:30pm every trading day.

**Objectives.** We use our flexible methodology to answer a number of questions.

1. Does the order book provide information in addition to what is captured by trades?
2. Is self excitation important once we account for the order book?
3. Is the impact of such information nonlinear?

4. Is the simple exponential kernel ( $L = 1$  in (2)) enough to capture the self excitation nature of event arrivals?

To do so, we test various functional restrictions within our framework. We test the hypotheses that  $h_0$  is constant against the unrestricted model, and that  $g_0$  is constant against the unrestricted model. To account for possible nonlinearities, we use one-hot-encoding to obtain an estimator of  $g_0$ . To verify whether the nonlinearity is important, we also allow  $g_0$  to be linear in the underlying raw variables, i.e. without mapping them into a higher dimensional space using one-hot-encoding. In this case, we do not restrict the coefficients to be nonnegative. The out of sample test procedure for all these specifications is discussed in Section 4.1.

**Data quality.** The data was obtained from the Lobster database querying the first ten levels of the order book. The Lobster dataset is based on Level 3 data. It means that for each given stock, every order is included, e.g. insertion, cancellation, execution of visible and hidden limit orders. Lobster construct a snapshot of the order book for any such orders. Details and nuisances about the dataset can be found in Section A.3.1 of the Appendix.

**Data size.** We conflate multiple updates into a single one if the exchange time stamp is the same. Then, for AMZN, we have more than 17 million book updates with unique time stamp for ten levels and in excess of 400 thousand trades over the whole sample of 42 available trading days (Table 4). We use the first three levels, and these usually account for about 60% of the updates. The number of parameters  $K$  to be estimated is at most 177. This means about 300 million data points. Our estimation procedure had no problems to load the data in RAM, and was rather fast (less than two hours to parse data, a day at the time, and estimate a model).

## 5.1 The Raw Covariates

We estimate (1) where  $X$  is a one-hot encoding of the raw covariates reported in Table 3. Some of these raw covariates are obtained applying exponential moving average (EWMA) filters to the data. This is the case if a smoothing parameter is specified in Table 3. The EWMA of a variable  $Z(t_i)$  with smoothing parameter  $\alpha$  is

$$EWMA(Z(t_i)) = \alpha EWMA(Z(t_{i-1})) + (1 - \alpha) Z(t_i) \quad (16)$$

where  $EWMA(Z(t_1)) = Z(t_1)$ ,  $i = 2, 3, 4, \dots$ . Here,  $t_1$  is the time of the first update in the variable  $Z$  at the start of each day. EWMA's are computed for each day. With abuse of notation we then set the raw covariate equal to  $EWMA(Z(t_i))$ . The raw covariates are similar to the ones in Mucciante and Sancetta (2022), though with some differences. Hence,

we briefly summarize their construction. The book volume imbalance at level  $j$  is defined as

$$\text{VolImb}_j = \frac{\text{BidSize}_j - \text{AskSize}_j}{\text{BidSize}_j + \text{AskSize}_j} \quad (17)$$

where  $\text{BidSize}_j$  is the bid size (quantity) at level  $j$ , and similarly for  $\text{AskSize}_j$ . This variable takes values in  $[-1, 1]$ . The trade imbalance is computed from the EWMA of the signed traded volume every time there is a trade. We then divide it by the EWMA of the unsigned volumes. The EWMA's parameter is  $\alpha = 0.98$  for both denominator and numerator. Durations are in seconds with nanosecond decimals. They are then passed to EWMA filters with parameters  $\alpha = 0.98$  and  $0.90$ . The spread is computed in basis points. After the application of EWMA's filters, our model (1) has 7 raw covariates per instrument plus a seasonal component. The seasonal is time of the day for each update between  $[09 : 30, 16 : 30]$  EST, standardized to be in  $[0, 1]$ . We also include information from two additional auxiliary instruments with no seasonal component. Then, the total number of raw covariates is 22. We then apply one hot encoding as discussed in Section 5.2. The total number of parameters to estimate becomes at most 168 once we add a constant. Adding a constant introduces perfect linear dependence. As discussed at the end of Section 3 this is not a problem from a computational point of view.

We define a reference time to be the time at which a book or trade update is sent for a traded instrument. A traded instrument is the one whose intensity we are modelling. Then, to limit the size of the data, after computing all the covariates, we sample them at the reference time only. Finally, to ensure that the covariates are predictable, we make them left continuous by lagging them as a very last step in the procedure. Failing to do so would lead to forward looking bias. To see this, let  $t_{i-1}$  and  $t_i$  be the time of the  $i-1$  and  $i$  order book update. The intensity at time  $t_i$  can only use the book update from time  $t_{i-1}$ . It will use the information from the  $i^{\text{th}}$  book update immediately after  $t_i$ . This implies that if we were to observe a trade at time  $t_i$ , the  $i^{\text{th}}$  book update could not be used to predict the trade. In summary, we ensure that the same conditions for live trading are reproduced in our estimation.

**Table 3:** Raw Covariates Used for Estimation. The column “Smoothing” reports the smoothing parameter used if an EWMA had been applied to the original variable.

Variables	Short Name	Smoothing
Seasonal	Seas	
Volume Imbalance Level 1	VolImb1	
Volume Imbalance Level 2	VolImb2	
Volume Imbalance Level 3	VolImb3	
Spread	Spread	
Trade Imbalance	TrdImb98	$\alpha = 0.98$
Durations	Dur98, Dur90	$\alpha = 0.98, 0.90$

## 5.2 One-Hot Encoding of Covariates

To automate the procedure, we opted for a simple rule that can be applied to all the covariates. Let  $q_x$  be the  $x\%$  quantile of a covariate based on the estimation sample. For all covariates, we consider the following bins:  $[-\infty, q_1)$ ,  $[q_1, q_{10})$ ,  $[q_{10}, q_{25})$ ,  $[q_{25}, q_{50})$ ,  $[q_{50}, q_{75})$ ,  $[q_{75}, q_{90})$ ,  $[q_{90}, q_{99})$ ,  $[q_{99}, \infty)$ . If for each covariate the quantiles are not unique, we take the set of unique quantile and construct bins accordingly. For example suppose that for the spread  $q_{50} = q_{75} = q_{90}$  while all other quantiles are unique. Then, the bins we use for the spread are  $[-\infty, q_1)$ ,  $[q_1, q_{10})$ ,  $[q_{10}, q_{25})$ ,  $[q_{25}, q_{50})$ ,  $[q_{50}, \infty)$ . Also note that binning in  $[-\infty, q_1)$  or  $[0, q_1)$  produces the same result for covariates that take nonnegative values. According to the aforementioned binning strategy, the total number of parameters to be estimated, including a parameter for a constant, is at most 177 when all quantiles are unique. However, due to nonuniqueness of some of the quantiles for the spread, the actual number of parameters to estimate becomes 168.

## 5.3 The Models

To understand the importance of the nonlinearity in the impact of order book variables and the self exciting nature of the trading events, we consider a number of model specifications within the current framework. We list these next.

E:  $h(t) = 1$ ,  $g(t) = X(t)'b$  (no self excitation,  $Z \mapsto X$  by one-hot encoding)

H01:  $h(t)$  as in (2) with  $L = 1$ ,  $g(t) = 1$  (no book information)

H02:  $h(t)$  as in (2) with  $L = 2$ ,  $g(t) = 1$  (no book information)

H1:  $h(t)$  as in (2) with  $L = 1$ ,  $g(t) = X(t)'b$  ( $Z \mapsto X$  by one-hot encoding)

H2:  $h(t)$  as in (2) with  $L = 2$ ,  $g(t) = X(t)'b$  ( $Z \mapsto X$  by one-hot encoding)

H1L:  $h(t)$  as in (2) with  $L = 1$ ,  $g(t) = Z(t)'b$  (linear raw covariates with  $b$  in the reals)

H2L:  $h(t)$  as in (2) with  $L = 2$ ,  $g(t) = Z(t)'b$  (linear raw covariates with  $b$  in the reals)

In the above, the raw covariates in Table 3 are denoted by  $Z$ . For the nonlinear models, these are mapped to  $X$  via one-hot encoding. The model is then estimated under a positivity constraint and under an upper bound constraint. In particular, model E corresponds to no self excitation, so that conditioning on  $g(t)$  the hazard functions is the one of an exponential distribution. Models H01 and H02 are Hawkes processes with kernel equal to the sum of one and two exponential functions, respectively. Models H1 and H2 are like H01 and H02, respectively, times  $g(t)$ , i.e. they include order book information. Models H1L and H2L are like H1 and H2, respectively, but only allow linear impact of the raw covariates. In this case the  $b$  coefficient is allowed to be negative so that the impact can have a negative sign. Then, to avoid a negative intensity, we impose a floor on the intensity when testing out of sample. Additional details on estimation constraints can be found in Section A.3.2 of the Appendix.

## 5.4 Results

Data summary statistics show that the size of the data is large across all instruments, though there is some degree of variation (Table 4). For example, during liquid periods, trades information is disseminated even every few microseconds.

We estimate our models for buy and sell events separately. The positivity constraint and the constraint on the sum of the coefficients lead to a relatively sparse estimator. For the model with  $L = 1$  in (2), which we defined as H1 in Section 5.3, the average number of nonzero estimated coefficients is roughly between 20% to 30% across the four stocks both for Any Trade Arrivals and the Large Trade Arrivals.

We inspect the number of nonzero coefficient for each covariate. Covariates for which the coefficients of the one-hot encoding are all zero do not enter the model and could be deemed as unimportant. We note that there is a reasonable level of consistency across the different stocks. The detailed results are reported in Tables 5 and 6.

We conduct the test of Section 4.1 in order to answer the questions from Section 5. For example, we write E-H1 to mean that the loglikelihood ratio is constructed as the likelihood of model E minus the loglikelihood of model H1 (see Section 5.3). Then, we compute the following test statistics:

1. H01-H1, H02-H1, and H02-H2: A large negative value of the test statistic implies that the order book provide information in addition to what is captured by the self exciting nature of event arrivals;
2. E-H1: a large negative value means that self excitation is important even after accounting for order book information;



**Table 4:** Sample Size Statistics. The total number  $N$  of events that correspond to Any Trade Arrivals and Large Trade Arrivals is reported together with the total number  $m$  of book updates with unique time stamps. The number of days is 42 for the period 01/Mar/2019 - 30/04/2019.

	N		m
	Any Trade Arrivals	Large Trade Arrivals	
AMZN	631,370	407,130	17,130,000
CSCO	295,220	107,930	27,943,000
DIS	493,100	292,820	24,938,000
KO	121,210	52842	13,956,000

3. H1L-H1, H2L-H1, H1L-H2, and H2L-H2: A large negative value means that the impact of the order book is nonlinear irrespective of the kernel chosen in (2);
4. H2-H1: A large positive value means that the simple exponential kernel is not sufficient to capture the self excitation of the trading events.

We find that accounting for both self excitation and order book information (models H1 and H2) is important (Points 1 and 2). We also find that the impact of of the order book is nonlinear (Point 3). Finally, the self exciting nature of trading events usually requires the use of a kernel more complex that the simple exponential one (Point 4). The details from the tests are in Table 7. The values of the test statistics are very large in absolute value because the sample size  $T$  is large (see the remarks about the power of the test just before Theorem 4).

We conclude mentioning that one-hot encoding can produce plots that are interpretable and regular even without any constraint (Figures 1 and 2).

## 6 Conclusion

This paper presented a Hawkes process augmented by order book information. The model is the standard Hawkes process with a multiplicative term which is a function of order book variables. We showed conditions under which the process is stationary, ergodic and strong mixing. We focused on a flexible estimation procedure suitable for large datasets where the results are however intuitively simple to interpret. Using theoretical results as well as simulation examples we were able to justify the procedure. The results showed that the convergence rates only deteriorate at a logarithmic rate with the number of parameters. The motivation for the algorithm presented here and its study is to allow us to use flexible techniques to model the impact of the covariates derived from the order book, possibly using very large datasets. In particular, we focus on one-hot encoding. This maps the original covariates into a higher dimensional space, a case covered by our methodology. Our application to four stocks traded

**Table 5:** Active Covariates for Any Trade Arrivals. The total number of estimated non zero coefficients (from one-hot encoding) for each covariate (as in Table 3) are reported. The total sum (Sum) of nonzero coefficients and their relative proportion (Proportion), out of the estimated 168 parameters, are also reported. The models considered are E, H1 and H2, as defined in Section 5.3.

	AMZN			CSCO				DIS			KO		
	E	H1	H2	E	H1	H2		E	H1	H2	E	H1	H2
AMZN_Dur90	2	3	3	2	2	2	DIS_Dur90	3	4	4	1	3	2
AMZN_Dur98	3	3	3	3	3	3	DIS_Dur98	4	4	4	1	3	3
AMZN_Seas	4	5	5	1	2	2	DIS_Seas	2	3	3	0	1	1
AMZN_Spread	4	4	4	0	1	1	DIS_Spread	1	1	1	1	1	1
AMZN_TrdImb98	3	3	3	0	0	0	DIS_TrdImb98	3	2	2	1	0	0
AMZN_VolImb1	0	0	0	0	0	0	DIS_VolImb1	4	4	4	3	3	3
AMZN_VolImb2	0	0	0	0	0	0	DIS_VolImb2	3	3	3	1	1	1
AMZN_VolImb3	0	0	0	3	3	3	DIS_VolImb3	2	2	2	3	4	4
CSCO_Dur90	1	1	1	4	4	4	KO_Dur90	0	0	0	4	4	4
CSCO_Dur98	3	3	3	5	5	5	KO_Dur98	1	1	1	3	4	4
CSCO_Spread	1	1	1	0	0	0	KO_Spread	0	0	0	1	1	1
CSCO_TrdImb98	1	0	0	2	2	2	KO_TrdImb98	0	0	0	2	2	2
CSCO_VolImb1	3	3	3	4	5	5	KO_VolImb1	1	1	1	5	5	5
CSCO_VolImb2	2	3	3	3	3	3	KO_VolImb2	0	0	0	3	3	3
CSCO_VolImb3	0	0	0	1	2	2	KO_VolImb3	0	0	0	1	1	1
SPY_Dur90	3	3	3	3	3	3	SPY_Dur90	2	2	2	4	4	4
SPY_Dur98	4	4	4	2	3	3	SPY_Dur98	1	1	1	3	3	3
SPY_Spread	2	2	2	2	2	2	SPY_Spread	2	2	2	2	2	2
SPY_TrdImb98	1	2	2	1	1	1	SPY_TrdImb98	1	1	1	2	2	2
SPY_VolImb1	3	3	3	2	3	3	SPY_VolImb1	1	3	3	3	3	3
SPY_VolImb2	1	1	1	1	1	1	SPY_VolImb2	0	0	0	1	1	1
SPY_VolImb3	0	1	1	1	1	1	SPY_VolImb3	0	0	0	2	1	1
Sum	41	45	45	40	46	46		31	34	34	47	52	51
Proportion	0.24	0.27	0.27	0.24	0.27	0.27		0.18	0.20	0.20	0.28	0.31	0.30

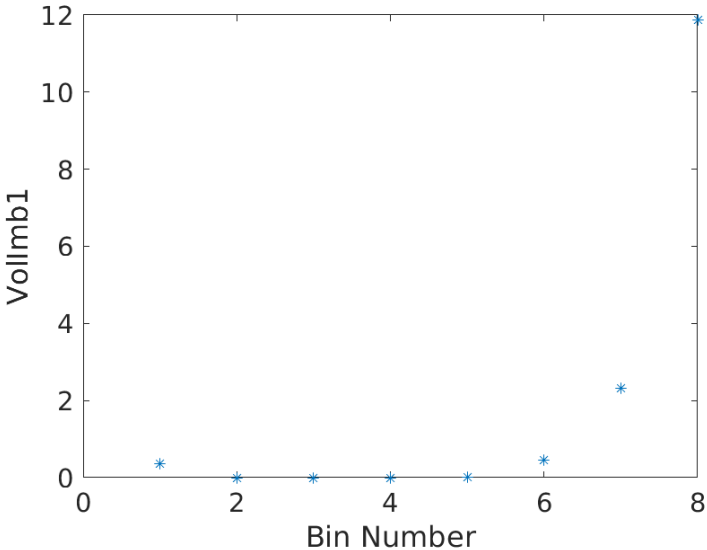
**Table 6:** Active Covariates for Large Trade Arrivals. The total number of estimated non zero coefficients (from one-hot encoding) for each covariate (as in Table 3) are reported. The total sum (Sum) of nonzero coefficients and their relative proportion (Proportion), out of the estimated 168 parameters, are also reported. The models considered are E, H1 and H2, as defined in Section 5.3.

	AMZN			CSCO				DIS			KO		
	E	H1	H2	E	H1	H2		E	H1	H2	E	H1	H2
AMZN_Dur90	2	2	2	2	2	2	DIS_Dur90	3	3	3	1	1	1
AMZN_Dur98	3	3	3	2	2	2	DIS_Dur98	3	3	3	2	2	2
AMZN_Seas	4	4	4	3	3	3	DIS_Seas	4	3	3	2	2	2
AMZN_Spread	4	4	4	1	0	0	DIS_Spread	1	1	1	1	1	1
AMZN_TrdImb98	2	2	2	0	0	0	DIS_TrdImb98	2	2	2	0	0	0
AMZN_VolImb1	4	4	4	0	0	0	DIS_VolImb1	4	4	4	1	1	1
AMZN_VolImb2	0	0	0	0	0	0	DIS_VolImb2	2	3	3	1	1	1
AMZN_VolImb3	0	0	0	2	1	1	DIS_VolImb3	1	1	1	2	2	2
CSCO_Dur90	0	1	1	4	4	4	KO_Dur90	0	0	0	4	3	3
CSCO_Dur98	3	3	3	3	3	3	KO_Dur98	1	1	1	3	3	4
CSCO_Spread	1	1	1	0	0	0	KO_Spread	0	0	0	0	0	0
CSCO_TrdImb98	0	0	0	1	0	0	KO_TrdImb98	0	0	0	1	0	0
CSCO_VolImb1	2	2	2	3	3	3	KO_VolImb1	0	1	1	3	3	3
CSCO_VolImb2	2	2	2	2	2	2	KO_VolImb2	0	0	0	1	0	0
CSCO_VolImb3	0	0	0	1	1	1	KO_VolImb3	0	0	0	1	0	0
SPY_Dur90	3	3	3	3	3	3	SPY_Dur90	2	2	2	3	3	3
SPY_Dur98	4	4	4	3	3	3	SPY_Dur98	1	1	1	3	3	3
SPY_Spread	2	2	2	2	2	2	SPY_Spread	2	2	2	2	2	2
SPY_TrdImb98	1	1	1	1	1	1	SPY_TrdImb98	1	1	1	2	2	2
SPY_VolImb1	3	3	3	3	3	3	SPY_VolImb1	1	2	2	3	3	3
SPY_VolImb2	1	1	1	1	1	1	SPY_VolImb2	0	0	0	1	0	0
SPY_VolImb3	1	1	1	1	0	0	SPY_VolImb3	0	0	0	0	0	0
Sum	42	43	43	38	34	34		28	30	30	37	32	33
Proportion	0.25	0.26	0.26	0.23	0.20	0.20		0.17	0.18	0.18	0.22	0.19	0.20

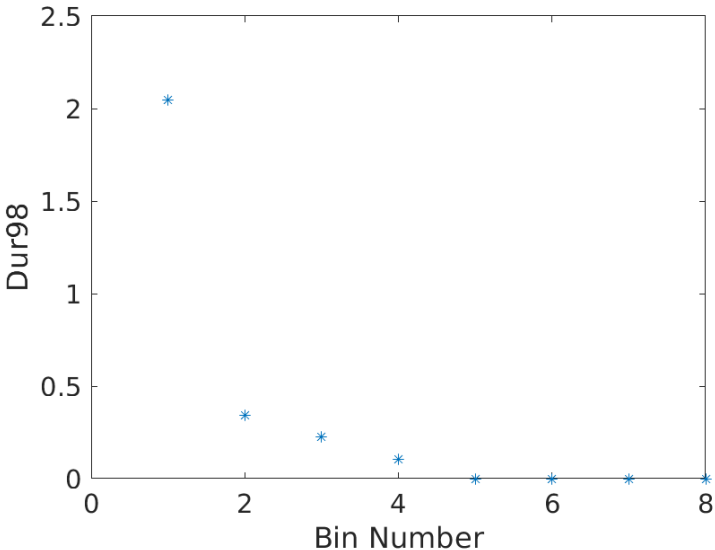
**Table 7:** Test of Model Restrictions. The values for the test statistic in (12) are reported. The test is constructed so that the last 5 trading days are used as test and the previous days for estimation. Models are defined in Section 5.3. The null hypothesis is that two models perform the same out of sample. A large positive value means that the null is rejected and the first model performs better than the second. The reverse applies for large negative values. The statistic is standard normal, and the “simple” models are almost always rejected.

		Any Trade Size				Large Trade Size			
		AMZN	CSCO	DIS	KO	AMZN	CSCO	DIS	KO
E-H1	Buy	-109.88	-70.20	-93.44	-44.41	-81.40	-22.48	-66.48	-16.57
	Sell	-112.77	-72.18	-86.61	-43.36	-82.20	-26.36	-59.26	-20.23
H01-H1	Buy	-20.28	-7.84	-6.43	-12.83	-37.01	-35.13	-14.08	-21.27
	Sell	-15.53	-8.32	-9.65	-39.54	-35.27	-41.36	-21.21	-38.94
H02-H1	Buy	-20.96	-3.10	39.23	-4.11	-20.53	-33.75	-4.24	-21.34
	Sell	-15.28	-6.82	35.15	-24.35	-18.58	-39.07	-12.37	-38.00
H2-H1	Buy	-73.86	-44.02	63.92	27.65	41.51	11.13	20.55	8.91
	Sell	-54.83	49.30	59.82	18.56	42.20	10.39	27.04	7.38
H1L-H1	Buy	-27.75	-13.95	-7.61	-11.54	-27.35	-25.92	-10.68	-23.18
	Sell	-24.68	-15.76	-2.95	-39.38	-26.78	-32.51	-8.68	-41.81
H2L-H1	Buy	-27.99	-12.27	29.98	2.42	-20.67	-25.14	-5.88	-23.05
	Sell	-24.42	-17.14	33.22	-22.65	-19.54	-31.02	-3.14	-40.92
H02-H2	Buy	-19.02	-2.77	-1.07	-18.27	-38.17	-34.68	-14.76	-22.12
	Sell	-15.09	-9.02	-5.28	-34.93	-36.76	-40.34	-26.29	-40.28
H1L-H2	Buy	-26.70	-13.71	-40.41	-24.05	-34.81	-26.44	-16.76	-23.96
	Sell	-24.57	-17.45	-38.13	-40.26	-34.61	-33.22	-16.58	-42.91
H2L-H2	Buy	-26.95	-12.03	-3.63	-10.28	-29.10	-25.79	-12.24	-23.88
	Sell	-24.31	-18.82	-0.90	-32.08	-28.33	-32.07	-11.16	-43.16

**Figure 1:** CSCO Volume Imbalance Level 1. The model is for  $L = 1$  in (2) for CISCO and Any Trade Arrivals of buy trades. The estimated coefficients  $b_k$  for VolImb1 are plotted on the Y-axis as a function of  $k$ , which is the bin number out of 8 bins based on quantiles (see Section 5.2). Bin numbers greater than 4 correspond to positive values of VolImb1.



**Figure 2:** CSCO Duration (Dur98). The model is for  $L = 1$  in (2) for CISCO and Any Trade Arrivals of buy trades. The estimated coefficients  $b_k$  for Dur98 are plotted on the Y-axis as a function of  $k$ , which is the bin number out of 8 bins based on quantiles (see Section 5.2). For example, bin number 1 corresponds to Dur98 smaller that the 1% quantile.



on the NYSE showed the importance of using order book information. In summary, our sample testing procedure shows that nonlinearity of the order book adds value to the self exciting nature of high frequency event arrivals.

## References

- [1] Alaya, M.Z., S. Bussy, S. Gaïffas, and A. Guilloux (2019) Binarisity: A Penalization for One-Hot Encoded Features in Linear Supervised Learning. *Journal of Machine Learning Research* 20, 1-34.
- [2] Bacry, E., I. Mastromatteo and J.-F. Muzy (2015) Hawkes Processes in Finance. *Market Microstructure and Liquidity* 1, No.1.
- [3] Bauwens, L. and N. Hautsch (2009) Modelling Financial High Frequency Data Using Point Processes. In T.G. Andersen, R.A. Davis, J.-P. Kreiss and T. Mikosch (eds.), *Handbook of Financial Time Series*, 953-982. New York: Springer.
- [4] Beck, A. and L. Tetruashvili (2013) On the Convergence of Block Coordinate Descent Type Methods. *SIAM Journal on Optimization* 23, 2037–2060.
- [5] Brémaud, P. (1981) *Point Processes and Queues: Martingales Dynamics*. New York: Springer.
- [6] Brémaud, P. and L. Massoulié (1996) Stability of Nonlinear Hawkes Processes. *Annals of Probability* 24, 1563-1588.
- [7] Cartea, A., S.N. Cohen and S. Labyad (2021) Gradient-based Estimation of Linear Hawkes Processes with General Kernels. <https://arxiv.org/abs/2111.10637>
- [8] Cheysson, F. and G. Lang (2022) Spectral estimation of Hawkes processes from count data. Forthcoming in the *Annals of Statistics*, <https://arxiv.org/abs/2003.04314>.
- [9] Cont, R., A. Kukanov and S. Stoikov (2014) The Price Impact of Order Book Events. *Journal of Financial Econometrics* 12, 47-88.
- [10] Da Fonseca, J. and R. Zaatour (2014) Hawkes Process: Fast Calibration, Application to Trade Clustering, and Diffusive Limit. *Journal of Futures Markets* 34, 548-579.
- [11] Daley, D.J. and D. Vere-Jones (2003) *An Introduction to the Theory of Point Processes*, Volume II. New York: Springer.
- [12] Engle, R.F. and J.R. Russell (1998) Autoregressive Conditional Duration: A New Model for Irregularly Spaced Transaction Data. *Econometrica* 66, 1127-1162.

- [13] Filimonov, V and D. Sornette (2015) Apparent Criticality and Calibration Issues in the Hawkes Self-Excited Point Process Model: Application to High-Frequency Financial Data. *Quantitative Finance* 15, 1293-1314.
- [14] Fosset, A., J.-P. Bouchaud and Michael Benzaquen (2020) Endogenous Liquidity Crises. <https://arxiv.org/abs/1912.00359>.
- [15] Friedman, J., T. Hastie, H. Höfling, and R. Tibshirani (2007) Pathwise Coordinate Optimization. *Annals of Applied Statistics* 1, 302-332.
- [16] Gaïffas, S. and A. Guillaoux (2012) High-Dimensional Additive Hazards Models and the Lasso. *Electronic Journal of Statistics* 6, 522-546.
- [17] Gao, X. and L. Zhu (2018) Functional Central Limit Theorems for Stationary Hawkes Processes and Application to Infinite-Server Queues. *Queueing Systems* 90, 161–206.
- [18] Grossman, S.-J. and J.E. Stiglitz (1980) On the Impossibility of Informationally Efficient Markets. *American Economic Review* 70, 393-408.
- [19] Hall, A.D. and N. Hautsch (2007) Modelling the Buy and Sell Intensity in a Limit Order Book Market. *Journal of Financial Markets* 10, 249-286.
- [20] Hawkes, A.G. (1971) Spectra of Some Self-Exciting and Mutually Exciting Point Processes. *Biometrika* 58, 83–90.
- [21] Hawkes, A.G. and D. Oakes (1974) A Cluster Process Representation of a Self-Exciting Process. *Journal of Applied Probability* 11, 493-503.
- [22] Huang, R. and T. Polak (2011) LOBSTER: The Limit Order Book Reconstructor. Technical Report, School of Business and Economics, Humboldt Universität zu Berlin.
- [23] Kallneberg, O. (1997) *Foundations of Modern Probability*. New York: Springer.
- [24] Kercheval, A.N., Y. Zhang (2015) Modelling High-Frequency Limit Order Book Dynamics with Support Vector Machines. *Quantitative Finance* 15, 1-15.
- [25] Kirchner, M. (2017) An Estimation Procedure for the Hawkes Process. *Quantitative Finance* 17, 571-595.
- [26] MacKenzie, D. (2017) A Material Political Economy: Automated Trading Desk and Price Prediction in High - Frequency Trading. *Social Studies of Science* 47, 172-194 .
- [27] Morariu-Patrichi, M. and M.S. Pakkanen (2022) State-Dependent Hawkes Processes and their Application to Limit Order Book Modelling. *Quantitative Finance* 22, 563-583.

- [28] Mounjid, O., M. Rosenbaum and P. Saliba (2019) From Asymptotic Properties of General Point Processes to the Ranking of Financial Agents. <https://arxiv.org/abs/1906.05420>.
- [29] Mucciante, L. and A. Sancetta (2022) Estimation of a High Dimensional Counting Process Without Penalty for High Frequency Events. *Econometric Theory*: <https://doi.org/10.1017/S0266466622000238>.
- [30] Nishiyama, Y. (2000) Weak Convergence of Some Classes of Martingales with Jumps. *Annals of Probability* 28, 685-712.
- [31] Ogata, Y. (1978) The Asymptotic Behaviour of the Maximum Likelihood Estimator for Stationary Point Processes. *Annals of the Institute of Statistical Mathematics* 30, 243-261.
- [32] Ogata, Y. and H. Akaike (1982) On Linear Intensity Models for Mixed Doubly Stochastic Poisson and Self-Exciting Point Processes. *Journal of the Royal Statistical Society B* 44, 102-107.
- [33] Sancetta, A. (2016) Greedy Algorithms for Prediction. *Bernoulli* 22, 1227-1277.
- [34] Sancetta, A. (2018) Estimation for the Prediction of Point Processes with Many Covariates. *Econometric Theory* 34, 598-627.89-107.
- [35] Tsybakov, A.B. (2003) Optimal Rates of Aggregation. *Proceedings of COLT-2003, Lecture Notes in Artificial Intelligence*, 303-313.
- [36] Wu, P., M. Rambaldi, J.-F. Muzy and E. Bacry (2019) Queue-Reactive Hawkes Models for the Order Flow. <https://arxiv.org/abs/1901.08938>.
- [37] Zhang, Y., J. Duchi and M. Wainwright (2015) Divide and Conquer Kernel Ridge Regression: A Distributed Algorithm with Minimax Optimal Rates. *Journal of Machine Learning Research* 16, 3299-3340.



# Appendix

In Section A.1, we give further information on the interpretation of the model in terms of exogenous and endogenous information arrival. Section A.2 proves all the results stated in the paper. Section A.3 includes further details regarding the empirical application.

## A.1 Model Interpretation: Exogenous Information Versus Endogenous Market Activity

We expand on the model representation in terms of an immigration birth process (Hawkes and Oakes, 1974). In our context, immigration means exogenous information arrival, while birth represents endogenously generated information. When  $g_0$  is constant, exogenous information is information that is independent of arrivals, i.e. independent of  $N$ , hence strictly exogenous. Endogenous information at time  $t$  is the information that has been generated by the trading activity up to time  $t$ , i.e.  $(N(s))_{s < t}$ . Here this interpretation is generalized.

Write  $\kappa(t) := \sum_{l=1}^L d_{0,l} e^{-a_{0,l} t}$ . This is the kernel of the process in (2). Let  $T_n$  be the  $n^{\text{th}}$  event arrival, and suppose that this is the latest arrival time before time  $t$ . Then,  $g_0(t) c_0$  represents the unit rate of arrival of immigrants at time  $t$ . Each immigrant has descendants who procreate ad infinitum. If  $g_0$  were constant, the immigration arrival would be strictly exogenous, as in the standard Hawkes process. Here, the intensity of immigrants arrivals depends on what we call the environment, and it is captured by the process  $g_0$  at time  $t$ .

At time  $t$  there are  $n$  independent individuals in the population. Nobody dies, but they all age. At time  $t$ , the  $i^{\text{th}}$  individual has age  $t - T_i$  and the unit intensity at which they give birth is  $g_0(t) \kappa(t - T_i)$ . Given that  $\kappa$  is decreasing, an older individual has lower probability of giving birth. However, their chances can be increased if  $g_0(t)$  is higher than normal. Hence, the introduction of the process  $g_0$  makes the individuals in the populations dependent, unlike the standard Hawkes process. Conditional on the environment, i.e. conditioning on  $g_0(t)$ , the probability of giving birth is independent across the individuals in the population. Conditioning on  $g_0(t)$ , the intensity of new arrivals in the the population, either via immigration or birth is  $g_0(t) c_0 + g_0(t) \sum_{i=1}^n \kappa(t - T_i)$ . This is exactly (1) once we substitute the definition of  $\kappa$ .

Suppose that  $T_0$  is the time of arrival of an immigrant. During their infinite lifetime, their expected total number of offspring is  $n_{g_0} := \mathbb{E}_{T_0} \int_{T_0}^{\infty} g(t) \kappa(t - T_0) dt$ , where  $\mathbb{E}_{T_0}$  is expectation conditional on the information at time  $T_0$ . If  $g_0$  is constant,  $n_{g_0}$  does not depend on  $T_0$ . Then, the number  $n_{g_0}$  is called branching ratio. When less than one, the probability of extinction is one. This is necessary to ensure that the effect of any event arrival eventually dies out so that the process is stationary. If  $g_0$  is not constant, but uniformly bounded, we still have  $n_{g_0} < 1$

if  $\sup_t |g_0(t)|_\infty \left( \sum_{l=1}^L d_{0,l}/a_{0,l} \right)^{-1} < 1$ , where  $\sum_{l=1}^L d_{0,l}/a_{0,l} = \int_0^\infty \kappa(t) dt$ .

In the context of the present paper, exogenous news are represented by the immigrants. However, information is costly and the goal of an informed participant is to trade revealing as little information as possible (Grossman and Stiglitz, 1980). Trading is also costly, so that the trading strategy of both informed and uninformed traders must depend on the environment, e.g. the state of the order book. Optimal execution strategies used to minimize private information transmission and trading cost needs to take into account the order book and market microstructures. In a frictionless market where cost of information is zero, prices would adjust immediately and not result in any prolonged trading activity. In the context of a Hawkes process, this means that each immigrant would have no descendants and trading is only kept alive by exogenous news events. For examples, this is the case when  $\min_l a_{0,l} \rightarrow \infty$  in (2). On the opposite side of the spectrum, information is never fully absorbed by the market and trading is kept alive by the descendants and an equilibrium price is never reached. Clearly, this is unnatural. Eventually information is absorbed by the market and trading is only kept alive by exogenous information transmission. This means that  $n_{g_0} < 1$ . The extent to which  $n_{g_0}$  is close to one tells us about the level of endogenous market interaction that spawned out of any news. Equivalently, it represents the average fraction of endogenously generated events. Clearly a market where information is not quickly absorbed and trading is just the result of endogenous interaction (numerous descendants from each immigrant) is a market prone to fragility. Filomonov and Sornette (2015) call this market reflexivity.

## A.2 Proofs

The intensity  $\lambda_0(t) = \lambda_0(\omega, t)$  is a continuous time stochastic process, i.e. a function of two variables  $t \geq 0$  and  $\omega \in \Omega$  where  $(\Omega, \mathcal{B}, P)$  is a probability space and for each  $t$ ,  $\lambda_0(t, \cdot)$  is measurable on  $\Omega$ . Similarly  $h_0(t) = h_0(\omega, t)$  and  $X(t) = X(\omega, t)$ . The covariate process  $X$  and the baseline intensity are predictable processes. These quantities are supposed to be left continuous with right hand limits. For ease of notation, we may freely switch between  $\lambda_0(t)$  and  $h_0(t)g_0(t)$  and compactly write  $\lambda_0 = h_0g_0$ . To make the notation simpler and more readable, we shall write  $Ph_0g_0$  to mean  $\int_\Omega h_0(\omega, 0)g_0(\omega, 0)dP(\omega)$  and similarly for other quantities.

Given that all the quantities are finite dimensional, both  $Q_T$  and  $R_T$  (in (6) and (7)) are Fréchet differentiable. Throughout to ease notation, for any function  $f$  that is bounded below by a constant, we use  $c_f$  to indicate that constant. If bounded above by a constant, we use  $C_f$  to indicate such constant. If  $\mathcal{F}$  is a class of functions bounded below and/or above by constants, we denote such constants by  $c_{\mathcal{F}}$  and  $C_{\mathcal{F}}$  respectively. Finally, to avoid making reference to constants, we use  $\lesssim$  when the left hand side (l.h.s.) is bounded by a constant times the right hand side (r.h.s.) and  $\asymp$  when the l.h.s. is bounded below and above by constants

times the r.h.s..

We prove Theorem 1 first, followed by Theorem 3 and then 2, as the latter easily follows from the former. We then prove Theorem 4.

### A.2.1 Proof of Theorem 1

We shall use as much as possible the notation in Brémaud and Massoulié (1996), BM96 henceforward. Throughout, we view  $N$  as an element in  $\mathbb{M}$ , the space of Radon measures on  $\mathbb{R}$ . Let  $\omega$  be a Poisson point process associated with the arrival times  $(\tau_i)_{i \geq 1}$ . Theorem 4 and Lemma 3 in BM96 show that, iteratively for  $n = 0, 1, 2, \dots$ , given  $\omega$  and an intensity  $\lambda^n$  we can construct a point process  $N^n$  that has such intensity (BM96, eq. 14). We use the same notation as in BM96, hence here  $N^n$  and  $\lambda^n$  are sequences and should not be confused with  $n$  powers. The intensity is constructed to satisfy  $\lambda^{n+1}(t) = \phi\left(\int_{(-\infty, t)} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}(t-s)} dN^n(s)\right)$  for  $\phi: \mathbb{R} \rightarrow [0, \infty)$  and such that  $\phi(x) \leq \alpha + \beta x$  for some positive constants  $\alpha, \beta$ . In our case  $\phi$  is replaced by a predictable process  $\phi_t^n(x) = g_0^n(t)(c_0 + x)$  where, by the Regularity Conditions,  $g_0^n(t) = f_0(W(t), S_t N^n)$ . It also follows that because  $\omega$  is stationary for the shift operator, so is  $N^n$  (BM96, first and last paragraphs on p.1273 and 1274, respectively). Hence,  $g_0^n(t)$  is stationary for every iteration  $n$  by Condition 1. Using the upper bound for  $\phi_t^n(\cdot)$  and stationarity of  $\lambda^n$ , we have that, for  $t = 0$ ,

$$\begin{aligned} \mathbb{E}\lambda^{n+1}(0) &= \mathbb{E}\phi_0^n\left(\int_{(-\infty, 0)} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}s} dN^n(s)\right) \\ &\leq |g_0^n|_\infty \left(c_0 + \mathbb{E} \int_{(-\infty, 0)} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}s} dN^n(s)\right) \\ &= |g_0^n|_\infty \left(c_0 + \mathbb{E}\lambda^n(0) \int_{(0, \infty)} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}s} ds\right). \end{aligned}$$

This is in the same form as the first display on page 1575 of BM96. It is then easy to see that the above is finite if  $|g_0^n|_\infty \int_0^\infty \sum_{l=1}^L d_{0,l} e^{-a_{0,l}s} ds < 1$  which is the case if  $|g_0^n|_\infty \sum_{l=1}^L \frac{d_{0,l}}{a_{0,l}} < 1$ . Given that the intensity is not explosive for any  $n$ , ergodicity of  $N^n$  follows from the discussion on page 1573 of BM96. To show that as  $n \rightarrow \infty$  the intensity converges to  $\lambda_0$  we use the same argument based on monotonicity, as done in the last steps of the proof of Theorem 4 of BM96 on page 1575. This requires that  $\lambda^n$  and  $N^n$  are increasing in  $n$ . Following step by step the argument in BM96, this is true, here, if  $g_0^n(t)$  is increasing in  $n$ . This is the case by the regularity conditions.

We then use Theorem 1 in BM96 to see that the stationary distribution is unique. For  $x \geq 0$ , use  $\phi_t(x) = g_0(t)(c_0 + x)$  instead of fixed  $\phi(x)$ , as done there. We know that  $\phi_t(\cdot)$  is Lipschitz with Lipschitz constant  $|g_0|_\infty$  satisfying  $|g_0|_\infty \int_0^\infty \sum_{l=1}^L d_{0,l} e^{-a_{0,l}s} ds < 1$ . Then

we can see that their proof follows through based on the remarks above. This shows the stationarity of the process.

The process is  $\alpha$ -mixing with exponentially decaying mixing coefficients by Theorem 1 in Cheysson and Lang (2022).

The fact that the process restricted to  $\mathcal{A}$  will converge eventually to a stationary point process follows from Point b. in Theorem 1 of BM96. The result holds for any process whose initial condition satisfies  $\sup_{t \geq 0} \varepsilon_v(t) < \infty$  and  $\lim_{t \rightarrow \infty} \varepsilon_v(t) = 0$  almost surely for all  $v > 0$ , where  $\varepsilon_v(t) := \int_{t-v}^t \int_{(-\infty, 0)} \sum_{l=1}^L d_{0,l} e^{-a_{0,l}(r-s)} dN(s) dr$ . Clearly, this is the case for the initial condition  $\mathcal{A} := \{N(t) = 0 : t \leq 0\}$  that we are considering.

### A.2.2 Proof of Theorem 3

At first we state the following result on the moments of the process.

**Lemma 1** *Under the Regularity Conditions, for any finite interval  $[r, s]$ ,  $\mathbb{E}N^p([r, s])$  for any  $p < \infty$ , where  $N([r, s]) := \int_r^s dN(t)$ . Moreover,  $\mathbb{E}\lambda_0^4(t) < \infty$  uniformly in  $t \geq 0$ .*

**Proof.** This follows from a trivial modification of Lemmas 1 and 2 in Zhu (2013). Lemma 1 in Zhu (2013) says that nonlinear Hawkes processes as in the proof of Theorem 1, and that start at zero (i.e. conditioning on the empty past), have moment of all order. By stationarity, Lemma 2 in Zhu (2013) extend the result to the same Hawkes processes without conditioning on the empty past. Following the proof of Lemma 1 in Zhu (2013), we can see that the lemma also applies to our model as long as  $|g_0|_\infty \leq B < \infty$ . In the proof of Lemma 1 in Zhu (2013), we just need to replace his  $\alpha$  and  $\lambda(0)$  with our  $B$  and  $c_0 B$ , respectively. Hence, our Hawkes process has moments of all orders.

Finally, Lemma 15 in Guo and Zhu (2018) says that for a standard linear Hawkes process the intensity has finite fourth moment. Following the proof of Lemma 15 in Guo and Zhu (2018), using the fact that  $|g_0|_\infty \leq B < \infty$ , we see that the result applies to our process as well. The proof requires that  $\mathbb{E}N^4([r, s]) < \infty$  for some suitably small interval  $[r, s]$ . This is the case by the first statement of the present lemma. ■

Completing the square, we have that

$$R_T(g; h) = \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} - g \right)^2 d\mu - \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} \right)^2 d\mu - \frac{2}{T} \int_0^T \frac{g}{h} dM. \quad (\text{A.1})$$

Here,  $dM = dN - \lambda_0 d\mu$  so that  $M$  is a martingale. Recall that  $\mathcal{G} = \{g = X'b : |b|_1 \leq B, b_k \geq 0, k \leq K\}$ . By convexity,

$$\sup_{b: b \geq 0, |b|_1 \leq B} \left| \frac{1}{T} \int_0^T \sum_{k=1}^K b_k X_k dM \right| \leq B \max_{k \leq K} \left| \frac{1}{T} \int_0^T X_k dM \right|. \quad (\text{A.2})$$

The main ingredient in the proof is to show convergence of  $\max_{k \leq K} \sup_h \left| \frac{1}{T} \int_0^T (X_k/h) dM \right|$  to zero. This requires to control the oscillations of the process over finite partitions of the parameter space of elements in  $\mathcal{H}$ . To this end we use the following.

**Lemma 2** *Let  $h^{(j)}(t) = c^{(j)} + \sum_{l=1}^L h_l^{(j)}(t)$  where  $h_l^{(j)}(t) = \int_{(-\infty, t)} d_l^{(j)} e^{-a_l^{(j)}(t-s)} dN(s)$ ,  $j = 1, 2$  such that  $|c^{(1)} - c^{(2)}| \leq \epsilon$ ,  $|d_l^{(1)} - d_l^{(2)}| \leq \epsilon/(2L)$ ,  $|a_l^{(1)} - a_l^{(2)}| \leq \epsilon/(2\bar{d}L)$  where  $a_l^{(1)}, a_l^{(2)} \geq \underline{a} := \min\{a \in \mathcal{A}\}$  and  $d_l d_l' \leq \bar{d} := \max\{d \in \mathcal{D}\}$ . Under the Regularity Conditions, for any arbitrary but fixed  $h^{(1)}$  and  $h^{(2)}$  satisfying the above,*

$$\frac{\int_0^T |h^{(1)}(t) - h^{(2)}(t)|^2 \lambda_0(t) dt}{\epsilon^2} \lesssim T + \int_0^T \left[ \int_{(-\infty, t)} e^{-a(t-s)} dN(s) \right]^2 \lambda_0(t) dt \quad (\text{A.3})$$

for any  $a < \underline{a}$ . For any  $a \in (0, \underline{a})$  we have that  $\frac{1}{T} \int_0^T \left[ \int_{(-\infty, t)} e^{-a(t-s)} dN(s) \right]^2 \lambda_0(t) dt \rightarrow C_a$  in probability, where  $C_a$  is a finite constant.

**Proof.** Write

$$\left| h^{(1)}(t) - h^{(2)}(t) \right| \leq \left| c^{(1)} - c^{(2)} \right| + \sum_{l=1}^L \left| h_l^{(1)}(t) - h_l^{(2)}(t) \right|. \quad (\text{A.4})$$

By the mean value theorem and basic inequalities,

$$\begin{aligned} \left| h_l^{(1)}(t) - h_l^{(2)}(t) \right| &\leq \int_{(-\infty, t)} \left| d_l^{(1)} - d_l^{(2)} \right| e^{-a_l^{(1)}(t-s)} dN(s) \\ &\quad + \left| a_l^{(1)} - a_l^{(2)} \right| \max_{\tau \in [0, 1]} \int_{(-\infty, t)} d_l^{(2)} \exp\{-a_{l, \tau}(t-s)\} (t-s) dN(s) \end{aligned}$$

where  $a_{l, \tau} = a^{(2)} + \tau(a_l^{(1)} - a_l^{(2)})$ . Now note that  $e^{-x}$  is a decreasing function of  $x \in \mathbb{R}$  and  $e^{-x} x \leq C e^{-(1-\epsilon)x}$  for any  $\epsilon > 0$  and some fixed  $C < \infty$ . Hence, For any  $a < \underline{a}$ , the r.h.s. is bounded above by a constant multiple of

$$\left( \left| d_l^{(1)} - d_l^{(2)} \right| + \left| a_l^{(1)} - a_l^{(2)} \right| d_l^{(2)} \right) \int_{(-\infty, t)} e^{-a(t-s)} dN(s) \leq \frac{\epsilon}{L} \int_{(-\infty, t)} e^{-a(t-s)} dN(s).$$

Inserting this bound in (A.4) gives the first statement in the lemma.

For any  $r < t$ , define the process  $Y(r, t) := \left[ \int_{(r, t)} e^{-a(t-s)} dN(s) \right]^2$  which is a measurable function of  $N$ . We claim that  $\mathbb{E}|Y(-\infty, t)|^2$  and  $\mathbb{E}\lambda_0^2(t)$  are finite. Moreover, both quantities are stationary by Theorem 1. Then, Lemma 2 in Ogata (1978) says that  $\frac{1}{T} \int_0^T Y(-\infty, t) \lambda_0(t) dt$  satisfies the ergodic theorem and converges in probability to a finite constant, say  $C_a$ . This would prove the second statement in the lemma. Hence, we need to show that the claims relating to the finite moments are true.

We show that  $\mathbb{E}|Y(0, t)|^2 < \infty$  and  $\mathbb{E}\lambda_0^2(t) < \infty$ , uniformly in  $t \geq 0$ . By stationarity of  $N$  and monotonicity this implies that  $\sup_{t>0} \mathbb{E}|Y(-\infty, t)|^2 < \infty$ . Select a constant  $\Delta > 0$  and partition the interval  $(0, t]$  into  $n(t)$  subintervals  $\Delta_i := ((i-1)\Delta, i\Delta]$  of size  $\Delta$  with the  $n(t)$  interval equal to  $\Delta_{n(t)} := ((n(t)-1)\Delta, t]$ . Define  $Y_i = \left[ \int_{\Delta_i} e^{-a(t-s)} dN(s) \right]^2$ ; for ease of notation, in  $Y_i$ , we are suppressing the dependence on  $t$ . Note that  $Y_i \leq e^{-2a(t-\Delta_i)} N^2(\Delta_i)$  if  $i < n(t)$  and  $Y_{n(t)} \leq N^2(\Delta_{n(t)})$ . As in the main text, with some abuse of notation,  $N(\Delta_i) := \int_{\Delta_i} dN(t)$ . Hence,  $\mathbb{E}(Y(-\infty, t))^4 = \mathbb{E}\left(\sum_{i=1}^{n(t)} Y_i\right)^4 \leq \left(1 + \sum_{i=1}^{n(t)-1} e^{-2a(t-\Delta_i)}\right)^4 \mathbb{E}N^4(\Delta_1)$  using stationarity. By Lemma 1,  $N(\Delta_1)$  has moments of all orders. Moreover, it is not difficult to see that the sum in the parenthesis on the r.h.s. of the inequality is finite uniformly in  $t$ . Hence, the first claim follows. Again, by Lemma 1,  $\mathbb{E}\lambda_0^2(t) < \infty$ , which is the last claim we needed to prove. Hence, the proof of the present lemma is concluded. ■

We can then show the uniform convergence of the martingale process.

**Lemma 3** *Under the Regularity Conditions,*

$$\mathbb{E} \sup_{h \in \mathcal{H}} \max_{k \leq K} \left| \int_0^T \frac{X_k}{h} dM \right| \lesssim \sqrt{T \ln(1+K)}.$$

**Proof.** This is an application of Point (ii) in Theorem 2.5 in Nishiyama (2000). It is based on a number of conditions that we shall verify. For each  $\epsilon \in (0, \delta]$ , with  $\delta > 0$  to be defined momentarily, we need to find a partition  $\Pi(\epsilon) = \bigcup_{l=1}^{N(\epsilon, \Psi)} \Psi(\epsilon, l)$  of  $\Psi = \mathcal{C} \cup \mathcal{A} \cup \mathcal{D}$  where  $N(\epsilon, \Psi)$  is an integer valued function increasing in  $\epsilon$ , and such that  $N(\delta, \Psi) = 1$ ;  $N(\epsilon, \Psi)$  should not be confused with the point process  $N$ . We also need to show that for constants  $C_1$  and  $C_2$  with probability going to one, we have that

$$\int_0^T \max_{h \in \mathcal{H}} \max_{k \leq K} \left( \frac{X_k}{h} \right)^2 \lambda_0 d\mu \leq C_1 T \tag{A.5}$$

and

$$\sup_{\epsilon \in (0, \delta]} \max_{1 \leq l \leq N(\epsilon)} \int_0^T \max_{\psi, \phi \in \Psi(\epsilon, l)} \max_{k \leq K} \frac{|X_k|^2}{\epsilon^2} \left( \frac{1}{h_\psi} - \frac{1}{h_\phi} \right)^2 \lambda_0 d\mu \leq C_2 T. \tag{A.6}$$

At first, we show (A.5). To this end, note that  $|X_k|_\infty \leq 1$ , and that  $h$  is bounded below by  $c_{\mathcal{H}}$ . Then, the l.h.s. of (A.5) is bounded above by  $\int_0^T c_{\mathcal{H}}^{-2} \lambda_0 d\mu \leq 2c_{\mathcal{H}}^{-2} P \lambda_0 T$  with probability going to one; the factor 2 can be reduced to any arbitrary number greater than one. Hence,  $C_1 = 2c_{\mathcal{H}}^{-2} P \lambda_0$ .

We now define the partition  $\Pi(\epsilon)$  and use it to verify (A.6). For each  $\epsilon > 0$  we define a partition of the parameter space  $\mathcal{C} = [\underline{c}, \bar{c}] = \bigcup_{l=1}^{N(\epsilon, \mathcal{C})} \mathcal{C}_l$ ,  $\mathcal{D} = [0, \bar{d}]^L = \bigcup_{l=1}^{N(\epsilon, \mathcal{D})} \mathcal{D}_l$ , and  $\mathcal{A} = [\underline{a}, \bar{a}]^L = \bigcup_{l=1}^{N(\epsilon, \mathcal{A})} \mathcal{A}_l$  such that  $|c^{(1)} - c^{(2)}| \leq \epsilon$ ,  $\left\{ \left| d_l^{(1)} - d_l^{(2)} \right| \leq \frac{\epsilon}{2L} : l = 1, 2, \dots, L \right\}$ ,  $\left\{ \left| a_l^{(1)} - a_l^{(2)} \right| \leq \frac{\epsilon}{2dL} : l = 1, 2, \dots, L \right\}$  for any parameters belonging to the same partition of  $\mathcal{C}$ ,  $\mathcal{D}$  and  $\mathcal{A}$ , respectively. To avoid trivialities in the notation, we are supposing that for each

$l = 1, 2, \dots, L$ , the parameters are constrained in the same interval, e.g.  $a_l \in [\underline{a}, \bar{a}]$ . By the same steps used to establish (A.5), the l.h.s. of (A.6) is bounded above by

$$\sup_{\epsilon \in (0, \delta]} \max_{1 \leq l \leq N(\epsilon)} \int_0^T \max_{\psi, \phi \in \Psi(\epsilon, l)} \frac{1}{\epsilon^2} \left( \frac{h_\psi - h_\phi}{c_{\mathcal{H}}^2} \right)^2 \lambda_0 d\mu.$$

By Lemma 2 we can choose  $C_2 \asymp c_{\mathcal{H}}^{-4} (1 + C_a)$ , where  $C_a$  is as defined there. The partition we use for  $\mathcal{H}$  has cardinality  $N(\epsilon, \mathcal{H}) = N(\epsilon, \mathcal{C}) N(\epsilon, \mathcal{D}) N(\epsilon, \mathcal{A}) \leq \left(\frac{\bar{c}-\underline{c}}{\epsilon}\right) \left(\frac{2L\bar{d}}{\epsilon}\right)^L \left(\frac{2L(\bar{a}-\underline{a})}{\epsilon d}\right)^L \leq (\delta/\epsilon)^{1+2L}$ , where  $\delta = \max\{\bar{c} - \underline{c}, 2L\bar{d}, 2L(\bar{a} - \underline{a})/d\}$ . This partition needs to be multiplied by  $K$  because  $\{X_k : k = 1, 2, \dots, K\}$  is a family of processes with exactly  $K$  elements. Hence, Theorem 2.5 Point (ii) in Nishiyama (2000) says that

$$\mathbb{E} \sup_{h \in \mathcal{H}} \max_{k \leq K} \left| \int_0^T \frac{X_k}{h} dM \right| \lesssim \sqrt{C_2 T} \int_0^\delta \sqrt{\ln \left( \max \left\{ e, K \left( \frac{\delta}{\epsilon} \right)^{1+2L} \right\} \right)} d\epsilon + \frac{C_1 T}{\delta \sqrt{C_2 T}},$$

where, for convenience, we slightly changed the definition of the entropy function and we also simplified the statement of Theorem 2.5 Point (ii) in Nishiyama (2000) because (A.5) and (A.6) hold with probability going to one. The integral is bounded above by a constant multiple of  $\delta \sqrt{\ln(1+K)}$  and the result follows because  $\delta$  is fixed. ■

We can now prove, one by one the points in the statement of Theorem 3.

**Proof of Point 1.** This follows from (A.1), (A.2) and Lemma 3 noting that  $h \geq c_{\mathcal{H}} > 0$ .

**Proof of Point 2.** The argument is standard (van der Vaart and Wellner, 2000, proof of Theorem 3.2.5). For any  $g, g' \in \mathcal{G}$  and  $h \in \mathcal{H}$  define  $\Delta(g, g^h; h) = \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} - g \right)^2 d\mu - \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} - g^h \right)^2 d\mu$ , and  $d^2(g, g') = \frac{1}{T} \int_0^T (g - g')^2 d\mu$ . Note that  $\Delta(g, g^h; h) \geq \frac{1}{4} d^2(g, g^h)$  if  $d^2(g, g^h) \geq 4 \frac{1}{T} \int_0^T \left( \frac{\lambda_0}{h} - g^h \right)^2 d\mu$  (van der Vaart and Wellner, 2000, Problem 3.4.5). Clearly, if this is not the case,  $d^2(g, g^h) \leq 4A_T$  and there is nothing more to prove. Assuming that this is the case, given that  $\hat{g}^h$  minimizes  $R_T(g; h)$ , the event  $d^2(\hat{g}^h, g^h) > \epsilon$  is contained in the event

$$\sup_{g \in \mathcal{G}: d^2(g, g^h) > \epsilon} \left[ R_T(g^h; h) - R_T(g; h) \right] \geq 0.$$

Adding and subtracting  $\Delta(g, g^h; h)$  and using the lower bound in terms of  $d^2(g, g^h)$  deduce that the event in the above display is contained in the event

$$\sup_{d^2(g, g^h) > \epsilon} \left| R_T(g^h; h) - R_T(g; h) - \Delta(g, g^h; h) \right| \geq \epsilon/4.$$

By (A.1),

$$R_T(g; h) - R_T(g^h; h) - \Delta(g, g^h; h) = -\frac{2}{T} \int_0^T \frac{g - g^h}{h} dM.$$

From (A.2) and Lemma 3 deduce that there is a finite constant  $C$  such that the r.h.s. is less than  $C \times B\sqrt{T^{-1} \ln K}$  with probability going to one. Hence, choosing  $\epsilon = 4C \times B\sqrt{T^{-1} \ln K}$  we deduce that the probability of the event  $\{d^2(\hat{g}^h, g^h) > \epsilon\}$  goes to zero when  $\Delta(\hat{g}^h, g^h; h) \geq \frac{1}{4}d^2(\hat{g}^h, g^h)$  and this concludes the proof of Point 2.

**Proof of Point 3.** Using the previously defined notation and the fact that  $g_0 = \lambda_0/h_0$ , we have that  $d^2(g^{\hat{h}}, g_0) = \frac{1}{T} \int_0^T \left(\frac{\lambda_0}{\hat{h}} - g^{\hat{h}}\right)^2$ . Adding and subtracting  $\lambda_0/\hat{h}$  inside the square on the r.h.s., and using a simple inequality, we have that

$$d^2(g^{\hat{h}}, g_0) \leq \frac{2}{T} \int_0^T \left(\frac{\lambda_0}{\hat{h}} - g^{\hat{h}}\right)^2 d\mu + \frac{2}{T} \int_0^T (\hat{h} - h_0)^2 \left(\frac{\lambda_0}{\hat{h}h_0}\right)^2 d\mu. \quad (\text{A.7})$$

Similarly we deduce that

$$\frac{1}{T} \int_0^T \left(\frac{\lambda_0}{\hat{h}} - g\right)^2 d\mu \leq \frac{2}{T} \int_0^T \left(\frac{\lambda_0}{h_0} - g\right)^2 d\mu + \frac{2}{T} \int_0^T (\hat{h} - h_0)^2 \left(\frac{\lambda_0}{\hat{h}h_0}\right)^2 d\mu. \quad (\text{A.8})$$

Taking  $\inf_{g \in \mathcal{G}}$  on both sides and by definition of  $g^{\hat{h}}$ , the above is equal to

$$\frac{1}{T} \int_0^T \left(\frac{\lambda_0}{\hat{h}} - g^{\hat{h}}\right)^2 \leq \frac{2}{T} \int_0^T (\hat{h} - h_0)^2 \left(\frac{\lambda_0}{\hat{h}h_0}\right)^2 d\mu, \quad (\text{A.9})$$

where we have used the fact that  $g_0 = \lambda_0/h_0 \in \mathcal{G}$  and that the inf in the first term on the r.h.s. of (A.8) is attained at  $g_0$ . Inserting the above display in (A.7) and recalling that  $\hat{h} \geq c_{\mathcal{H}} > 0$ ,

$$d^2(g^{\hat{h}}, g_0) \leq \frac{4}{T} \int_0^T (\hat{h} - h_0)^2 \left(\frac{g_0}{c_{\mathcal{H}}}\right)^2 d\mu = O_P\left(B\sqrt{\frac{\ln(1+K)}{T}}\right),$$

where the r.h.s. follows by assumption. Incidentally, using (A.9) we have also shown that  $\frac{4}{T} \int_0^T \left(\frac{\lambda_0}{\hat{h}} - g^{\hat{h}}\right)^2 d\mu = O_P\left(B\sqrt{\frac{\ln(1+K)}{T}}\right)$  so that we can just apply the result from Point 2 with  $A_T = O_P\left(B\sqrt{\frac{\ln(1+K)}{T}}\right)$  and deduce Point 3 from the triangle inequality:  $d(\hat{g}^h, g_0) \leq d(\hat{g}^h, g^h) + d(g^h, g_0)$ .



### A.2.3 Proof of Theorem 2

For notational simplicity it is tacitly assumed that  $g$  and  $h$  are in  $\mathcal{G}$  and  $\mathcal{H}$ . Subtracting  $Q_T(h_0, g_0)$  from  $Q_T(h, g)$ ,

$$Q_T(h, g) - Q_T(h_0, g_0) = \frac{1}{T} \int_0^T (\lambda_0 - hg)^2 d\mu - \frac{2}{T} \int_0^T hgdM,$$

where  $dM(t) = dN(t) - \lambda_0(t)dt$ . By similar arguments as for Lemma 3, we have that

$$\mathbb{E} \sup_{h \in \mathcal{H}} \max_{k \leq K} \left| \int_0^T hX_k dM \right| \lesssim \sqrt{T \ln(1+K)}.$$

This can be established noting that the analogue of (A.5) to the present case follows because  $hX_k \lesssim \int_{(-\infty, t)} e^{-a(t-s)} dN(s)$  for any  $a < \underline{a}$  as Lemma 2. The r.h.s. has finite expectation, as shown in the proof of Lemma 2. The analogue of (A.6) also follows using the arguments in the proof of Lemma 3. In consequence, the minimizers of  $Q_T(h, g)$  w.r.t.  $h$  and  $g$  minimize asymptotically  $\frac{1}{T} \int_0^T (\lambda_0 - hg)^2 d\mu$ . By the Regularity Conditions, and similar arguments as in the proof of Point 2 in Theorem 3, we also deduce the consistency rates of the estimator.

### A.2.4 Proof of Theorem 4

Using the definition of  $\epsilon_T$  we have that

$$\frac{L_T^{(1)} - L_T^{(2)}}{\sqrt{T\hat{\sigma}_T^2}} = \frac{\int_0^T \left( \ln(\hat{\lambda}^{(1)}) - \ln(\hat{\lambda}^{(2)}) \right) dM}{\sqrt{T\hat{\sigma}_T^2}} + \frac{\epsilon_T}{\sqrt{T\hat{\sigma}_T^2}}$$

where  $dM(t) = dN(t) - \lambda_0(t)dt$ . The second term on the r.h.s. is  $o_p(1)$  and can be disregarded. This is only true if  $\hat{\sigma}_T^2 > 0$  with probability going to one. By the assumptions in the statement of the theorem, it is not difficult to see that this is the case using ergodicity of  $N$ . Let  $\Delta_i = ((i-1)P\lambda_0, iP\lambda_0]$ ,  $i = 1, 2, \dots, n$  and to avoid trivialities suppose that  $nP\lambda_0 = T$ . We also note that  $P\lambda_0 = |\Delta_i|$ , where the r.h.s. is the Lebesgue measure of  $\Delta_i$ . Then,

$$\frac{1}{\sqrt{T}} \int_0^T \left( \ln(\hat{\lambda}^{(1)}) - \ln(\hat{\lambda}^{(2)}) \right) dM = \frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i \quad (\text{A.10})$$

where  $Y_i := |\Delta_i|^{-1/2} \int_{\Delta_i} \ln\left(\frac{\hat{\lambda}^{(1)}}{\hat{\lambda}^{(2)}}\right) dM$ . By construction  $\mathbb{E}_{i-1} Y_i = 0$ , where  $\mathbb{E}_{i-1}$  is expectation conditioning on  $(Y_j)_{j \leq i-1}$ . Given that  $(Y_i)_{i \geq 1}$  is stationary and ergodic by Theorem 1, for the martingale central limit theorem to apply to (A.10), it is sufficient that  $\mathbb{E}|Y_i|^2 < \infty$ . Let  $f = \ln\left(\frac{\hat{\lambda}^{(1)}}{\hat{\lambda}^{(2)}}\right)$  so that  $\mathbb{E}|Y_i|^2 = |\Delta_i|^{-1} \mathbb{E} \left( \int_{\Delta_i} f dM \right)^2$ . The r.h.s. is equal to  $|\Delta_i|^{-1} \mathbb{E} \int_{\Delta_i} f^2 \lambda_0 d\mu$  by standard isometry. By stationarity, this is  $Pf^2\lambda_0 < \infty$ . To see this, use Holder inequality

$Pf^2\lambda_0 \leq Pf^4P|\lambda_0|^2$  and note that  $Pf^4 \lesssim P\hat{\lambda}^{(1)} + P\hat{\lambda}^{(2)}$  using a simple bound on the logarithm and the fact that the intensities  $\hat{\lambda}^{(k)}$  are bounded away from zero. Using the same arguments as at the end of the proof of Lemma 2, it is easy to deduce that  $P\hat{\lambda}^{(k)} < \infty$ . This shows that (A.10) converges to a normal random variable with mean zero and variance  $\mathbb{E}Y_i^2 = \mathbb{E}f^2\lambda_0$ . Then,  $\frac{1}{\sqrt{n}} \sum_{i=1}^n Y_i / \sqrt{\mathbb{E}f^2\lambda_0}$  has variance one, hence it is asymptotically standard normal in distribution. By ergodicity of the counting process,  $\mathbb{E}f^2\lambda_0 = \hat{\sigma}_T^2 + o_p(1)$  so that the assertion in (15) is proved.

### A.3 Empirical Study: Additional Details

We first discuss some of the nuisances and challenges of the dataset we use. We then define the bins used for the construction of the one-hot encoding of each variable. We conclude with details regarding the parameters restrictions used in the estimation.

#### A.3.1 Data Conflation and Limits of Our Dataset

Exchanges broadcast a number of information in the same packet via the electronic communication network. This information is received and consumed by any algorithm that subscribes to a given feed on the network. The exchange timestamp does not allow to know what was sent in a packet, as this is often a function of network capacity. Hence, it is not possible to know how the information was received. Moreover, the exchange timestamp and order id does not allow us to synchronize data across multiple instruments with certainty. Historically, this has been a clear problem with Level 2 data. However, we are not aware of the extent to which such problem is mitigated by the use of Level 3 data on the NYSE. In depth knowledge usually requires analysis of properly collected data in live trading. Here, we made the assumption that all orders with the same exchange time stamp are sent and consumed by an algorithm at the same time. We also synchronized different stocks on the basis of their exchange timestamp. Finally, our data were queried for 10 levels. It is not clear how this could be done in live trading. In practice, the algorithm subscribing to a Level 3 feed is flooded by all messages and conflation is necessary. This conflation may differ substantially from what has been done in the present paper.

No study that uses purchased data can escape the above problems. However, to reassure the reader, we note that in a previous version of this paper we carried out similar analysis using CME Level 2 data collected by a proprietary trading firm colocated in the Aurora data centre in Chicago. This data had nanosecond arrival timestamp (time at which the data would have been seen by an algorithm) together with an identifier for the order in which the information could be found within the same packet. This allowed for perfect synchronization across products and knowledge of the end of packet time stamp. It is reassuring that with such

dataset, we obtained results consistent with the data used in the present version of the paper. We do not report the results on this high quality proprietary CME Level 2 data because it dates back to 2014 and market microstructure might have changed in the meantime.

### A.3.2 Parameters' Restrictions for Estimation

We impose the restriction that the linear coefficients  $b_k$  are in  $[0, 10\beta]$  where  $\beta$  is as in Section 11. This implies that  $B = 10K\beta$ . We multiply by 10 to avoid skinking coefficients too much. For the models that are linear in the raw covariates we allow the linear coefficients  $b_k$  to be in  $[-10\beta, 10\beta]$ , as otherwise we cannot capture negative impact. In this case, the intensity is not guaranteed to be nonnegative. Hence, when testing, we impose a lower bound on the intensity equal to  $\text{eps}$  for the models linear in the raw covariates;  $\text{eps}$  is machine epsilon.

Let  $q_{\text{Dur},10}$  and  $q_{\text{Dur},50}$  be the 10% and 50% quantile of the trades durations. The parameters in the estimation of  $h$  in (2) are restricted as follows,  $c \in [10^{-3}, 10] / q_{\text{Dur},50}$ ,  $d_l \in [0, 1] / q_{\text{Dur},50}$  and  $a_l \in [10^{-2}, 10] / q_{\text{Dur},10}$ . Scaling by  $q_{\text{Dur},10}$  and  $q_{\text{Dur},50}$  ensures that we keep the correct order of magnitude irrespective of how time is measured. Here, time is measured in seconds with nine decimal places and for example  $q_{\text{Dur},10}$  tends to be in the order of  $10^{-4}$ .