

# Hawkes Processes Modeling, Inference, and Control: An Overview\*

Rafael Lima<sup>†</sup>

**Abstract.** Hawkes processes are a type of point process that models self-excitement among time events. They have been used in a myriad of applications, ranging from finance and earthquakes to crime rates and social network activity analysis. Recently, a variety of different tools and algorithms have been presented at top-tier machine learning conferences. This work aims to give a broad view of recent advances in Hawkes process modeling and inference suitable for a newcomer to the field. The parametric, nonparametric, deep learning, and reinforcement learning approaches are broadly discussed, along with the current research challenges for the topic and the real-world limitations of each approach. Illustrative application examples in the modeling of retweeting behavior, earthquake aftershock occurrence, and malaria outbreak modeling are also briefly discussed.

**Key words.** Hawkes processes, point processes, machine learning

**AMS subject classifications.** 68T99, 62M20, 60G55

**DOI.** 10.1137/21M1396927

---

## Contents

<b>1</b>	<b>Introduction</b>	<b>332</b>
<b>2</b>	<b>Theoretical Background</b>	<b>333</b>
2.1	Multivariate Marked Temporal Point Processes . . . . .	334
2.2	Hawkes Processes . . . . .	335
2.3	Spatiotemporal HPs . . . . .	338
2.4	Simulation Algorithms . . . . .	339
<b>3</b>	<b>Parametric HPs</b>	<b>339</b>
3.1	Enhanced and Composite Triggering Kernels . . . . .	340
3.2	Scalability . . . . .	342
3.3	Training . . . . .	344
<b>4</b>	<b>Nonparametric HPs</b>	<b>347</b>
4.1	Frequentist Nonparametric HPs . . . . .	347
4.2	Acceleration of Impact Matrix Estimation through Matching of Cumulants . . . . .	347
4.3	Online Learning . . . . .	349
4.4	Bayesian Nonparametric HPs . . . . .	350

---

\*Received by the editors February 4, 2021; accepted for publication (in revised form) January 25, 2022; published electronically May 9, 2023. Any opinions, findings, and conclusions expressed in this work are those of the author and do not necessarily reflect the views of Samsung R&D Institute Brazil.

<https://doi.org/10.1137/21M1396927>

<sup>†</sup>Samsung R&D Institute Brazil, Campinas, 13083-730, Brazil (rafael.goncalves.lima@gmail.com).

<b>5</b>	<b>Neural Network Based HPs</b>	<b>352</b>
5.1	Self-Attentive and Transformer Models . . . . .	358
5.2	Graph Convolutional Networks . . . . .	358
<b>6</b>	<b>Further Approaches</b>	<b>359</b>
6.1	Sparse Gaussian Processes . . . . .	359
6.2	Stochastic Differential Equation . . . . .	360
6.3	Graph Properties . . . . .	360
6.4	Epidemic HPs . . . . .	360
6.5	Popularity Prediction . . . . .	360
<b>7</b>	<b>Stochastic Control and Reinforcement Learning of HPs</b>	<b>360</b>
7.1	Stochastic Optimal Control (SOC) . . . . .	361
7.2	Reinforcement Learning . . . . .	363
7.3	Imitation Learning . . . . .	364
<b>8</b>	<b>Real-World Data Limitations</b>	<b>364</b>
8.1	Synchronization Noise . . . . .	364
8.2	Sequences with Few Events . . . . .	365
8.3	Sequences with Missing Data . . . . .	365
<b>9</b>	<b>Application Examples</b>	<b>365</b>
9.1	Retweet . . . . .	365
9.2	Earthquake Aftershocks . . . . .	366
9.3	Malaria Outbreak Forecasting . . . . .	366
9.4	Financial Modeling . . . . .	367
<b>10</b>	<b>Comparisons with Other Temporal Point Process Approaches</b>	<b>367</b>
<b>11</b>	<b>Current Challenges for Further Research</b>	<b>367</b>
<b>12</b>	<b>Conclusions</b>	<b>368</b>
	<b>Acknowledgments</b>	<b>368</b>
	<b>References</b>	<b>368</b>

**1. Introduction.** Point processes are tools for modeling the arrival of time events. They have been broadly used to model both natural and social phenomena related to the arrival of events in a continuous-time setting, such as the queuing of customers in a given store, the arrival of earthquake aftershocks [63, 32], the failure of machines at a factory, the request for packages over a communication network, and the death of citizens in ancient societies [17].

Predicting, and thus being able to effectively intervene in, all these phenomena is of huge commercial and/or societal value, and thus there has been intensive investigation of the theoretical foundations of this field.

Hawkes processes (HPs) [29] are a type of point process that models self- and mutual excitation, i.e., when the arrival of an event makes future events more likely to happen. They are suitable for capturing epidemic, clustering, and faddish behavior

within social and natural time-varying phenomena. The excitation effect is represented by an additional function term within the intensity of the process (i.e., the expected arrival rate of events): the triggering kernel, which quantifies the influence of events of a given process in the self- and mutual triggering of its associated intensity functions. Much of HP research has been devoted to modeling the triggering kernels, handling issues of scalability to a large number of concurrent processes and large quantities of data, and addressing speed and tractability of the inference procedure.

Regarding the modeling of the triggering kernels, one method involves the assumption that they can be defined by simple parametric functions, such as one or multiple exponentials, Gaussians, Rayleigh, Mittag-Leffler [12] functions, and power laws. Much of the work dealing with this type of approach is concerned with enriching these parametric models [42, 84, 91, 90, 20], scaling them for high dimensions, i.e., multivariate processes of large dimensions [7, 44], dealing with distortions related to restrictions on the type of available data [92], and proposing adversarial losses as a complement to the simple maximum likelihood estimation (MLE) [94].

Another way of modeling the triggering function is by assuming that it is represented by a finite grid, in which the triggering remains constant along each of its subintervals. In this piecewise constant (or nonparametric) approach, most notably developed in [59, 8], the focus has been on speeding up the inference through parallelization and online learning of the model parameters [96, 2].

A more recent approach, enabled by the rise to prominence of deep learning models and techniques, accompanied by the increase in computational power and availability of data over recent years, involves modeling the causal triggering effect through the use of neural network models, most notably RNNs, LSTMs, and GANs. These models allow for less bias and more flexibility than the parametric models in modeling the triggering kernel, while taking advantage of the numerous training and modeling techniques developed by the booming connectionist community.

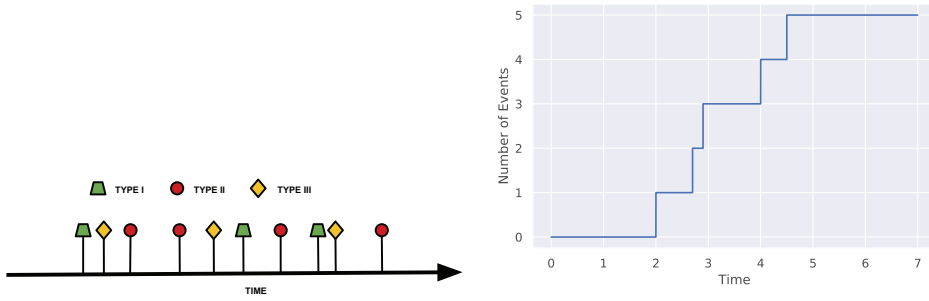
In addition, regarding the control of self-exciting point processes, i.e., the modification of the process parameters into more desirable configurations while taking into account an associated “control cost” to the magnitude of these modifications, recent works make use of either dynamic programming (continuous Hamilton–Jacobi–Bellman equation based approach) [99, 98], Kullback–Leibler divergence penalization (a.k.a. “information bottleneck”) [83], or reinforcement learning based or imitation learning based techniques [80, 47].

Although there have been some interesting reviews and tutorials regarding domain-specific applications of HPs in finance [30, 7] and social networks [69], a broad view of the inference and modeling approaches is still lacking. A work similar to ours is [93], which, although very insightful, lacks coverage of important advances such as the previously mentioned control approaches, as well as the richer variants of neural network based models. Furthermore, a concise coverage of the broader class of temporal point processes is given in [71], while reviews for parametric spatiotemporal formulations of the HP are given in [67] and [97].

In what follows, we introduce the mathematical definitions involving HPs, then carefully describe advances in each of the aforementioned approaches, and then finish with a summary, along with some other considerations.

**2. Theoretical Background.** In this and the next section, we present the mathematical definitions used throughout the remaining sections of the paper. HPs considered were originally introduced in [29] and [28].

In the present work, we restrain ourselves to the marked temporal point process



**Fig. 1** Left: Intuitive diagram of a marked temporal point process (MTPP) with three types of events (marks). Right: Example of a counting process on the time interval  $[0, 7]$ .

(MTPP), i.e., the point process in which each event is defined by a time coordinate and a mark (or label). An intuitive example of an MTPP is shown in Figure 1.

The key definitions below are those of counting process, intensity function, triggering kernel, impact matrix, (log-)likelihood, covariance, Bartlett spectrum, higher-order moments, and branching structure.

**2.1. Multivariate Marked Temporal Point Processes.** Realizations of univariate MTPPs, here referred to by  $\mathcal{S}$ , are one or more sequences of events  $e_i$ , each a function of the time coordinate  $t_k$  and the mark  $m_k$ , such as

$$(2.1) \quad \mathcal{S} = \{(t_0, m_0), \dots, (t_S, m_S)\},$$

where  $S$  is the total number of events. Marks may represent, for example, a specific user in a social network or a specific geographic location. For more complex problems, such as in the check-in times prediction of [95], a composite mark may represent a user of interest and a specific location.

An easy way to generalize this notation would be to refer to multiple realizations of multivariate MTPPs as  $\mathcal{S} = \{\mathcal{S}_{i,j}\}$ , where  $i$  refers to the dimension of the process, while  $j$  refers to the index of the sequence. Now, regarding only the purely temporal portion of the process, i.e., the time coordinates  $t_k$ , it is also common to express them by means of a counting process  $N(t)$ , which is simply the cumulative number of event arrivals up to time  $t$ :

$$(2.2) \quad \int_{0^-}^t dN_s,$$

where

- $dN_{t_k} = 1$  if there is an event at  $t_k$ ;
- $dN_{t_k} = 0$  otherwise.

This is illustrated in Figure 1.

Associated to each temporal point process, there is an intensity function, which is the expected rate of arrival of events,

$$(2.3) \quad \lambda(t)dt = E \{dN_t = 1\},$$

which may or may not depend on the history of past events. Such dependence results in a so-called conditional intensity function (CIF),

$$(2.4) \quad \lambda(t)dt = E \{dN_t = 1|\mathcal{H}\},$$

where  $\mathcal{H}$  is the history of all events up to time  $t$ :

$$(2.5) \quad \mathcal{H} : \{t_{i,j} \in \mathcal{S} | t_{i,j} < t\}.$$

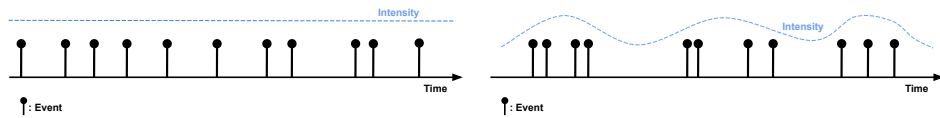
This concept will be further discussed in the next section.

**2.2. Hawkes Processes.** The simplest example of a temporal point process is the homogeneous Poisson process (HPP), in which the intensity is a positive constant:

$$(2.6) \quad \lambda(t) = \mu$$

for  $\mu \in \mathbb{R}^+$ .

In the case of an inhomogeneous Poisson process (IPP), the intensity  $\lambda(t)$  is allowed to vary. Both the HPP and the IPP are shown in Figure 2.



**Fig. 2** Illustrative examples of HPPs (left) and IPPs (right). The events are represented by black dots.

HPPs and IPPs have in common the fact that each consecutive event interval sampled from the intensity function is independent of the previous ones. When analyzing several natural phenomena, one may wish to model how events in each dimension  $i$  of the process—which may be representing a specific social network user, an earthquake shock at a given geographical region, or a percentage jump in the price of a given stock, just to cite a few examples—affect the arrival of events in all the dimensions of the process, including its own.

In particular, we are interested in the cases where the arrival of one event makes further events more likely to happen, which is reflected as an increase in the value of the intensity function after the time of that event. When this increase happens in the intensity function of the same  $i$ th dimension as the event, the effect is called self-excitation. When the increase happens in the intensity of other dimensions, we refer to it as mutual excitation.

HPs model self-excitation in an analytical expression for the intensity through the insertion of an extra term that is designed to capture the effect of all the previous events of the process in the current value of the CIF. For a univariate HP with constant background rate, we have

$$(2.7) \quad \lambda_{HP}(t) = \underbrace{\mu}_{\text{baseline intensity}} + \underbrace{\sum_{t_i < t} \phi(t - t_i)}_{\text{self-excitation term}},$$

while, for the multivariate case with dimension  $D$ , we have both self-excitation ( $\phi_{ii}(t)$ ) and mutual excitation terms ( $\phi_{ij}(t)$  s.t.  $(i \neq j)$ ):

$$(2.8) \quad \lambda_{HP}^i(t) = \mu_i + \sum_{j=1}^D \sum_{t_{ij} < t} \phi_{ij}(t - t_{ij}).$$

The assumptions of

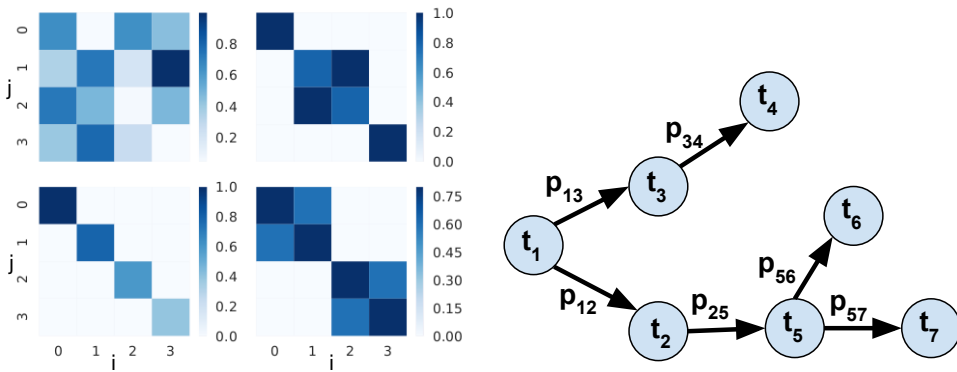
1. causality:  $\phi(t) = 0 \quad \forall t < 0;$
2. positivity:  $\phi(t) \geq 0 \quad \forall t \geq 0$

are usually held for all  $\phi_{ij}(t)$ .<sup>1</sup>

Some works consider a time-dependent background rate  $\mu(t)$ , to account for the seasonal aspect of some phenomena such as disease outbreaks or crime rates [24, 79, 52, 24]. This is further discussed in section 9.

Now consider the case that the kernel matrix  $\Phi(t) = [\phi_{ij}(t)]_{i,j=0}^{d,n}$  can be factored into  $\Phi(t) = \alpha \odot \kappa(t)$ , with  $\alpha = [\alpha_{ij}]_{i,j=0}^{d,n}$  and  $\kappa(t) = [\kappa_{ij}(t)]_{i,j=0}^{d,n}$ , where “ $\odot$ ” corresponds to the Hadamard (elementwise) product.

$\alpha$ , representing the impact matrix, can implicitly capture a myriad of different patterns of self- and mutual excitation, as exemplified in Figure 3. This factorization is particularly convenient when adding penalization terms related to network properties to the loglikelihood-based loss of a multivariate HP, such as in [91] and [51].



**Fig. 3** Left: Four examples of  $4 \times 4$  impact matrices  $\alpha$ . Each  $\alpha_{ij}(t)$  has the corresponding value indicated by the color scale on the side. Right: Illustrative example of the concept of branching structure. A given edge  $t_i \rightarrow t_j$  means that  $t_i$  triggered  $t_j$  with the probability  $p_{ji}$ .

To be of practical value, realizations of HPs are constrained to having a finite number of events for any subinterval of the simulation horizon  $[0, T]$ . This corresponds to the kernel function (or kernel matrix) satisfying the following stationarity condition:

$$(2.9) \quad \mathbf{Spr}(\|\Phi(t)\|) = \mathbf{Spr}(\{\|\phi_{ij}(t)\|\}_{1 \leq i, j \leq D}) < 1,$$

where  $\|\phi(t)\|$  corresponds to  $\|\phi(t)\| = \int_0^\infty \phi(t)dt$  and  $\mathbf{Spr}(\cdot)$  corresponds to the spectral radius of the matrix, i.e., the largest value among its eigenvalues.

If this stationarity condition is satisfied, the process will reach a weakly stationary state, i.e., when the properties of the process, most notably its “moments,” vary only as a function of the relative distance, here referred to as “ $\tau$ ”, of its points.

The first-order moment, or statistics of the HP, is defined as

$$(2.10) \quad \Lambda_i = E\{\lambda_{HP}^i(t)\} = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \lambda_{HP}^i(t)dt = (\mathbb{I} - \|\Phi(t)\|)^{-1} \mu_i,$$

<sup>1</sup>The case in which  $\phi(t) < 0$  for  $t \geq 0$  is referred to as an “inhibiting process” and is not usually considered in HP work.

while the second-order statistics, or stationary covariance, is defined as

$$(2.11) \quad \nu^{ij}(t' - t)dt dt' = E\{dN_t^i dN_{t'}^j\} - \Lambda_i \Lambda_j dt dt' - \epsilon_{ij} \Lambda_i \delta(t' - t) dt,$$

where  $\epsilon_{ij}$  is 1 if  $i = j$ , and 0 otherwise, while  $\delta(t)$  refers to the Dirac delta distribution.

The Fourier transform of this stationary covariance is referred to as the *Bartlett spectrum*. Sometimes a different transform, the Laplace transform, is used for the same purpose. In Hawkes’s seminal paper [29], high importance is given to the fact that, when assuming some specific parametric functions for the excitation matrix  $\Phi(t)$ , it is possible to find simple formulas for the covariance of the process in the frequency domain. One example is the univariate case for  $\phi(t)$  defined as the frequently used “parametric exponential kernel”:  $\phi(t) = \alpha e^{-\beta t}$  for  $\alpha, \beta \in \mathbb{R}$ .

For this choice, we have

$$(2.12) \quad \nu^*(s) = \mathcal{L}\{\nu\}(s) = \frac{\alpha\mu(2\beta - \alpha)}{2(\beta - \alpha)(s + \beta - \alpha)} \quad (s \in \mathbb{C}),$$

where  $\mathcal{L}\{\cdot\}(s)$  refers to the Laplace transform.<sup>2</sup> The detailed steps of this computation can be found in [29].

Going beyond the first- and second-order statistics, it is also possible to define statistics of higher orders; see [50, 17]. Although they become less and less intuitive and tractable as their order increases, [2] makes use of third-order statistics,  $K^{ijk}$ , in a specific application of the generalized method of moments for modeling the impact matrix of multivariate HPs. It is defined as

$$(2.13) \quad K^{ijk} dt = \int \int_{\tau, \tau' \in \mathbb{R}^2} \left( \mathbb{E}(dN_t^i dN_{t+\tau}^j dN_{t+\tau'}^k) \right. \\ \left. - 2\mathbb{E}(dN_t^i) \mathbb{E}(dN_{t+\tau}^j) \mathbb{E}(dN_{t+\tau'}^k) - \mathbb{E}(dN_t^i dN_{t+\tau}^j) \mathbb{E}(dN_{t+\tau'}^k) \right. \\ \left. - \mathbb{E}(dN_t^i dN_{t+\tau'}^k) \mathbb{E}(dN_{t+\tau}^j) - \mathbb{E}(dN_{t+\tau}^j dN_{t+\tau'}^k) \mathbb{E}(dN_t^i) \right)$$

for  $1 \leq i, j, k \leq D$ , and it is connected to the skewness of  $N_t$ .

Now, regardless of the function family chosen for modeling  $\Phi(t)$  and  $\mu$ , its fitness will be computed by measuring its likelihood over a set of sequences similar to the set of sequences used for training the model. Let  $\mathcal{S}$  be a set of  $M$  sequences, each with a total number of  $N_j$  events, considered over the interval  $[0, T]$ , such that

$$(2.14) \quad \mathcal{S} = \{\mathcal{S}^j\}_{j=1}^M = \left\{ \left[ (t_1^j, m_0^j), \dots, (t_{N_j}^j, m_{N_j}^j) \right] \right\}_{j=1}^M,$$

with  $m_k^j \in \{1, 2, \dots, D\} \forall j, k \in \mathbb{Z}_+$ .

Let  $\mathbb{F}$  be a family of multivariate HPs with dimension  $D \geq 1$  and parametric exponential kernels assumed for the shape of the excitation functions, such that the CIF  $\lambda_i^j(t)$  of the  $i$ th node defined over the sequence  $\mathcal{S}^j$  is given by<sup>3</sup>

$$(2.15) \quad \lambda_i^j(t) = \mu_i + \sum_{t_k^j \leq t} \alpha_{m_k^j i} e^{-\beta_{m_k^j i} (t - t_k^j)} \quad (t_k^j \in \mathcal{S}^j).$$

<sup>2</sup>The Laplace transform  $\mathcal{L}\{f\}(s)$  of a function  $f(t)$  defined for  $t \geq 0$  is computed as  $\mathcal{L}\{f\}(s) = \int_0^\infty f(t)e^{-st}$  for some  $s \in \mathbb{C}$ .

<sup>3</sup>Equivalent definitions of  $\mathbb{F}$  can be given to families of HPs defined by other types of HPs, such as those with power-law kernels, or those with the corresponding CIF modeled by a recurrent neural network.

Given parameter vectors

$$\boldsymbol{\mu} = \{\mu_m\}_{m=1}^D \in \mathbb{R}_+^D, \quad \boldsymbol{\theta} = \{(\alpha_{mn}, \beta_{mn})\}_{m=1, n=1}^{D, D} \in \mathbb{R}_+^{2D^2},$$

the likelihood function given in the logarithmic form, i.e., the loglikelihood, of a multivariate HP over a set  $\mathcal{S}$  of  $M$  sequences considered over the interval  $[0, T]$ , is given by

$$(2.16) \quad llh_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\theta}, \mathbb{F}) = \sum_{j=1}^M \left( \sum_{i=1}^D \sum_{k \in N_j} \log \lambda_i^j(t_k^j) - \underbrace{\sum_{i=1}^D \int_0^T \lambda_i^j(t) dt}_{\text{Compensator}} \right).$$

Thus, the goal of modeling an HP over  $\mathcal{S}$  is the act of finding vectors  $\boldsymbol{\mu}$  and  $\boldsymbol{\theta}$  such that

$$(2.17) \quad (\boldsymbol{\mu}, \boldsymbol{\theta}) = \operatorname{argmax} llh_{\mathcal{S}}(\boldsymbol{\mu}, \boldsymbol{\theta}, \mathbb{F}).$$

A more rigorous and complete derivation of (2.17) can be found in [9].

Another concept which is of relevance in some inference methods is that of a branching structure ( $\mathcal{B}$ ), which defines the ancestry of each event in a given sequence, i.e., specifies the probability that the  $i$ th event  $t_i$  was caused by the effect of a preceding event  $t_j$  in the CIF ( $p_{ji}$  for  $0 \leq i < n$ ) or by the baseline intensity  $\mu$  ( $p_{0i}$ ). The probabilities  $p_{ji}$  and  $p_{0i}$  can be given by

$$(2.18) \quad p_{ji} = \frac{\phi(t_i - t_j)}{\lambda(t_i)} \quad (\text{for } j \geq 1) \quad \text{and} \quad p_{0i} = \frac{\mu}{\lambda(t_i)}.$$

As an example, consider Figure 3, in which event  $t_1$ , the first of the event series, causes events  $t_2$  and  $t_3$ ; event  $t_2$  causes event  $t_5$ ; event  $t_3$  causes event  $t_4$ ; and event  $t_5$  causes events  $t_6$  and  $t_7$ . The corresponding branching structure  $\mathcal{B}_E$  implied by these relations among the events has an associated probability given by

$$(2.19) \quad p(\mathcal{B}_E) = p_{01} * p_{12} * p_{13} * p_{25} * p_{34} * p_{56} * p_{67}.$$

**2.3. Spatiotemporal HPs.** While the original HP formulation solely takes into account the temporal dependencies among consecutive phenomena, a line of work [24, 79, 52, 24] has focused on exploring additional spatially dependent self-excitation effects, which have been shown to be of importance in modeling crime rates, drug overdoses, and infectious diseases, among others.

The CIF of a spatiotemporal HP with time-invariant background rate is given by

$$(2.20) \quad \lambda_{STHP}(\dagger, t) = \underbrace{\mu(\dagger)}_{\text{baseline intensity}} + \underbrace{\sum_{t_i < t} \phi(\dagger - \dagger_i, t - t_i)}_{\text{self-excitation term}},$$

where  $\{\dagger_1, \dagger_2, \dots, \dagger_N\}$  denotes the history of locations of events and  $\{t_1, t_2, \dots, t_N\}$  again denotes the history of timestamps of these events.

For a comprehensive treatment of spatiotemporal self-exciting effects, the reader may refer to [30].



**2.4. Simulation Algorithms.** Regarding the experimental aspect of HPs, synthetic data may be generated through the following methods:

1. **Ogata's Modified Thinning Algorithm** [62]: This starts by sampling the first event at time  $t_0$  from the baseline intensity. Then each posterior event  $t_i$  is obtained by sampling it from an HPP with intensity fixed as the value calculated at  $t_{i-1}$ , and then either
  - accepting it with probability  $\frac{\lambda(t_i)}{\lambda^*(t_{i-1})}$ , where  $\lambda^*(t_{i-1})$  is the value of the intensity at time  $t_{i-1}$ , while  $\lambda_{t_i}$  is the value calculated through (2.7), or
  - rejecting it and proceeding to resample a posterior event candidate.
2. **Perfect Simulation** [60]: This derives from the fact that the HP may be seen as a superposition of Poisson processes. It proceeds by sampling events from the baseline intensity, taken as the initial level, and then sampling levels of descendant events for each of the events sampled at the previous level. To each event  $t_{0,i}$  sampled from the baseline intensity, we associate an IPP with intensity defined as  $\phi(t - t_{0,i})$ , and then sample its descendant events. Next, we take each sampled descendant event and associate it with its corresponding IPP, and so on, until all the levels have been explored over the simulation horizon  $[0, T]$ .

---

**Algorithm 2.1** Ogata's Modified Thinning Algorithm (Univariate Case)

---

```

Input  $\mu, \phi(t), T$ 
Define  $t = 0$ 
Sample  $t_1$  from exponential distribution with rate  $\mu$ 
Update  $t = t + t_1$ 
Define  $n = 1$ 
while  $t < T$  do
   $\lambda_n = \mu + \sum_{i=1}^n \phi(t - t_n)$ 
  Sample  $t_{n+1}$  from exponential distribution with rate  $\lambda_n$ 
   $\lambda_{n+1} = \mu + \sum_{i=1}^n \phi(t + t_{n+1} - t_n)$ 
  Sample  $u$  from uniform distribution over  $[0, 1]$ 
  if  $\frac{\lambda_{n+1}}{\lambda_n} < u$ , then
    Update  $t = t + t_{n+1}$ 
    Update  $n = n + 1$ 
  end if
end while
return  $\{t_i\}_{i=1}^n$ 

```

---

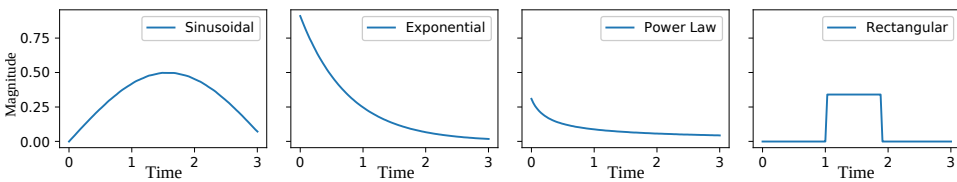
More detailed, step-by-step descriptions of each of these two algorithms are shown in pseudocode format, in Algorithms 2.1 and 2.2, for the case of univariate HPs. In the next section, we focus on the HP models that assume simple parametric forms for the excitation functions, along with their variants, which we will refer to as “parametric HPs.”

**3. Parametric HPs.** In this section, we discuss the HP models that assume simple parametric forms for the excitation functions. Figure 4 shows some examples of commonly used functions for parametric HPs. Much recent work concerns incrementing these models for dealing with specific aspects of certain domains, such as social networks [103, 42], audio streaming [84], and medical check-ups [92], among others.

**Algorithm 2.2** Perfect Simulation of Hawkes Processes (Univariate Case)

```

Input  $\mu, \phi(t), T$ 
Define  $j = 0$ 
Simulate HPP with  $\lambda = \mu$  over  $[0, T]$  to obtain  $\{t_j^i\}_{i=1}^{n_j}$ 
while  $\exists t_j^i (\forall i \leq n_j)$  do
  for  $(i = 1 ; i \leq n_j ; i++)$  do
    Simulate IPP with  $\lambda = \phi(t - t_j^i)$  over  $t \in [0, T]$  to obtain  $\{t_{(j+1),k}^i\}_{k=1}^{n_{j+1}^i}$ 
  end for
  Update  $n_{j+1} = \sum_{i=1}^{n_j} n_{j+1}^i$ 
  Update  $\{t_{(j+1)}^i\}_{i=1}^{n_{j+1}} = \bigcup_{i=1}^{n_j} \{t_{(j+1),k}^i\}_{k=1}^{n_{j+1}^i}$ 
  Update  $j = j + 1$ 
end while
return  $\bigcup_{i=0}^j \{t_i^i\}_{i=1}^{n_i}$  (After sorting)
  
```



**Fig. 4** Four examples of parametric HP kernels ( $\phi(t)$ ). Each of them is used to model a different type of interaction among events of a given HP.

In the following subsections, we give a broad view of how the parametric models of HPs are used and improved upon through a series of ideas. We have divided the focus of recent research on parametric HPs into three different strategies, accompanied by working examples. The three strategies are as follows:

1. Enhancing and composing simple parametric kernels to adapt the model to specific modeling situations and datasets.
  2. Improving the scalability of parametric HP models for multivariate cases with many nodes and sequences with many jumps.
  3. Improving robustness of training for worst-case scenarios and defective data.
- Further examples of each strategy are also briefly mentioned in section 11.

**3.1. Enhanced and Composite Triggering Kernels.** As a way of modeling the daily oscillations of the triggering effects on Twitter data, [42] proposes a time-varying excitation function for HPs. The probability  $\mathbb{P}$  of getting a retweet over the time interval  $[t, t + \delta t]$ , with small  $\delta$ , is modeled as

$$(3.1) \quad \mathbb{P}(\text{Retweet in } [t, t + \delta t]) = \lambda(t)\delta t,$$

in which the time-dependent rate is dependent on previous events as

$$(3.2) \quad \lambda(t) = p(t) \sum_{t_i < t} d_i \phi(t - t_i),$$

with  $p(t)$  the infectiousness rate,  $t_i$  the time corresponding to the  $i$ th retweet arrival, and  $d_i$  the number of followers of the  $i$ th retweeting individual.

Furthermore, the memory kernel  $\phi(s)$ , a probability distribution for the time intervals between a tweet by the followee and its retweet by the follower, has been shown to be heavily tailed in a variety of social networks [103]. It is fitted to the empirical data by the function

$$\phi(s) = \begin{cases} 0 & \text{for } s < 0, \\ c_0 & \text{for } 0 \leq s \leq s_0, \\ c_0(s/s_0)^{-(1+\theta)} & \text{for } s > s_0, \end{cases}$$

where the parameters  $c_0$ ,  $s_0$ , and  $\theta$  are known.

The model is defined so that the daily cycles of human activity are naturally translated into cycles of retweet activity. The time dependence of the infectious rate is, therefore, defined as

$$(3.3) \quad p(t) = p_0 \left\{ 1 - r_0 \sin \left( \frac{2\pi}{T_m} (t + \phi_0) \right) \right\}^{\tau_m} \sqrt{e^{-(t-t_0)}}.$$

The parameters  $p_0$ ,  $r_0$ ,  $\phi_0$ , and  $\tau_m$  correspond to the intensity, the relative amplitude of the oscillation, its phase, and the characteristic time of popularity decay, respectively. They are fitted through a least square error (LSE) minimization procedure over the aggregation of retweet events over time bins of length  $\delta t$ .

Another improvement over the traditional parametric forms for HPs involves the addition of a nonlinearity into the expression for the CIF,

$$(3.4) \quad \lambda_{HP}(t) = g \left( \mu + \sum_{t_i < t} \phi(t - t_i) \right),$$

in which  $g(\cdot)$  correspond to a so-called *link function*, e.g., sigmoid:

$$(3.5) \quad g(x) = \frac{1}{1 + e^{-x}}.$$

The work in [84] proposes a procedure for simultaneously learning  $g(\cdot)$ ,  $\mu$ , and  $\phi(t) = \alpha\kappa(t)$ , with  $\kappa(t)$  taken as  $e^{-t}$ , for assuring convergence of the algorithm.

The procedure is formulated using a moment-matching idea over a piecewise-constant approximation for  $g(\cdot)$ , which leads to the definition of the objective function as a summation:

$$(3.6) \quad \min_{g \in \mathcal{G}, \mathbf{W}} \frac{1}{n} \sum_{i=1}^n \left( N_i - \int_0^{t_i} g(w \cdot x_t) dt \right)^2,$$

with  $w = (\mu, \alpha)^T$  and  $x_t = (1, \sum_{t_i \in \mathcal{H}_t} \kappa(t - t_i))^T$ .

The algorithm is run by recursively updating the estimates  $\hat{w}$  with the Isotron algorithm (see [38]), and  $\hat{g}$  with a projected gradient descent step.

Theoretical bounds for the approximation error of the method are also given, along with extensions of the algorithm for general point processes, monotonically decreasing nonlinearities, low-rank processes, and multidimensional HPs.

Another possible enhancement for modeling the parametric HP is the use of a composition of Gaussian kernels of different bandwidths, as in [91]. The core idea is that the maximum nonzero frequency component of the kernel is bounded as the same value of the intensity function, since that is simply a weighted sum of the basis

functions. Therefore, for every value of the tolerance  $\xi$ , it is possible to find a frequency value  $\omega_0$  such that

$$(3.7) \quad \int_{\omega_0}^{\infty} |\hat{\lambda}(\omega)| d\omega \leq \xi.$$

From this  $\omega_0$ , the method then defines the triggering function  $\phi(t)$  as a composition of  $\tilde{D}_\phi$  Gaussian functions, with  $\tilde{D}_\phi$  equally spaced values of bandwidth over the interval  $[0, \omega_0]$ .

The estimation of the value of the intensity function, for transforming into the frequency domain, is done through a kernel density estimation with Gaussian kernels of bandwidth fixed according to Silverman's rule of thumb,

$$(3.8) \quad h = \sqrt[5]{\frac{4\hat{\sigma}^5}{3n}} \approx \frac{1.06\hat{\sigma}}{\sqrt[5]{n}},$$

where  $\hat{\sigma}$  is the standard deviation of the time intervals, and  $n$  is the number of events of a given sequence.

After the Gaussian functions are defined, it remains to estimate the model coefficients  $\Theta$ , with the impact matrix now an impact tensor  $A = \{\alpha_{ijk}\}$ , through a convex surrogate loss penalized by parameters related to the sparsity of the matrix, the temporal sparsity of the kernels, and the pairwise similarity,

$$(3.9) \quad \underset{\Theta \geq 0}{\operatorname{argmin}} -\mathcal{L}_\Theta + \gamma_S \|A\|_1 + \gamma_G \|A\|_{1,2} + \gamma_P E(A),$$

where

- $\|A\|_1 = \sum_{i,j,k} |\alpha_{ijk}|$  is the  $L^1$ -norm of the tensor, which is related to its temporal sparsity; this causes the excitation functions to go to zero at infinity, therefore maintaining the stability of the process;
- $\|A\|_{1,2} = \sum_{i,j} \|\{\alpha_{ij1}, \dots, \alpha_{ij\tilde{D}_\phi}\}\|_2$  is related to the sparsity over the  $\tilde{D}_\phi$  basis functions of a given node of the process and enforces the local independence of the process;
- $E(A) = \sum_i \sum_{i' \in \mathcal{C}_i} \|\alpha_i - \alpha_{i'}\|_F^2 + \|\alpha^i - \alpha^{i'}\|_F^2$  is a coefficient to enforce the pairwise similarity of the process, in which  $\mathcal{C}_i$  corresponds to the cluster to which node  $i$  belongs,  $\|\cdot\|_F$  is the Frobenius norm,  $\alpha_i = \{\alpha_{ijk}\}$  for fixed  $i$ , and  $\alpha^i = \{\alpha_{ijk}\}$  for fixed  $j$ . This means that if  $i$  and  $i'$  are similar types of events, then their mutual excitation effects should be similar as well;
- $\gamma_S$ ,  $\gamma_G$ , and  $\gamma_P$  are coefficients to be tuned for the model.

The estimation is done through an expectation-maximization procedure close to those of [59] and [106], which first randomly initializes the impact tensor and the vector of baseline intensities  $\mu$ , and then iterates as follows:

1. Estimate the probability that each event was generated by each of the compositional basis kernels, as well as the baseline intensity.
2. Average the probabilities over all events of all training sequences for updating the coefficients of each basis function and the baseline intensity.

The two steps are repeated until the parameter estimates converge.

**3.2. Scalability.** Another setting in which a parametric choice of kernels is highly convenient is the scalability of the inference procedure for high-dimensional networks and sequences with a large number of events, which occur in several domains, such as social interaction data, which is simultaneously large (i.e., large numbers of posts),

high-dimensional (numerous users), and structured (i.e., the users' interactions are not random but, instead, present some regularities).

One interesting inference method in this direction is the work presented in [44], which achieves a complexity  $O(nD)$ , with  $n$  the number of events comprising the process history and  $D$  the dimension of the impact matrix of the process.

The method used, entitled *scalable low-rank Hawkes processes* (SLRHP), takes advantage of the memoryless property of the exponential and the underlying regularity of large networks connected to social events: The memoryless property (which means that, in HPs with exponential excitation functions, the effect of all past events on the intensity value of a given point can be computed just from the time of the last event before that point) speeds up the intensity computing portion of the inference procedure iterations, while the underlying regularity of large impact matrices associated with the social phenomena allows the dynamics of large-dimensional HPs to be captured by impact matrices of much smaller magnitudes.

The baseline rates and excitation functions of the model are then defined using a low-rank approximation

$$(3.10) \quad \mu_i(t) = \sum_{j=1}^E P_{ij} \tilde{\mu}_j,$$

$$(3.11) \quad \phi_{mi}(t) = \sum_{j,l=1}^E P_{ij} P_{ml} \tilde{\phi}_{lj}(t),$$

in which  $P \in \mathbb{R}_+^{D \times E}$  is a projection matrix from the original  $D$ -dimensional space to a low-dimensional space  $E$  ( $E \ll D$ ). This projection can also be seen as a low-rank approximation of the excitation function matrix  $\Phi$ , in which

$$(3.12) \quad \Phi = P \tilde{\Phi} P^T.$$

Since  $E \ll D$ , the formulated low-rank approximated inference algorithm SLRHP manages to do the following:

1. capture a simplified underlying regularity imposed on the inferred intensity rates' parameters by adopting sparsity-inducing constraints on the model parameters;
2. lower the number of parameters for both the baseline rates and the excitation kernels, with the  $D$  natural rates and  $D^2$  triggering kernels being lowered to  $r$  and  $E^2$ , respectively. This advantage is diminished slightly by the additional cost of inferring the  $(D \times E)$ -sized projection matrix  $P$ .

Another way of dealing with the scalability issues of multivariate HPs, in terms of both the total number of events in the sequences and the number of nodes, is through the mean-field treatment, as described in [6]. Compared with the SLRHP method, which focuses on reducing the dimensionality of the underlying network, the mean-field treatment focuses on finding closed-form expressions for the approximate estimations involved in the optimization defined over the network in its full dimensionality. The key step of this method is to consider that the arrival intensity of each node of the process is in a wide sense stationary, which implies the stability condition for the excitation matrix, and fluctuates only slightly around its mean value.

This last assumption, called the *mean-field hypothesis*, posits that if  $\lambda^i(t)$  corresponds to the intensity of the  $i$ th node and  $\hat{\Lambda}^i$  corresponds to the empirical estimator

of the first-order statistics of that node,

$$(3.13) \quad \tilde{\Lambda}^i = \frac{N_T^i}{T},$$

where  $N_T^i$  corresponds to the total number of events that have arrived at node  $i$  up to the final time of the simulation horizon  $[0, T]$ , then we have that

$$(3.14) \quad \frac{|\lambda^i(t) - \tilde{\Lambda}^i|}{\tilde{\Lambda}^i} \ll 1 \quad \forall t \in [0, T].$$

The condition defined by (3.14) is met when (i)  $\|\phi(t)\| \ll 1$ , independently of the shape of  $\phi(t)$ ; (ii) the dimensionality of the multivariate HP is sufficiently high; and also (iii)  $\phi(t)$  changes sufficiently slowly, so that the influence of past events averages to a near constant value.

From this, we can recover the parameters  $\theta^i$  from the intensity function  $\lambda_t^i$ , and the intensity function from the first-order statistics  $\tilde{\Lambda}^i$ , as

$$(3.15) \quad \log \lambda_t^i \simeq \log \tilde{\Lambda}^i + \frac{\lambda^i(t) - \tilde{\Lambda}^i}{\tilde{\Lambda}^i} - \frac{(\lambda^i(t) - \tilde{\Lambda}^i)^2}{2(\tilde{\Lambda}^i)^2}$$

and

$$(3.16) \quad \lambda^i(t) = \mu^i + \int_{0^-}^t \sum_{j=1}^D \phi^{ji}(t) dN_t^i.$$

The method yields mean-field estimates for the parameters with error that decays proportionally to the inverse of the final time  $T$  of the sequences:

$$(3.17) \quad \mathbb{E}(\theta^i) \approx \theta_{MF}^i,$$

$$(3.18) \quad \mathbf{cov}(\theta^{ji}, \theta^{j'i'}) \sim \frac{1}{T},$$

where  $\theta_{MF}^i$  refers to the mean-field estimator of the parameters.

**3.3. Training.** Elaborating further on some difficulties involved in inferring parametric kernels from real data, an interesting method for truncated sequences is described in [92].

In the case of learning HP parameters from real data, one often has to deal with sequences which are only partially observed, i.e., the time event arrivals are only available over a finite time window.

This poses a challenge concerning the robustness of the learning algorithms, since the triggering pattern from unobserved events is not considered: The inference deals with the error induced by computing intensity values over finite time windows, i.e., by computing intensity values along simulation horizons  $[0, T]$ ; the closed-form equation would assume that its value at 0 is simply the baseline rate  $\mu$ .

From the expression for the CIF,

$$(3.19) \quad \lambda_{HP}(t) = \mu + \sum_{t_i < t} \phi(t - t_i)$$

for  $t \in [0, T]$ , we have that, taking some value  $T' \in [0, T]$ , we may split the triggering effect term into two parts:

$$(3.20) \quad \lambda_{HP}(t) = \mu + \sum_{t_i < T'} \phi(t - t_i) + \sum_{T' \leq t_i < T} \phi(t - t_i)$$

for  $t \in [T', T]$ .

If we are observing the sequences only over the interval  $[T', T]$ , the second term is implicitly ignored, which may lead to severe degradation of the learning procedure, especially in the case that the excitation functions decay slowly.

The method proposed handles this issue through a *sequence-stitching method*. The trick is to sample “candidate predecessor events” and choose the most likely one from its similarity w.r.t. the observed events. The augmented sequences can then be used for the actual HP parameters’ learning algorithm.

In practice, this means that, given  $M$  sequences  $\mathcal{S}_m$  realized over  $[T', T]$ , the method would *not* learn from the regular MLE formula

$$(3.21) \quad \theta^* = \operatorname{argmax}_{\theta} llh(\{\mathcal{S}_m\}_{m=1}^M, \theta),$$

but instead from some expression which takes into account the expected influence of unobserved predecessor events

$$(3.22) \quad \theta_{SDC}^* = \operatorname{argmax}_{\theta} \mathbb{E}_s \mathcal{H}_{T'} llh([\mathcal{S}, \mathcal{S}_m], \theta),$$

where  $\mathcal{H}_{T'}$  corresponds to the distribution over all possible sequences of events happening before time  $T'$ . In practice, this expectation is computed not over the real distribution, but over some finite number of (relatively few) samples, such as 5 or 10. This finite sample approximation converts (3.22) into

$$(3.23) \quad \theta_{SDC}^* = \operatorname{argmax}_{\theta} \sum_{\mathcal{S}_{stitch} \in \mathcal{K}} p(\mathcal{S}_{stitch}) llh(\mathcal{S}_{stitch}, \theta),$$

where  $p(\mathcal{S}_{stitch})$  is the probability of the stitched sequence obtained from the concatenation of the original observed sequence and one of the sample predecessor candidate sequences. This probability is obtained by normalizing over similarity values over the candidate sequences obtained from some similarity function  $\mathbb{S}(\cdot, \cdot)$  of the form

$$(3.24) \quad \mathbb{S}(\mathcal{S}_k, \mathcal{S}) = \sqrt[\psi]{e^{-\|f(\mathcal{S}_k) - f(\mathcal{S})\|^2}},$$

where  $f(\cdot)$  is some feature of the event sequence, and  $\psi \in \mathbb{R}^+$  is some scale parameter. By fixing the excitation pattern matrix  $\kappa$  as being composed of exponentials  $\kappa(t) = e^{-\beta t}$  and imposing sparsity constraint  $\|\alpha\|_1 = \sum_{i,j} |\alpha_{ij}|$  over the impact matrix, the equivalent problem,

$$(3.25) \quad \theta^* = \operatorname{argmax}_{\mu \geq 0, \alpha > 0} \sum_{\mathcal{S}_{stitch} \in \mathcal{K}} p(\mathcal{S}_{stitch}) llh(\mathcal{S}_{stitch}, \theta) + \gamma \|\alpha\|_1,$$

is shown to be solved through expectation-maximization updated for both  $\mu$  and  $\alpha$ . The method is more suitable for slowly decaying excitation patterns, in which the influence of the unobserved events is more prominent. In the case of exponentials with large values for the decay factor  $\beta$ , the improvement margins mostly vanish.

Another interesting improvement in the training procedure of MLE for HP parametric functions involves the complementary use of adversarial and discriminative learning, as in [94]. Although adversarial training has gained ever-increasing relevance in neural network based models in the last few years due to the popularization of generative adversarial networks (GANs) [25] and their variants, as will be discussed further in section 5, keeping the assumption of a simple parametric shape for the excitation function is a way to insert domain-specific knowledge into the inference procedure.

The key idea of this complementary training is that, while the discriminative loss, here defined as the mean squared error (MSE) between discretized versions of predicted and real sequences, tends to direct the parameter updates toward smoother prediction curves, the adversarial loss tends to push the temporally evenly distributed sampled sequences toward more realistic-looking curves.

For gradient descent based updates, a discretization of the point process is carried out to approximate the predictions through a recursive computation of the integral of the intensity function (the *compensator* portion of the loglikelihood function) in a closed-form expression. The parameters of the complementary training are actually initialized by a purely MLE procedure, which was found to be more insensitive to initial points. The *MLE + GAN* training updates then follow. The full procedure can be summarized by the following steps:

1. Subdivide the  $M$  original sequences on  $[0, T]$  into training (on  $[0, T^{tn}]$ ), validation (on  $[T^{tn}, T^{vd}]$ ), and test (on  $[T^{vd}, T]$ ) portions, using previously defined parameters  $T^{tn}$  and  $T^{vd}$ .
2. Initialize the parameters of the model through the “purely” MLE procedure.
3. Sample  $M$  sequences from the model over the interval  $[0, T^{tn}]$ .
4. By choosing a specific parameter shape for the excitation function and binning both the original and the simulated sequences over equally spaced intervals, define the closed-form expression for the MSE (discriminative) loss  $\mathcal{L}_{MSE}$  over all the sequences and dimensions of the process.
5. Define the GAN (adversarial) loss of the model over the sequences as

$$(3.26) \quad \mathcal{L}_{GAN} = \begin{cases} E_{\mathcal{S}^{tn} \mathbb{P}(\mathcal{S}^{tn})} [F_W(Y_{\theta_S}(\mathcal{S}^{tn}))] - E_{\mathcal{S}^{tn} \mathbb{P}(\mathcal{S}^{tn})} [F_W(Y_{\theta_S}(\mathcal{S}^{tn}))] \\ \text{for the critic network } F_W, \\ -E_{\mathcal{S}^{tn} \mathbb{P}(\mathcal{S}^{tn})} [F_W(Y_{\theta_S}(\mathcal{S}^{tn}))] \\ \text{for the training model (generator),} \end{cases}$$

where  $\mathbb{P}(\mathcal{S}^{tn})$  corresponds to the underlying probability distribution that we assume has generated the original sequences  $\mathcal{S}^{tn}$ ,  $F_W\{\cdot\}$  refers to a neural network that can compute the so-called Wasserstein distance, a metric for difference among distributions which will be further explained in section 5, and  $Y_{\theta_S}(\cdot)$  corresponds to the parametric model for sequence generation.

6. Compute the joint loss for the MLE and GAN portions as

$$(3.27) \quad \mathbf{L}_{MLE+GAN} = \gamma_{GAN} \mathbf{L}_{MLE} + (1 - \gamma_{GAN}) \mathbf{L}_{GAN} \text{ for } \gamma_{GAN} \in [0, 1].$$

7. Compute gradients of  $\mathbf{L}_{MLE+GAN}$  over each parameter of model  $Y_{\theta_S}$ .
8. Update parameter estimates of the training model with  $\eta \nabla_{\theta_S}(\mathbf{L}_{MLE+GAN})$  for some learning rate  $\eta$ .
9. Repeat steps 3 to 8 until convergence.



The method is an interesting combination of the enriched dynamical modeling from the adversarial training strategy with the robustness over small training sets of the parametric-based MLE estimation for HPs.

In this section, we have provided a comprehensive analysis of progress in HP modeling and inference for excitation functions assumed to be of a simple parametric shape, along with their compositions and variants. In the next section, we will discuss advances in nonparametric HP excitation function strategies.

**4. Nonparametric HPs.** Nonparametric HPs consider that a rigid and simple parametric assumption for the triggering kernel may not be enough to capture all the subtleties of the excitation effects that could not be retrieved from the data. They may be broadly divided into two main approaches:

1. Frequentist.
2. Bayesian.

We will briefly discuss their variants in this section.

**4.1. Frequentist Nonparametric HPs.** The frequentist approach to HP modeling and inference consists in assuming that the excitation function (or matrix) can be defined as a binned grid (or a set of grids), in which the values of the functions are taken as piecewise constant inside each bin, and the width of the bin is (hopefully) expressive enough to model the local variations of the self-excitation effect.

They were first developed in [46], [5], and [8]. In the case of [46], the final values of the bins were found by solving a discretized ordinary differential equation, implied by the branching structure of the discretized triggering kernel and background rate over the data, through iterative methods. The approach in [5] and [8], on the other hand, recovers the piecewise constant model by exploiting relations, in the frequency domain, among the triggering kernel, the background rate, and the second-order statistics of the model, also obtained in a discretized way over the data.

The increased expressiveness of this type of excitation model incurs two main drawbacks:

- The bin division grid concept is close to that of a histogram over the distance among events, which usually requires much larger datasets to lead to accurate predictions, in contrast to parametric models, which behave better on shorter and fewer sequences but most likely underfit on large sequence sets.
- The time of the inference procedure may also be much larger, since it involves sequential binning computation procedures which cannot take advantage of the Markov property of parametric functions such as the exponentials.

Two improvements, to be discussed in this subsection, deal with these exact drawbacks through

- acceleration of the computations over each sequence and/or over each bin of the excitation matrix/function;
- reduction of the times of binning computational procedures through a so-called online update of the bin values.

**4.2. Acceleration of Impact Matrix Estimation through Matching of Cumulants.** One acceleration strategy is developed in [2], which replaces the task of estimating the excitation functions directly by estimating their cumulative values, i.e., their integrated values from zero up to infinity, which is enough to quantify the causal relationships among the nodes. That is, instead of estimating  $\phi_{ij}(t)$  for each node,

the method estimates a matrix  $\|\Phi(t)\| = \{\|\phi_{ij}(t)\|\}$ , in which

$$(4.1) \quad \|\phi_{ij}(t)\| = \int_0^\infty \phi_{ij}(t) dt \quad \forall (i, j) \in D \times D.$$

The method, called nonparametric Hawkes process cumulant (NPHC), then proceeds to compute, from the sequences, moment estimates  $\hat{\mathbf{R}}$  up to the third order. It then finds some estimate  $\|\hat{\Phi}(t)\|$  of this cumulant matrix which minimizes the  $L^2$  squared error between these estimated moments and the actual moments  $\mathbf{R}(\|\Phi(t)\|)$ , which are uniquely determined from  $\|\Phi(t)\|$ :

$$(4.2) \quad \left\| \hat{\Phi}(t) \right\| = \operatorname{argmin}_{\|\Phi(t)\|} \|\mathbf{R}(\|\Phi(t)\|) - \hat{\mathbf{R}}\|^2.$$

This  $L^2$  minimization comes from the fact that, by defining

$$(4.3) \quad \mathbf{V} = \left( \mathbb{I}^D - \left\| \hat{\Phi}(t) \right\| \right)^{-1},$$

one may express the first-, second-, and third-order moments of the process as

$$(4.4) \quad \Lambda^i = \sum_{m=1}^D V^{im} \mu^m,$$

$$(4.5) \quad \nu^{ij} = \sum_{m=1}^D \Lambda^m V^{im} V^{jm},$$

$$(4.6) \quad K^{ijk} = \sum_{m=1}^D (V^{im} V^{jm} \nu^{km} + V^{im} \nu^{jm} V^{km} + \nu^{im} V^{jm} V^{km} - 2\Lambda^m V^{im} V^{jm} V^{km}),$$

and thus we may find an estimator  $\hat{\mathbf{V}} = \operatorname{argmin}_{\mathbf{V}} \mathbf{L}_{NPHC}(\mathbf{V})$ , with  $\mathbf{L}_{NPHC}(\mathbf{V})$  defined as

$$(4.7) \quad \mathbf{L}_{NPHC}(\mathbf{V}) = (1 - \gamma_{NPHC}) \|\mathbf{K}^c(\mathbf{V}) - \hat{\mathbf{K}}^c\|_2^2 + \gamma_{NPHC} \|\boldsymbol{\nu}(\mathbf{V}) - \hat{\boldsymbol{\nu}}\|_2^2,$$

where  $\gamma_{NPHC}$  is a weighting parameter,  $\|\cdot\|_2^2$  is the Frobenius norm, and  $\mathbf{K}^c = \{K^{ijj}\}_{1 \leq i, j \leq D}$  is a two-dimensional compression of the tensor  $\mathbf{K}$ . From this expression, by setting

$$(4.8) \quad \gamma_{NPHC} = \frac{\|\hat{\mathbf{K}}^c\|_2^2}{\|\hat{\mathbf{K}}^c\|_2^2 + \|\hat{\boldsymbol{\nu}}\|_2^2},$$

one may arrive at

$$(4.9) \quad \left\| \hat{\Phi}(t) \right\| = \mathbb{I}^D - \hat{\mathbf{V}}^{-1}.$$

The estimates of moments in the algorithm are actually computed through truncated and discretized (binned) countings along a single realization of the process and, since the real-data estimates are usually not symmetric, the estimates are averaged along positive and negative axes.

Also, for  $D = 1$ , the estimate  $\|\hat{\Phi}(t)\|$  can be estimated solely from the second-order statistics. For higher-dimensional processes, it is the skewness of the third-order moment that uniquely fixes  $\|\hat{\Phi}(t)\|$ .

**4.3. Online Learning.** Another improvement of the method consists in updating the parameters of the discretized estimate of the excitation function through a single pass over the event sequence, i.e., an online learning procedure [96].

In the case of the referred algorithm, the triggering function is assumed to

1. be positive,

$$(4.10) \quad \phi(t) \geq 0 \quad \forall t \in \mathbb{R};$$

2. have a decreasing tail, i.e.,

$$(4.11) \quad \sum_{k=m}^{\infty} (t_k - t_{k-1}) \sup_{x \in (t_{k-1}, t_k]} |f(y)| \leq \zeta_f(t_{i-1}) \quad \forall i > 0$$

for some bounded and continuous  $\zeta_f : \mathbb{R}^+ \mapsto \mathbb{R}^+$  such that  $\lim_{t \rightarrow \infty} \zeta_f(t) = 0$ ;

3. belong to a reproducing kernel Hilbert space (RKHS), which here is used as a tool for embedding similarity among high-dimensional and complex distributions into lower-dimensional ones.

The method proceeds by taking the usual expression for the loglikelihood function,

$$(4.12) \quad llh_{\tilde{T}}(\boldsymbol{\lambda}) = - \sum_{d=1}^D \left( \int_0^{\tilde{T}} \lambda_d(s) ds - y_{d,k} \log \lambda_d(t_k) \right),$$

and optimizing instead over a discretized version of it,

$$(4.13) \quad \begin{aligned} llh_{\tilde{T}}(\boldsymbol{\lambda}) &= \sum_{d=1}^D \sum_{k=1}^{M(t)} \left( \int_{\chi_{k-1}}^{\chi_k} \lambda_d(s) ds - y_{d,k} \log \lambda_i(t_k) \right) \\ &= \sum_{d=1}^D \Delta L_{d,t}(\lambda_d), \end{aligned}$$

with  $(t_1, \dots, t_{n(t)})$  denoting the event arrival times over an interval  $[0, \tilde{T}]$  and with a partitioning  $\{0, \chi_1, \dots, \chi_{M(t)}\}$  of this interval  $[0, \tilde{T}]$  such that

$$(4.14) \quad \chi_{k+1} = \min_{t_i \geq \chi_k} \{ \iota * \lfloor \chi_k / \iota \rfloor + \iota, t_i \}$$

for some small  $\iota > 0$ . The discretized version can then be expressed as

$$(4.15) \quad \begin{aligned} llh_{\tilde{T}}^{(\iota)}(\boldsymbol{\lambda}) &= \sum_{d=1}^D \sum_{k=1}^{M(t)} \left( (\chi_k - \chi_{k-1}) \lambda_d(\chi_k) - y_{d,k} \log \lambda_i(\chi_k) \right) \\ &= \sum_{d=1}^D \Delta L_{d,\tilde{T}}^{(\iota)}(\lambda_d). \end{aligned}$$

The optimization procedure is carried out at each slot of the  $M(t)$  partition, taking into account the following items:

- A truncation over the intensity function effect, i.e.,

$$(4.16) \quad \phi(t) = 0 \quad \forall t > t_{max},$$

to simplify the optimization of the integral portion of the loss. The error over this approximation is shown to be bounded by the decreasing tail assumption.

**Table 1** Comparison of the computational complexity of parametric and nonparametric HP estimation methods, extracted from [2].  $Iter$  is the number of iterations of the optimization procedure,  $\tilde{D}_\phi$  is the number of composing basis kernels of  $\phi(t)$ ,  $D$  is the dimensionality of the multivariate HP,  $n_{max}$  is the maximum number of events per sequence, and  $M$  is the number of components of the discretization applied to  $\phi(t)$ . Complexities are taken from [2] and [96].

Method	Total complexity
ODE HP [106]	$\mathcal{O}(Iter * \tilde{D}_\phi(n_{max}^3 D^2 + M * (n_{max} D + n_{max}^2)))$
Granger causality HP [91]	$\mathcal{O}(Iter * \tilde{D}_\phi n_{max}^3 D^2)$
Wiener-Hopf eq. HP [8]	$\mathcal{O}(n_{max} D^2 M + D^4 M^3)$
NPHC [2]	$\mathcal{O}(n_{max} D^2 + Iter * D^3)$
Online learning HP [96]	$\mathcal{O}(Iter * D^2)$

**Table 2** Performance comparison of several multivariate HP estimation methods in the MemeTracker [45] dataset, extracted from [2]. The relative error between a ground truth impact matrix  $\alpha = \{\alpha_{ij}\}$  and its estimate  $\hat{\alpha} = \{\hat{\alpha}_{ij}\}$  is simply  $\sum_{i,j} |\alpha_{ij} - \hat{\alpha}_{ij}| / |\alpha_{ij}| \mathbf{1}_{\{\alpha_{ij} \neq 0\}} + |\hat{\alpha}_{ij}| \mathbf{1}_{\{\alpha_{ij} = 0\}}$ . ( $\mathbf{1}_{\{\cdot\}}$  is the indicator function.)

Performance metric	Multivariate HP estimation method			
	ODE HP [106]	Granger causality HP [91]	ADM4 [105]	NPHC [2]
Relative error	0.162	0.19	0.092	0.071
Estimation time (s)	2944	2780	2217	38

- A Tikhonov regularization over the coefficients  $\mu_d$  and  $\phi_{d,d}$ , which is simply the addition of weighted  $\|\mu_d\|^2$  and  $\|\phi_{d,d}\|^2$  terms to the loss function, to keep their resulting values small.
- A projection step for the triggering function optimization part, to keep them all positive.

The recent improvements over frequentist nonparametric HP estimation focus on two main strategies:

- Speeding up inference through replacement of the excitation matrix as objective by the matrix of cumulants, which is shown to be enough to capture the mutual influence among each pair of nodes.
- An online learning procedure, which uses some assumptions on the kernels (positive, decreasing tail RKHS) to recover estimates of the HP parameters over a single pass on some partitioning of the event arrival timeline.

A comparison of the complexity of several parametric and frequentist nonparametric HP estimation methods is shown in Table 1, and performance metrics are shown in Table 2. In general, it is possible to see the focus of more recent methods, such as NPHC and the online learning approach, on reducing the complexity per iteration of the resulting estimation procedure through approximation assumptions on the underlying model.

**4.4. Bayesian Nonparametric HPs.** Another nonparametric treatment of HPs revolves around the assumption that the triggering kernel and the background rates can be modeled by distributions (or mixtures of distributions) from the so-called exponential family, which, through their conjugacy relationships, allow for closed-form computations of the sequential updates in the model. These were mainly proposed in [20], [18], and [102].

In [20], HPs are used for modeling the clustering of document streams, captur-

ing the dynamics of arrival time patterns, used together with textual content-based clustering.

It seems logical that news and other media-related information sources revolving around a given occurrence, such as a natural catastrophe, a political action, or a celebrity scandal, are related not only in terms of word content, but also their time occurrences, as journalists tend to release more and more content on a topic of high public interest, but will slow the pace of publication as interest gradually vanishes or shifts toward other subjects.

The main idea is to unite both Bayesian nonparametric inference, which is a scalable clustering method that allows for new clusters to be added as the number of samples grows, and HPs. The corresponding Bayesian nonparametric model, the Dirichlet process, captures the diversity of event types, while the HPs capture the temporal dynamics of the event streams.

A Dirichlet process  $DP(\alpha, G_0)$  can be roughly described as a probability distribution over probability distributions. It is defined by a concentration parameter  $\alpha$ , proportional to the level of discretization (“number of bins”) of the underlying sampled distribution, and a base distribution  $G_0$ , which is the distribution to be discretized. As an example, for  $\alpha$  equal to 0, the distribution is concentrated at a single value, while, in the limit as  $\alpha$  goes to infinity, the sampled distribution becomes continuous.

The corresponding hybrid model, the Dirichlet–Hawkes process (DHP), is defined by

- $\mu$ , an intensity parameter;
- $\mathbb{P}_0^{DHP}(\theta_{DHP})$ , a base distribution over a given parameter space  $\theta_{DHP} \in \Theta_{DHP}$ ;
- a collection of excitation functions  $\phi_{\theta_{DHP}}(t, t')$ .

After an initial time event  $t_1$  and an excitation function parameter  $\theta_{DHP}^1$  are sampled from these base parameters  $\mu$  and  $\mathbb{P}_0^{DHP}(\theta_{DHP})$ , respectively, the DHP is then allowed to alternate between the following methods:

1. Sampling new arrival events  $t_i$  from the current value of  $\theta_{DHP}$  for the excitation function, with probability

$$(4.17) \quad \frac{\mu}{\mu + \sum_{i=1}^{n-1} \phi_{\theta_{DHP}^i}(t_n, t_i)}.$$

2. Sampling a new value for  $\theta_{DHP}$  with probability

$$(4.18) \quad \frac{\sum_{i=1}^{n-1} \phi_{\theta_{DHP}^i}(t_n, t_i)}{\mu + \sum_{i=1}^{n-1} \phi_{\theta_{DHP}^i}(t_n, t_i)}.$$

In this way, we are dealing with a superposition of HPs in which the arrival events tend toward processes with higher intensities, i.e., the preferential attachment, but which also allows for diversity, since there is always a nonzero probability of sampling a new HP from the baseline intensity  $\mu$ .

By defining the excitation functions  $\phi_\theta$  as a summation of parametric kernels

$$(4.19) \quad \phi_{\theta_{DHP}}(t_i, t_j) = \sum_{l=1}^K \alpha_{\theta_{DHP}}^l \kappa_{DHP}^l(t_i - t_j),$$

the model can be made even more general.

The approach in [18], besides the random histogram assumption for the triggering kernel, similar to the frequentist case, also considers the case of the kernel being defined by a mixture of beta distributions, with the model being updated through a sampling procedure (Markov chain Monte Carlo).

The work in [102] proposes a Gamma distribution over the possible values of  $\mu$  with the triggering kernel modeled as

$$(4.20) \quad \phi(\cdot) = \frac{\mathcal{GP}(\cdot)^2}{2},$$

where  $\mathcal{GP}(\cdot)$  is a Gaussian process [66].

These assumptions allow for closed-form updates over the posterior distributions over the background rate and the triggering kernel, which are claimed to be more scalable and efficient than the plain binning of the events.

In the next section, we will explore neural architectures, which were introduced for the modeling and generation of HPs through a more flexible representation of the effect of past events on the intensity function.

**5. Neural Network Based HPs.** In this section, we discuss the neural network based formulations of HP modeling. The main idea is to capture the influence of past events on the intensity function in a nonlinear, and thus hopefully more flexible, way.

This modeling approach makes use of recurrent models, which, in their simplest formulation, encode sequences of states

$$(5.1) \quad (z_0^s, z_1^s, \dots, z_N^s)$$

and outputs

$$(5.2) \quad (z_0^o, z_1^o, \dots, z_N^o),$$

in a way such that each state  $z_{i+1}^s$  can be obtained by a composition of the immediately preceding state  $z_i^s$  and a so-called hidden state  $h_i$  that captures the effect of the other past states,

$$(5.3) \quad h_i = \sigma_h(W_s z_i^s + W_h h_{i-1} + b_h),$$

$$(5.4) \quad z_i^o = \sigma_o(W_o z_i^o + b_o),$$

in which  $W_o$ ,  $W_s$ ,  $W_h$ ,  $b_h$ , and  $b_o$  are parameters to be fitted by the optimization procedure, while  $\sigma_h$  and  $\sigma_o$  are nonlinearities such as a sigmoid or a hyperbolic tangent function.

In the case of HPs, the state to be modeled is the intensity function along a sequence of time event arrivals, and an additional assumption is that its intensity value decays exponentially between consecutive events.

In the case of most neural network based models, the inference also counts with a mark distribution for the case of marked HPs, in which a multinomial or some other multiclass distribution is fitted together with the recurrent intensity model.

Arguably the first such type of model, the recurrent marked temporal point process (RMTTP) [19], jointly models marked event data using a recurrent neural network (RNN) with exponentiated output for modeling the intensity.

For sequences of the type  $\{t_i, y_i\}_{i=1}^N$ , in which  $t_i$  corresponds to the time of the  $i$ th event arrival, while  $z_i^o$  refers to the type of event or mark, we have a hidden cell  $h_i$  described by

$$(5.5) \quad \mathbf{h}_i = \max \{W_o z_i^o + W^t t_i + W_h h_{i-1} + b_h, 0\}$$

and a CIF defined as a function of this hidden state,

$$(5.6) \quad \lambda(t) = \exp(\mathbf{v}^t \mathbf{h}_i + w_t(t - t_i) + b^t),$$

while a  $K$ -sized mark set can have its probability modeled by a Softmax distribution:

$$(5.7) \quad P(z_{i+1}^o = k | \mathbf{h}_i) = \text{Softmax}(k, \mathbf{V}_k^{z^o} \mathbf{h}_i + b_k^{z^o}) = \frac{\exp(\mathbf{V}_k^{z^o} \mathbf{h}_i + b_k^{z^o})}{\sum_{k=1}^K \exp(\mathbf{V}_k^{z^o} \mathbf{h}_i + b_k^{z^o})}.$$

The likelihood of the whole sequence can then be defined as a product of conditional probability density functions for each event,

$$(5.8) \quad llh(\{t_i, z_i^o\}_{i=1}^N) = \prod_{i=1}^N f(t_i, z_i^o),$$

with

$$(5.9) \quad f_\lambda(t) = \lambda(t) \exp\left(-\int_{t_n}^t \lambda(t) dt\right),$$

where  $t_n$  is the latest event that occurred before time  $t$ . From this  $f(t)$ , we may estimate the time of the next event as

$$(5.10) \quad t_{i+1} = \int_{t_i}^\infty t f_\lambda(t) dt.$$

This allows us to optimize the parameters of the RNN model over the loss equal to this likelihood function, composed by these conditional density functions.

Given a set of  $M$  training sequences  $\mathcal{S}^j = \{t_i^j, y_i^j\}_{i=1}^{N^j}$ , we want to optimize the weight parameters over a loss function defined as

$$(5.11) \quad llh(\{\mathcal{S}^j\}_{j=1}^M) = \sum_j^M \sum_i^{N^j} \log \left( P(z_{i+1}^{o,j} | \mathbf{h}_i) + \log f_\lambda(t_{i+1}^j - t_i^j | \mathbf{h}_i) \right).$$

This optimization procedure is usually done through the backpropagation through time (BPTT) algorithm, which proceeds by “unrolling” the RNN cells for a fixed number of steps, then calculating the cumulative loss along all these steps together with the gradients over each of the  $W$ ’s,  $w$ ’s,  $v$ ’s, and  $b$ ’s, then updating these parameters with a predefined learning rate until convergence.

One improvement on the RNN-based modeling approach is described in [88], referred to as “time series event sequence” (TSES), which consists of treating the mark sequences as being derived from another RNN model, instead of the multinomial distribution. This RNN for the marks is then jointly trained with the RNN for event arrival times.

Another improvement over this neural network based modeling approach is the neural HP [56], which uses a variant of the basic RNN, called long short-term memory (LSTM) [33], applied to the intensity function modeling.

The neural HP models the intensity function of a multitype event sequence by associating each  $k$ th event type with a corresponding intensity function  $\lambda_k(t)$ , such that

$$(5.12) \quad \lambda_k(t) = f_k(w_k^T \mathbf{z}^h(t)),$$

with

$$(5.13) \quad \mathbf{z}^h(t) = \mathbf{o}_i \odot (2g(2\mathbf{c}(t) - 1)).$$

The variables  $\mathbf{o}_i$  and  $\mathbf{c}(t)$  are defined through the following update rules:

$$(5.14) \quad \mathbf{c}_{i+1} = f_{i+1} \odot \mathbf{c}(t_i) + \mathbf{i}_{i+1} \odot \mathbf{z}_{i+1}.$$

The variables  $\mathbf{i}_{i+1}$ ,  $\mathbf{f}_{i+1}$ ,  $\mathbf{z}_{i+1}$ , and  $\mathbf{o}_{i+1}$  are defined similarly to the gated variables of a standard LSTM cell; see the original paper for their full definitions. The value of the  $\mathbf{c}(t)$  is assumed to decay exponentially among consecutive events as

$$(5.15) \quad \mathbf{c}(t) = \bar{\mathbf{c}}_{i+1} + (\mathbf{c}_{i+1} - \bar{\mathbf{c}}_{i+1}) \exp(-\delta_{i+1}(t - t_i))$$

for

$$(5.16) \quad \bar{\mathbf{c}}_{i+1} = \bar{f}_{i+1} \odot \bar{\mathbf{c}}(t_i) + \bar{\mathbf{i}}_{i+1} \odot \mathbf{z}_{i+1},$$

$$(5.17) \quad \delta_{i+1} = g(\mathbf{W}_\delta \mathbf{k}_i + \mathbf{U}_\delta \mathbf{h}(t_i) + \mathbf{d}_\delta).$$

The  $W$ 's,  $U$ 's, and  $d$ 's of the model are trained so as to maximize the loglikelihood over a set of sequences. Compared with the previous RNN-based model, in the neural HP the following hold:

1. The baseline intensity  $\mu_k$  is not implicitly considered constant, but instead is allowed to vary.
2. The variations of the cell intensity are not necessarily monotonic, because the influences of each event type on the cell values may decay at a different rate.
3. The sigmoid functions along the composition equations allow for an enriched behavior of the intensity values.

All these features contribute to an increased expressiveness of the model. Besides, as in the regular LSTM models, the “forget” gates  $\mathbf{f}_{i+1}$  are trained so as to control how much influence the past values of  $\mathbf{c}(t)$  will have on its present value, thus allowing the model to possess a “long-term” memory.

Another variant of the RNN-based HPs, introduced in [89], models the baseline rate and the history influence as separate RNNs. The baseline rate is taken as a time series, with its corresponding RNN updating its state at equally spaced intervals, such as five days. The RNN modeling the influence of the history of past events in future ones updates its state at each event arrival. This has been shown to increase the time and mark prediction performance, as demonstrated in Table 3.

Both the background rate time series  $\{\mu(t)\}_{t=1}^T$  and the marked event sequence  $\{m_i, t_i\}_{i=1}^N$  are modeled by LSTM cells:

$$(5.18) \quad (\mathbf{h}^\mu(t), \mathbf{z}_c^\mu(t)) = \text{LSTM}_\mu(\mu(t), \mathbf{h}^\mu(t-1) + \mathbf{z}_c^\mu(t-1)),$$

$$(5.19) \quad (\mathbf{h}^m(i), \mathbf{z}_c^m(i)) = \text{LSTM}_m(\mathbf{m}_i, \mathbf{h}^m(i) + \mathbf{z}_c^m(i-1)).$$



The  $\mathbf{h}$  and  $c$  states correspond to the hidden state and the long-term dependency terms, respectively, similarly to the neural HP. Both terms are concatenated in a single variable  $\mathbf{z}_e(t)$  for jointly training both RNN models:

$$(5.20) \quad \mathbf{z}_e(t) = \tanh(\mathbf{W}_f [\mathbf{h}^\mu(t), \mathbf{h}^m] + \mathbf{b}_f),$$

$$(5.21) \quad \mathbf{U}(t) = \text{Softmax}(\mathbf{W}_U \mathbf{z}_e(t) + \mathbf{b}_U),$$

$$(5.22) \quad \mathbf{u}(t) = \text{Softmax}(\mathbf{W}_u [\mathbf{z}_e(t), \mathbf{U}(t)] + \mathbf{b}_u),$$

$$(5.23) \quad z_s = \mathbf{W}_s \mathbf{z}_e(t) + b_s,$$

with  $U$  and  $u$  denoting the main event types and subtypes, respectively, and  $z_s$  denoting the composed timestamp of each event. The loss over which the model is trained is defined in a cross-entropy way as

$$(5.24) \quad \sum_{j=1}^N \left( -\mathbf{W}_U(j) \log(U(t, j)) - w_u(j) \log(u(t, j)) - \log(f(z_s(t, j)|h(t-1, j))) \right),$$

with

$$(5.25) \quad f(z_s(t, j)|h(t-1, j)) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(z_s(t, j) - z_{\bar{s}}(t, j))^2}{2\sigma^2}},$$

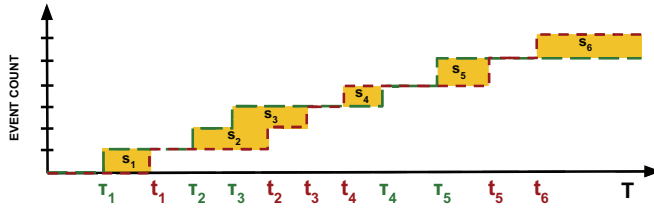
where  $z_{\bar{s}}(t, j)$  is the model predicted output for the corresponding event  $z_s(t, j)$ .

The model weights are then jointly trained, over the total loss function and under some correction for the frequency ratio of each event type, for both the background rate time series values and the event arrivals RNN, and are shown to outperform more “rigid” models.

Now, for the generation of HP sequences, both RMTTP and neural HPs have in common the fact that they intend to model the intensity function of the underlying process, so that new sequences may be sampled in a way that reproduces the behavior of the original dataset. This intensity modeling, however, has three main drawbacks:

1. It may be unnecessary, since the sequences may be simply produced by unrolling cells of corresponding RNN models.
2. The sequences from these intensity-modeling approaches are sampled using a “thinning algorithm” [62], which may result in slowed-down simulations in the case of repeatedly rejected event intervals.
3. These methods are trained by maximizing the loglikelihood over the training sequences, which is asymptotically equivalent to minimizing the KL-divergence over original and model distributions. This MLE approach is not robust in the case of multimodal distributions.

The model in [86] proposes approximating a generative model for generating event sequences by using an alternative metric of difference among distributions, the Wasserstein (or earth-moving) distance, discussed in section 3.



**Fig. 5** Intuition behind the distance metric  $|\cdot|_*$  between two given event sequences  $\{t_i\}$  and  $\{\tau_i\}$ . Based on a diagram in [86].

In the model, called Wasserstein generative adversarial temporal point process (WGANTPP), this Wasserstein loss is shown to be equal to

$$(5.26) \quad L' = \left[ \frac{1}{m} \sum_{i=1}^m F_w(G_\Theta(\mathcal{S}_s^i)) - \frac{1}{m} \sum_{i=1}^m F_w(\mathcal{S}_r^i) \right]$$

$$(5.27) \quad + \nu \sum_{i,j=1}^m \frac{|f_w(\mathcal{S}_r^i) - F_w(G_\Theta(\mathcal{S}_r^j))|}{|\mathcal{S}_r^i - G_\Theta(\mathcal{S}_s^j)|_*},$$

where the second term, along with the constant  $\nu$ , corresponds to the so-called Lipschitz constraints, related to the continuity of the models.

$\{\mathcal{S}_r^i\}_{i=1}^m$  are real-data sequences, while  $\{\mathcal{S}_s^i\}_{i=1}^n$  are sequences sampled from an HPP with a rate  $\lambda_{HPP}$  which is simply the expected arrival rate over all training sequences. The generator network  $G_\Theta$  and discriminator  $F_w$  are defined as

$$(5.28) \quad G_\Theta(\mathcal{S}_r) = \tilde{\mathcal{S}} = \{t_1, \dots, t_n\},$$

with

$$(5.29) \quad t(i) = g_G^x(f_G^x(h(i))) \text{ and } h(i) = g_G^h(f_G^h(z_s, h(i-1))),$$

$$(5.30) \quad F_w(\tilde{\mathcal{S}}) = \sum_{i=1}^n a(i),$$

and with

$$(5.31) \quad a(i) = g_D^a(f_D^a(h(i))) \text{ and } h(i) = g_D^h(f_D^h(t_i, h(i-1))),$$

with the  $g$ 's defined as nonlinearities and  $\tilde{\mathcal{S}}$  as some example time event sequence. The  $f$ 's are linear transformations, as in a standard RNN cell, with their corresponding weight matrices and bias vectors to be tuned using a stochastic gradient descent procedure.

The distance metric  $|\cdot|_*$  of two sequences  $\{t_i\}$  and  $\{\tau_i\}$ , for the case of purely temporal point processes in  $[0, T)$ , can be shown to be equivalent to

$$(5.32) \quad \sum_{i=1}^n |t_i - \tau_i| + (m - n) \times T - \sum_{i=n+1}^m \tau_i,$$

which has an intuitive graphical interpretation, as shown in Figure 5.

As previously discussed, the generator is trained so as to “fool” the discriminator, while the discriminator is trained so as to distinguish generated sequences from those of real data. This adversarial training procedure is roughly equivalent to gradient updates with opposite signs over their respective parameters: positive sign for the discriminator and negative sign for the generator.

At the end of the training procedure, one hopefully produces a generator network capable of producing sequences virtually indistinguishable from the real-data sequences.

This method, however, consists of training a network for generating entire sequences, and so the generator model learned may not accurately generate conditional output sequences from input sequences. Another model, described in [87], deals with this task by generating, from partially observed sequences, the future events of those sequences conditioned on their history, i.e., instead of aiming to capture the underlying distribution of a set of full sequences, the model performs a “sample agnostic” in-sample prediction.

Analogously to one of the parametric models described in section 3, the learning procedure of this in-sample neural network based prediction model takes advantage of both types of divergence measures: MLE loss (or KL-divergence) and Wasserstein distance.

The former aims for a rigid and unbiased parameter matching between two given probabilistic distributions, which is sensitive to noisy samples and outliers, while the latter has biased parameter updates but is sensitive to underlying geometrical discrepancies among sample distributions. This combined loss is a way to balance both sets of priorities. In the case of long-term predictions, in which initial prediction errors propagate and magnify themselves throughout the whole stream, this joint loss was found to strengthen the effectiveness of the inference procedure.

The proposed model borrows from the seq2seq architecture [77] and aims to model endings of individual sequences conditioned on their partially observed history of initial events, inserting an adversarial component in the training to increase the accuracy of long-term predictions. A network, designated as generator, encodes a compact representation of the initial partial observation of the sequence and outputs a decoded remainder of this same sequence. That is, for a full sequence,

$$(5.33) \quad \{t_1, t_2, \dots, t_{n+m}\},$$

the seq2seq modeling approach learns a mapping

$$(5.34) \quad G_{\Theta}(\mathcal{S}^{1,n}) = \mathcal{S}^{m,n}$$

such that

$$(5.35) \quad \mathcal{S}^{1,n} = \{t_1, t_2, \dots, t_n\} \text{ and } \mathcal{S}^{m,n} = \{t_{n+1}, t_{n+2}, \dots, t_{n+m}\}.$$

This mapping is defined through a composition of RNN cells,

$$(5.36) \quad \mathbf{h}_i = \eta_g^h(f_A^h(t_i, \mathbf{h}_{i-1})) \text{ and } t_{i+1} = \eta_g^x(f_g^x(\mathbf{h}_i)),$$

with the  $\eta$ 's defined as nonlinear activation functions and the  $f$ 's as linear transformations with trainable weight matrices and bias vectors.

The learning procedure consists of tuning the RNN cells' parameters to maximize the conditional probability

$$(5.37) \quad P(\mathcal{S}^{m,n} | \mathcal{S}^{1,n}) = \prod_{i=n}^{n+m-1} P(t_{i+1} | h_i, t_1, \dots, t_i),$$

which is carried out by parameter gradient updates over the combined Wasserstein and MLE losses. The adversarial component of the training, the discriminator  $F_w^{S2S}(\cdot)$ , is modeled as a residual convolutional network (see [31]) with a 1-Lipschitz constraint, which is related to the magnitude of the gradients of the discriminative model. The full optimization problem then becomes

$$(5.38) \quad \min_{\theta} \max_w \underbrace{\sum_{l=1}^M F_w^{S2S}(\{\mathcal{S}_l^{1,n}, \mathcal{S}_l^{m,n}\}) - \sum_{l=1}^M F_w^{S2S}(\{\mathcal{S}_l^{1,n}, G_{\Theta}(\mathcal{S}_l^{1,n})\})}_{\text{Wasserstein loss}}$$

$$(5.39) \quad - \underbrace{\gamma_{LIP} \left| \frac{\partial F_w^{S2S}(\hat{x})}{\partial \hat{x}} - 1 \right|}_{\text{1-Lipschitz constraint}} - \underbrace{\gamma_{MLE} \log(\mathbb{P}_{\theta}(\mathcal{S}^{m,n} | \mathcal{S}^{1,n}))}_{\text{MLE loss}}.$$

Further works on RNN-based modeling of HPs can be found in [64] and [37].

**5.1. Self-Attentive and Transformer Models.** Another improvement for neural network based modeling, proposed in [100], involves a so-called self-attention strategy [81] to improve the accuracy of the resulting network. The  $i$ th event tuple  $(t_i, m_i)$  is embedded as a variable  $x_i$ ,

$$(5.40) \quad \mathbf{z}_i = \mathbf{t}p_m + \mathbf{p}e_{(m_i, t_i)},$$

which simultaneously encodes information about the event mark through

$$(5.41) \quad \mathbf{t}p_m = \mathbf{z}_e^m \mathbb{W}_E,$$

with  $\mathbf{z}_e^m$  a one-hot encoding vector of the mark and  $\mathbb{W}_E$  an embedding matrix, and information about the time interval among consecutive events through a sinusoidal-based positional encoding vector  $\mathbf{p}e_{(m_i, t_i)}$ , with its  $k$ th entry defined as

$$(5.42) \quad pe_{(m_i, t_i)}^k = \sin(\omega_k^i \times i + \omega_k^t \times t_i).$$

From this encoded variable  $\mathbf{x}_i$ , a hidden state  $\mathbf{h}_{u,i}$  is then defined for each category  $u$  of the marks, which captures the influence of all previous events:

$$(5.43) \quad \mathbf{h}_{u,i+1} = \frac{\left( \sum_{j=1}^i f(\mathbf{z}_{i+1}, \mathbf{z}_j) \mathbf{g}(\mathbf{z}_j) \right)}{\sum_{j=1}^i f(\mathbf{z}_{i+1}, \mathbf{z}_j)}.$$

Through a series of nonlinear transformations, the intensity  $\lambda_u(t)$  for the  $u$ th mark is then computed. A concurrently developed approach in [107] uses multiple attention layers to build a so-called transformer HP, which also surpasses the performance of RNN-based approaches in a series of datasets, as shown in Table 3.

**5.2. Graph Convolutional Networks.** A further improvement of neural HP models, described in [73], involves the graph properties of multivariate HPs, which may be embedded in a neural network modeling framework through the recently proposed graph convolutional networks (GCNs) [41].

The method is composed of a GCN module for capturing meaningful correlation patterns in a large set of event sequences, followed by a usual RNN module for modeling the temporal dynamics. In short, the time sequences are modeled as HPs,

**Table 3** Performance comparison of neural network based HP models: (a) loglikelihood averaged per number of events; (b) RMSE of predicted time interval; and (c) accuracy of mark prediction. The performance was measured over sequences from Retweet [103], MemeTracker [45], Financial [19], Medical Records [36], and Stack Overflow [45] datasets. The TSES method is a likelihood-free model, and so its entries are not evaluated for the Loglikelihood per Event section of the table. Values are obtained from [107].

Method \ Dataset	(a) Loglikelihood per event					(b) Time prediction RMSE			(c) Mark prediction accuracy		
	RT	MT	FIN	MIMIC-II	SO	FIN	MIMIC-II	SO	FIN	MIMIC-II	SO
RMTTPP [19]	-5.99	-6.04	-3.89	-1.35	-2.60	1.56	6.12	9.78	61.95	81.2	45.9
Neural HP [56]	-5.60	-6.23	-3.60	-1.38	-2.55	1.56	6.13	9.83	62.20	83.2	46.3
TSES [89]	-	-	-	-	-	1.50	4.70	8.00	62.17	83.0	46.2
Self-attentive HP [100]	-4.56	-	-	-0.52	-1.86	-	3.89	5.57	-	-	-
Transformer HP [107]	-2.04	0.68	-1.11	0.82	0.04	0.93	0.82	4.99	62.64	85.3	47.0

and the adjacencies among different processes are encoded as a graph. The novel (GCN+RNN) model is meant to extract significant local patterns from the graph. The output of the initial GCN network module, which is simply a matrix of the form  $\chi = [\mu, \mathbb{A}]$  that includes the baseline vector  $\mu$  and the adjacency matrix  $\mathbb{A}$ , is fed into the RNN-based module, there taken as an LSTM. Then the output of this RNN module is input into a further module, a fully connected layer, for calculating the changes  $d\mathbb{X}$  to be applied to the current parameter matrix  $\mathbb{X}$ . Then, after each training step  $T$ , the predicted value of this parameter matrix becomes

$$\mathbb{X}(T) = \mathbb{X}(T - 1) + d\mathbb{X}(T - 1).$$

The work in [95] proposes a model for check-in time prediction composed of an LSTM-based module in which the feature vector is designed to capture each relevant aspect of the problem: the event time coordinate  $t_i$ , an additional field to indicate whether the check-in occurred on a weekday or during the weekend, the Euclidean distance between a given check-in location and the location (l) of the previous check-in, the location type of the check-in (e.g., hotel, restaurant, etc.), the number of users overlapping with a given location, and the check-ins by friends of the user. This aggregates social, geographical, and temporal information in a single neural network based HP-like predictive model.

All these variants of neural network point process models allow for more flexible (nonlinear) representation of the effect of past events on future ones, besides putting at the inference procedure’s disposal a myriad of deep learning tools and techniques which have enjoyed a surge in popularity over recent years.

**6. Further Approaches.** In this section, we briefly review some recently proposed approaches which do not fit conveniently into any of the three previously discussed subgroups, but may be considered as bridges between the usual HP-related tasks and other mathematical subfields.

**6.1. Sparse Gaussian Processes.** By building on the modeling in [102], which considers the triggering function of the HP as a Gaussian process, the work in [101] proposes an approach involving sparse Gaussian processes for optimizing over a dataset. In this, the optimization of the likelihood is taken not over the samples, but over a set of much fewer so-called inducing points, which are also taken as latent variables, in order to result in a final model which is both expressive enough to capture the complexity of the dataset but also tractable enough to be useful and applicable to reasonably sized datasets.

**6.2. Stochastic Differential Equation.** Another way of modeling HPs is through a stochastic differential equation, as proposed in [43], in which the decay of the triggering kernel is taken as exponential, but its amplitude is defined as a stochastic process, taken as being either a geometric Brownian motion or an exponentiated version of the Langevin dynamics.

**6.3. Graph Properties.** Besides some previously discussed works which deal with the properties of the excitation matrices of multivariate HPs as terms to be optimized jointly with other parameters, such as the graph convolutional approaches and the sparsity inducing penalization terms of some parametric and nonparametric approaches, there have been several other optimization strategies taking into account other properties of these matrices.

The work [49] explicitly inserts considerations of the excitation matrix as a distribution over some types of randomly generated matrices into the optimization of the HP likelihood. [51] introduces a penalization term involving the proximity of the excitation matrix to a so-called connection matrix, defined to capture the underlying connectivity among the nodes of the multivariate HP, to the parameter optimization strategy. [53] introduces a weighted sum of the Wasserstein discrepancy and the so-called Gromov–Wasserstein discrepancy as a penalizing factor on the usual MLE procedure of HP estimation, with the intention of inducing both absolute and relational aspects among the nodes of the HP. [4] introduces, besides the sparsity inducing term, another term related to the resulting rank of the excitation matrix that is designed to induce resulting matrices composed of both few nonzero entries and also few independent rows.

**6.4. Epidemic HPs.** The work in [70] blends the HP excitation effect with traditional epidemic models over populations. By considering a time event as an infection, it models the diffusion of a disease by introducing an HP intensity function which is modulated by the size of the available population,

$$(6.1) \quad \lambda(t) = \left(1 - \frac{N_t}{\tilde{N}}\right) \left\{ \mu + \sum_{t_i < t} \phi(t - t_i) \right\},$$

where  $N_t$  denotes the counting process associated with the HP, while  $\tilde{N}$  is the total finite population size.

**6.5. Popularity Prediction.** The work in [58] proposes the use of HP modeling blended with other machine learning techniques, such as random forests, to obtain an associated so-called popularity measure, which is defined as the total number of events the underlying process is expected to generate as  $t \rightarrow \infty$ . This measure is treated as an outcome derived from features associated with some entity (e.g., social network user), such as number of friends, total number of posted statuses, and the account creation time.

In the next section, we will discuss models in which one not only wants to capture the temporal dynamics, but also wishes to influence it toward a certain goal, implicitly defined through a so-called reward function.

**7. Stochastic Control and Reinforcement Learning of HPs.** In this section, we briefly review some control strategies regarding HPs. In some cases, one may wish not only to be able to model the traces of some event sequences, or to capture the underlying distribution of said sequences, but also to try to influence their temporal

dynamics toward more advantageous ones. These cases are considered in work on stochastic control and reinforcement learning approaches for HPs.

This concept of advantageous dynamics is explained through the definition of a so-called reward function, which is defined in terms of specific, and sometimes application-specific, properties of the sequences. Most works related to the subject deal with social network applications, and one example of reward is the total time the post of a user stays at the top of the feeds of his/her followers. One type of reward which is not domain-specific is the dissimilarity among two sets of sequences, computed through mappings such as “kernel mean embeddings” [61]. In the case of imitation learning approaches, one still focuses on modeling only the HP, without steering it toward desirable behaviors. In these approaches, the reward function is simply defined by how well the samples of the model chosen to be adjusted approximate the samples of the original HP.

**7.1. Stochastic Optimal Control (SOC).** One example of this control approach is described in [99], in which the variable to be controlled is the time to post of a given user, implicitly defined as an intensity function, so as to maximize the reward function  $\mathbf{r}(t)$ , here computed as the total time that this user’s posts stay at the top of the feeds of his/her followers.

This “when-to-post” problem can be formulated as

$$\min_{u(t_0, t_f)} \mathbb{E}_{(N_i, M_i)(t_0, t_f)} \left[ \Omega(\mathbf{r}(t_f)) + \int_{t_0}^{t_f} \mathbf{L}(\mathbf{r}(\tau), u(\tau)) d\tau \right]$$

(7.1) subject to  $u(t) \geq 0 \quad \forall t \in (t_0, t_f]$ ,

where

- $i$  is the index of the broadcaster of the posts;
- $N_i(t)$  is the counting process of the  $i$ th broadcaster, with  $\mathbf{N}(t) = \{N_i(t)\}_{i=1}^n$  being an array of counting processes along all the  $n$  users of the network;
- $\mathbb{A} \in \{0, 1\}^{n \times n}$  is the adjacency matrix of the network;
- $\mathbf{M}_i(t) = \mathbb{A}^T \mathbf{N}(t) - \mathbb{A}_i N_i(t)$ , which means that  $M_i(t)$  is the sum of the counting processes of all users connected to user  $i$ , excluding user  $i$  him/herself;
- $t_0$  and  $t_f$  are, respectively, the starting and ending times of the problem horizon taken into consideration;
- $u(t) = \mu_i(t)$ , the controlled variable, is the baseline intensity of user  $i$ , to be steered toward the maximization of the reward function;
- $\Omega(\mathbf{r}(t_f))$  is an arbitrarily defined penalty function;
- $\mathbf{L}(\mathbf{r}(\tau), u(\tau))$  is a nondecreasing convex loss function defined w.r.t. the visibility of the broadcaster’s posts in each of his/her followers’ feeds.

The approach used in the problem is to define an optimal cost-to-go  $J(\mathbf{r}(t_f), \lambda(t), t)$ ,

(7.2) 
$$J(\mathbf{r}(t_f), \lambda(t), t) = \min_{u(t, t_f)} \mathbb{E}_{(N, M)(t, t_f)} \left[ \phi(\mathbf{r}(t_f)) + \int_t^{t_f} \mathbf{l}(\mathbf{r}(\tau), u(\tau)) d\tau \right],$$

and find the optimal solution through Bellman’s principle of optimality:

$$J(\mathbf{r}(t), \lambda(t), t) = \min_{u(t, t+dt)} \{ \mathbb{E}[J(\mathbf{r}(t+dt), \lambda(t+dt), t+dt)] + \mathbf{l}(\mathbf{r}(t), u(t)) dt \}.$$

For example, in the case of a broadcaster with one follower ( $\mathbf{r}(t) = r(t)$ ), if the penalty and loss functions are defined as

(7.3) 
$$\phi(r(t_f)) = \frac{1}{2} r^2(t_f)$$

and

$$(7.4) \quad \mathbf{L}(r(t), u(t)) = \frac{1}{2}s(t)r^2(t) + \frac{1}{2}qu^2(t)$$

for some positive significance function  $s(t)$  and some trade-off parameter  $q$ , which calibrates the importance of both visibility and number of posts, we set the derivative of  $J(r(t), \lambda(t), t)$  over  $u(t)$  to 0 and solve it to obtain the analytical solution

$$(7.5) \quad \mathbf{u}^*(t) = q^{-1}[J(r(t), \lambda(t), t) - J(0, \lambda(t), t)],$$

which is thus the optimal intensity a broadcaster must adopt to maximize visibility, constrained on the cost associated to the number of posts, along this follower's feed. Further derivations are provided for the more natural and general case, in which the broadcaster may have multiple followers.

An earlier version of this type of SOC-based approach to influencing activity in social networks can be found in [98], in which the goal is to maximize the total number of actions (or events) in the network. Analogously to the previously discussed algorithm, one may solve the continuous time version of the Bellman equation by defining an optimal cost-to-go  $J(\boldsymbol{\lambda}(t), t)$ , which here depends only on the intensities of the nodes and the time.

The control input vector  $\mathbf{u}(t)$  acts on the network by increasing the original vector of uncontrolled intensities

$$(7.6) \quad \boldsymbol{\lambda}(t) = \boldsymbol{\mu}_0 + \mathbb{A} \int_0^t \kappa(t-s)d\mathbf{N}(s),$$

with the equivalent rates of an underlying counting process vector  $d\mathbf{M}(s)$ , such that the new controlled intensity vector  $\boldsymbol{\lambda}^*(t)$  is now described by

$$(7.7) \quad \boldsymbol{\lambda}^*(t) = \boldsymbol{\mu}_0 + \mathbb{A} \int_0^t \kappa(t-s)d\mathbf{N}(s) + \mathbb{A} \int_0^t \kappa(t-s)d\mathbf{M}(s),$$

where  $\kappa(t) = e^{-\beta t}$  in the model.

Then, in the same way, by differentiating the equivalent  $J(\boldsymbol{\lambda}(t), t)$  over the control input, setting the corresponding expression to 0, and then defining

$$(7.8) \quad \mathbf{L}(\boldsymbol{\lambda}(t), \mathbf{u}(t)) = -\frac{1}{2}\boldsymbol{\lambda}^T(t)\mathbf{Q}\boldsymbol{\lambda}(t) + \frac{1}{2}\mathbf{u}^T(t)\mathbf{S}\mathbf{u}(t)$$

and

$$(7.9) \quad \Omega(\boldsymbol{\lambda}(t_f)) = -\frac{1}{2}\boldsymbol{\lambda}^T(t_f)\mathbf{F}\boldsymbol{\lambda}(t_f),$$

with previously defined symmetric weighting matrices  $\mathbf{Q}$ ,  $\mathbf{F}$ , and  $\mathbf{S}$ , we arrive at a closed-form expression for the optimal control intensity value,

$$(7.10) \quad \mathbf{u}^*(t) = -\mathbf{S}^{-1} \left[ \mathbb{A}^T \mathbf{g}(t) + \mathbb{A}^T \mathbf{H}(t) \boldsymbol{\lambda}(t) + \frac{1}{2} \text{diag}(\mathbb{A}^T \mathbf{H}(t) \mathbb{A}) \right].$$

$\mathbf{H}(t)$  and  $\mathbf{g}(t)$  can be computed by solving the differential equations

$$(7.11) \quad \dot{\mathbf{H}}(t) = (\beta \mathbf{I} - \mathbb{A})^T \mathbf{H}(t) + \mathbf{H}(t)(\beta \mathbf{I} - \mathbb{A}) + \mathbf{H}(t) \mathbb{A} \mathbf{S}^{-1} \mathbb{A}^T \mathbf{H}(t) \mathbf{Q},$$



$$(7.12) \quad \begin{aligned} \dot{\mathbf{g}}(t) &= [\beta \mathbf{I} - \mathbb{A}^T + \mathbf{H}(t) \mathbb{A} \mathbf{S}^{-1} \mathbb{A}^T] \mathbf{g}(t) - \beta \mathbf{H}(t) \boldsymbol{\mu}_0 \\ &+ \frac{1}{2} [\mathbf{H}(t) \mathbb{A} \mathbf{S}^{-1} - \mathbf{I}] \text{diag}(\mathbb{A}^T \mathbf{H}(t) \mathbb{A}), \end{aligned}$$

with final conditions  $\mathbf{g}(t_f) = 0$  and  $\mathbf{H}(t_f) = -\mathbf{F}$ . The solution is constant between two consecutive events and is recomputed at each event arrival.

**7.2. Reinforcement Learning.** The previously described SOC-based approaches have two main drawbacks:

- The functional forms of the intensities and mark distributions are constrained to be from a very restricted class, which does not include the state-of-the-art RNN-based HP models, such as those described in section 5.
- The objective function being optimized is also restricted to very specific classes of functions, so as to maintain the tractability of the problem.

To circumvent these drawbacks, some approaches have been proposed which combine more flexible and expressive HP models with robust stochastic optimization procedures independent from the functional form of the objective function.

One of these methods, called “deep reinforcement learning of marked temporal point processes” [80], considers a given set of possible actions and corresponding feedbacks, which are both expressed as temporal point processes jointly modeled by an RNN-based intensity model  $\lambda_{\theta}^*(t)$ ,

$$(7.13) \quad \lambda_{\theta}^*(t) = \exp(b_{\lambda} + w_t(t - t_i) + \mathbf{V}_{\lambda} \mathbf{h}_i),$$

with

$$\mathbf{h}_i = \tanh(\mathbf{W}_h \mathbf{h}_{i-1} + \mathbf{W}_1 \mathcal{T}_i + \mathbf{W}_2 \mathbf{y}_i + \mathbf{W}_3 \mathbf{z}_i \mathbf{W}_4 \mathbf{b}_i + \mathbf{b}_h),$$

where

$$\mathcal{T}_i = f_T(t_i - t_{i-1}) \text{ and } \mathbf{b}_i = f_b(1 - b_i, b_i),$$

and

$$\mathbf{y}_i = f_y(y_i) \text{ if } b_i = 0, \text{ and } \mathbf{z}_i = f_z(z_i) \text{ if } b_i = 1.$$

The term  $b_i$  is an indicator function to whether the  $i$ th event is an action or a feedback. By taking the weight matrices and bias vectors from all the linear transformations  $f$  of a parameter vector  $\theta$ , the algorithm wishes to update this vector with the gradients of each parameter over an expected reward function  $J(\theta)$ :

$$(7.14) \quad \theta_{l+1} = \theta_l + \eta_l \nabla_{\theta} J(\theta)|_{\theta=\theta_l},$$

$$(7.15) \quad \nabla_{\theta} J(\theta) = \mathbb{E}_{\mathcal{U}_T} p_{\mathcal{U}, \theta}^*(\cdot), \mathcal{F}_T} p_{\mathcal{F}, \phi}^*(\cdot) [R^*(T) \nabla_{\theta} \log \mathbb{P}_{\theta}(\mathcal{U}_T)],$$

where

$$(7.16) \quad p_{\mathcal{U}, \theta}^* = (\lambda_{\theta}^*, m_{\theta}^*)$$

is the joint conditional intensity and mark distribution for action events, and

$$(7.17) \quad p_{\mathcal{F}, \phi}^* = (\lambda_{\phi}^*, m_{\phi}^*)$$

is the joint conditional intensity and mark distribution for feedback events. The reward  $R(T)$  is defined over some domain-specific metric, which may involve the responsiveness of followers in a social network setting, or the effectiveness of memorization of words in a foreign language, in a spaced repetition learning setting.

**7.3. Imitation Learning.** Another way of using this reinforcement learning approach is through a technique called “imitation learning” [47]. The reasoning behind it is to treat the real-data sequences as having been generated by an expert and then, using RNN-based sequence generation, try to make these models approximate the real-data sequence as closely as possible.

Thus, the reward function that dictates the proportion by which the gradients of the parameters over each sequence are going to be considered is equal to how likely this given sequence is to be drawn from the underlying distribution over the real data. This similarity is computed through an RKHS.

The theory of the RKHS is very extensive, and it is not our goal to give a detailed account of it here. The key idea is that, to compute similarities among items of a given space, you compute inner products between them. Taking two items,  $x_1$  and  $x_2$ , and computing a so-called positive definite kernel (PDK)  $K(x_1, x_2)$  is equivalent to computing an inner product among these two items in a high-dimensional, and potentially infinite-dimensional, vector space. The PDK used in the paper mentioned above is the Gaussian kernel.

The reward function is then defined as

$$(7.18) \quad \hat{r}^*(t) \propto \frac{1}{L} \sum_{l=1}^L \sum_{i=1}^{N_T^{(l)}} \mathbb{K}(s_i^{(l)}, t) - \frac{1}{M} \sum_{m=1}^M \sum_{i=1}^{N_T^m} \mathbb{K}(t_i^{(m)}, t),$$

where

- $L$  is the number of expert trajectories;
- $M$  is the number of trajectories generated by the model;
- $\mathbb{K}(\cdot, \cdot)$  is the reproducing kernel operator;
- $s_i^{(l)}$  is the  $i$ th time coordinate of the  $l$ th expert trajectory;
- $t_i^{(m)}$  is the  $i$ th time coordinate of the  $m$ th model-generated trajectory.

The parameters of the model are then updated through a gradient descent based approach toward convergence, in which the model is expected to generate sequences indistinguishable from the real-world process.

**8. Real-World Data Limitations.** One key aspect of HP modeling which inevitably encompasses all the previously mentioned approaches is that of their applicability to real-world datasets. These may present a series of systematic issues on the training and testing sequences that could entirely hinder the generalization of the models. We now discuss some key issues, along with some recently proposed methods designed to handle each of them.

**8.1. Synchronization Noise.** Temporal data, especially multivariate data, may have event streams extracted from distributed sensor networks. A key challenge is that of synchronization noise, i.e., when each source is subject to an unknown and random time delay.

In this case, an inference procedure which neglects these time delays may ignore critical causal effects of some events on others, thus resulting in a poorly generalized model. The work in [78] deals specifically with this aspect of real-data HP modeling, and it proposes, for the exponential triggering functions HPs, the inclusion of the random time shift vector (one entry for each distinct event stream) as a parameter in the HP model, which results in an inference procedure of the following form:

$$(8.1) \quad \hat{z}, \hat{\theta} = \underset{z \in \mathbb{R}, \theta \geq 0}{\operatorname{argmax}} \log \mathbb{P}(\tilde{\mathbf{t}} | \mathcal{N}, \theta),$$

where  $\mathcal{N}$  is the random noise vector, and  $\theta$  corresponds to the parameters of the original exponential HP model.

**8.2. Sequences with Few Events.** In many domains, available data is scarce, and the event streams will be composed of too few events, which results in noisiness of the likelihood and, as a consequence, unreliability of the fitted HP model, which calls for strong regularization strategies over the objective functions to be optimized.

For this type of situation, one approach, presented in [72], deals with HPs with triggering functions defined as exponentials and also as a mixture of Gaussian kernels, as in [91]. Then the parameters left to search are the background rates vector  $\mu$  and the tensor of weightings  $\mathbf{A}$  for the excitation functions. The optimization is done through a variational expectation-maximization algorithm, which takes the distributions over these parameters as Gaussians and optimizes, through Monte Carlo sampling, over an evidence lower bound (ELBO) of their corresponding loglikelihoods over a set of sequences.

**8.3. Sequences with Missing Data.** Another issue in HP modeling concerns learning from incomplete sequences, i.e., streams in which one or more of the events are missing. For this type of problem, two rather distinct approaches were recently proposed:

1. The first approach, presented in [76], is applied to exponential and power-law HPs and consists of a Markov chain Monte Carlo based inference over a joint process implicitly defined by the product of the likelihoods of the observed events and of the so-called virtual event auxiliary variables, which are candidates for unobserved events. This virtual variable is weighted through a parameter  $\kappa$ , which is related to the percentage of missing events w.r.t. the total event count.
2. The second approach, introduced in [57], proposes finding the missing events over the sequences through importance weighting of candidate filling event subsequences generated by a bidirectional LSTM model built on top of the neural HP [56].

**9. Application Examples.** In this section, we apply the HP modeling learned so far to case studies for three different domains: retweeting behavior in social networks, earthquake aftershocks, and malaria outbreak forecasting. We also make a few remarks about the use of HPs in financial modeling. We hope this will encourage the reader to consider HPs as a modeling choice for a broad scope of applications.

**9.1. Retweet.** In [21], HPs are used jointly with latent Dirichlet allocation (LDA) [10] models for distinguishing between genuine and fake (i.e., artificially induced) retweeting of posts among Twitter users.

Consider a set of 2508 users, with each  $j$ th user corresponding to a sequence

$$(9.1) \quad \mathbf{RT}^j = \{(t_i^j, \mathcal{W}_i^j)\}_{i=0}^{N_j},$$

where  $t_i^j$  corresponds to the timestamp associated with the  $i$ th retweet from the  $j$ th user, and  $\mathcal{W}_i^j$  is the text content of the corresponding retweet.

The 2508 corresponding sequences are manually separated between *genuine* and *fake* users and labeled as such, based on a set of criteria (e.g., the content in most of said user's retweets contains spammy links and common spam keywords; multiple retweets from a given user contain promotional/irrelevant text; the user's biographical information is fabricated or contains promotional activity; or a large number of tweets or retweets are posted within a very short time window of just a few seconds).

The two resulting disjoint sets were used for training two LDA models,  $LDA_f$  and  $LDA_g$ , for modeling the topics of fake and genuine retweeters, respectively. Given a predefined number of possible topics, here set as 10, and a given text content  $\mathcal{W}_i^j$ , the LDA model outputs a 10-element vector with the probabilities of the  $\mathcal{W}_i^j$  corresponding to each of the 10 possible topics.

The 10-element vector  $\mathcal{V}_f$  from  $LDA_f$  is concatenated with the 10-element vector  $\mathcal{V}_g$  from  $LDA_g$  and, together with the baseline intensity  $\mu$  and the decay  $\beta$  of an exponential HP with  $\phi(t) = e^{-\beta t}$ , fitted over the  $t_i^j$ 's of each  $j$ th sequence, it forms a feature vector

$$(9.2) \quad \{\mathcal{V}_f^j, \mathcal{V}_g^j, \mu^j, \beta^j\},$$

which is then fed into a clustering algorithm that aims to correctly classify each of the 2508 retweet sequences between fake and genuine. The intuition behind this hybrid HP-LDA model is to use temporal features from the HP modeling together with context (written) features from the users to improve the resulting detection algorithm.

**9.2. Earthquake Aftershocks.** It is well known from the study of earthquake-related time series that a strong first seismic shock gives way to a series of weaker aftershocks, which occur in a very restricted time window [63].

For modeling this self-exciting property of the aftershock arrivals, [63] proposes a power-law self-triggering kernel

$$(9.3) \quad \phi_{PWL}(t, \theta_{PWL}) = \frac{K}{(t+c)^p},$$

with  $\theta_{PWL} = (K, c, p) \in \mathbb{R}_+^3$  as trainable parameters. This model, together with an additional baseline rate parameter  $\mu$ , is fitted over the temporal sequence  $\{t_1, t_2, \dots, t_n\}$  of aftershock timestamps with an MLE optimization procedure, such as the one in (2.17), in which, due to the simple parametric form of the equation for the intensity, the loglikelihood can be given in closed form.

**9.3. Malaria Outbreak Forecasting.** The work in [79] proposes one possible way of modeling the diffusion of malaria cases within a given population by assuming it behaves as an HP of time-dependent background rate

$$(9.4) \quad \mu(t) = \max \left( \mathcal{E} + \mathcal{J}t + \mathcal{Q} \cos \left( \frac{2\pi t}{\mathcal{X}} \right) + \mathcal{Z} \sin \left( \frac{2\pi t}{\mathcal{X}} \right), 0 \right),$$

with

$$(9.5) \quad \mathcal{X} = 365.25$$

and constants  $\mathcal{E}, \mathcal{J}, \mathcal{Q}$ , to account for the yearly seasonality of imported cases, as well as a Rayleigh-type triggering kernel

$$(9.6) \quad \phi(t - (t_i + \mathcal{D})) = v(t - (t_i + \mathcal{D}))e^{\frac{-\varrho(t - (t_i + \mathcal{D}))^2}{2}}$$

with

$$(9.7) \quad t > t_i + \mathcal{D} \quad \text{and} \quad v, \varrho \geq 0,$$

with an additional delay term  $\mathcal{D}$ .

This non-strictly-decaying choice of  $\phi(t)$  is due to the fact that a person is not most infectious right after being bitten by the disease-transmitting mosquito. The delay term accounts for the incubating period of the bitten person, until they become infectious.

The defined parameters are fitted to streams of reported cases from China and Eswatini through a modified MLE strategy, to account for the resulting nonconvexity of the underlying loglikelihood.

**9.4. Financial Modeling.** The suitability of HPs to model sudden jumps in continuous time, without the need for time discretization strategies, has made them broadly applicable to a number of intraday high-frequency financial applications, such as in the following works:

- [3] applies HPs with exponential kernels to study the self- and mutually exciting effects among credit default swap market shocks in several European countries.
- [22] proposes a generalized version of HP which accounts for different opening hours of the markets due to time zone differences. It goes on to study the self- and mutually exciting behavior of price jumps in the S&P 500 and the Euro Stoxx 50.
- [1] uses a multivariate HP with exponential kernels to model the arrival of trades and cancellations of a limit order book.

For a comprehensive treatment of HPs in finance, the reader should refer to two excellent survey works focused specifically on this topic [30, 7].

**10. Comparisons with Other Temporal Point Process Approaches.** HPs, and the simpler Poisson processes, have been the most prevalent choice for modeling time event sequences, but other approaches have been proposed that occasionally surpass the performance of HPs in some situations:

- Wold processes: These are the equivalent of an HP in which only the effect of the most recent event is considered in the computation of the intensity function. This Markovian aspect, regardless of the choice of the excitation function, has been shown in [23] to surpass the performance of several HP models for the estimation of networked processes.
- Intensity-free learning of interevent intervals: Another approach which has been recently introduced involves ignoring the intensity function completely and focusing on modeling the probabilistic distribution of the time intervals among consecutive events. This distribution is modeled using normalizing flows [68], which can be summed up as families of distributions with incremental complexities. The approach was introduced in [74, 75] and was shown to surpass state-of-the-art neural-based HP models in some large-sized datasets.
- Continuous-time Markov chain: In this model, the marks correspond to states which have fixed rates (intensities) associated to them. The transition times are sampled from these constant intensities. It has been used as a comparison baseline for some HP models, such as [20].

**11. Current Challenges for Further Research.** We mention the following challenges currently tackled by HP researchers:

- Enriching HP variants (parametric, nonparametric, neural) or blending them with other machine learning approaches, so as to make them suitable for spe-

cific situations. Work with multiarmed bandits [15], randomized kernels [35], graph neural networks for temporal knowledge graphs [27], and composition of HP-like point processes with warping functions defined over the time event sequences [90] fall into this category.

- Improving the speed of inference or sampling to reduce the time spent in model estimation, an aspect which may be critical for some real-world applications. The works of [34] in Bayesian mitigation of spatial coarsening, [104] in multiresolution segmentation for nonstationary HPs using cumulants, [48] on thinning of event sequences for accelerating inference steps, [54] on the use of Lambert-W functions for improving sequence sampling, and [13] on perfect sampling are examples of this aspect, as well as [65] on recursive computation of HP moments.
- How to properly evaluate and compare HP models: While there has been a lot of work proposing new approaches, the comparison among existing models is often biased or incomplete. The works of [82] on how to quantify the uncertainty of the obtained models, [85] on measuring goodness-of-fit, [55] on robust identification of HPs with controlled terms, and [11] on the rigorous comparison of networked point process models address this type of challenge.
- Theoretical guarantees, properties, and formulations of specific HP approaches, such as work done in [14] on strong mixing, [26] on the consistency of some parametric models, [16] on elementary derivations of HP momenta, and [39, 40] on field master equation formulation for HPs.

**12. Conclusions.** Hawkes processes are a valuable tool for modeling a myriad of natural and social phenomena. The present work has aimed to give a broad view, suitable for a newcomer to the field, of the inference and modeling techniques involved in the application of HPs in a variety of domains. The parametric, nonparametric, deep learning, and reinforcement learning approaches were broadly covered, as well as the current research challenges on the topic and the real-world limitations of each approach. Illustrative application examples in the modeling of retweeting behavior, earthquake aftershock occurrence, and malaria outbreak modeling were also briefly discussed, to motivate the applicability of HPs in both natural and social phenomena.

**Acknowledgments.** The author would like to thank Thanh Nguyen-Tang, as well as the anonymous reviewers, for constructive comments on earlier versions of this work.

#### REFERENCES

- [1] F. ABERGEL AND A. JEDIDI, *Long-time behavior of a Hawkes process-based limit order book*, SIAM J. Financial Math., 6 (2015), pp. 1026–1043, <https://doi.org/10.1137/15M1011469>. (Cited on p. 367)
- [2] M. ACHAB, E. BACRY, S. GAÏFFAS, I. MASTROMATTEO, AND J. MUZY, *Uncovering causality from multivariate Hawkes integrated cumulants*, J. Mach. Learn. Res., 18 (2017), pp. 192:1–192:28. (Cited on pp. 333, 337, 347, 350)
- [3] Y. AIT-SAHALIA, R. LAEVEN, AND L. PELIZZON, *Mutual excitation in Eurozone sovereign CDS*, J. Econometrics, 183 (2014), pp. 151–167. (Cited on p. 367)
- [4] E. BACRY, M. BOMPAIRE, S. GAÏFFAS, AND J.-F. MUZY, *Sparse and low-rank multivariate Hawkes processes*, J. Mach. Learn. Res., 21 (2020), art. 50. (Cited on p. 360)
- [5] E. BACRY, S. DELATTRE, M. HOFFMANN, AND J. MUZY, *Scaling limits for Hawkes processes and application to financial statistics*, Stoch. Process. Appl., 123 (2013), pp. 2475–2499. (Cited on p. 347)

- [6] E. BACRY, S. GAÏFFAS, I. MASTROMATTEO, AND J. MUZY, *Mean-field inference of Hawkes point processes*, J. Phys. A, 49 (2016), art. 174006. (Cited on p. 343)
- [7] E. BACRY, I. MASTROMATTEO, AND J.-F. MUZY, *Hawkes processes in finance*, Market Microstructure Liquidity, 1 (2015), art. 1550005. (Cited on pp. 333, 367)
- [8] E. BACRY AND J. MUZY, *First- and second-order statistics characterization of Hawkes processes and non-parametric estimation*, IEEE Trans. Inform. Theory, 62 (2016), pp. 2184–2202. (Cited on pp. 333, 347, 350)
- [9] T. BJÖRK, *An Introduction to Point Processes from a Martingale Point of View*, preprint, 2011. (Cited on p. 338)
- [10] D. M. BLEI, A. Y. NG, M. I. JORDAN, AND J. LAFFERTY, *Latent Dirichlet allocation*, J. Mach. Learn. Res., 3 (2003), pp. 993–1022. (Cited on p. 365)
- [11] G. BORGES, P. O. S. V. DE MELO, F. FIGUEIREDO, AND R. ASSUNCAO, *Networked point process models under the lens of scrutiny*, in Machine Learning and Knowledge Discovery in Databases (ECML PKDD, 2020), Part I, Springer, 2021, pp. 198–215. (Cited on p. 368)
- [12] J. CHEN, A. G. HAWKES, AND E. SCALAS, *A fractional Hawkes process*, in Nonlocal and Fractional Operators, L. Beghin, F. Mainardi, and R. Garrappa, eds., SEMA SIMAI Springer Seri. 26, Springer, 2021, pp. 121–131. (Cited on p. 333)
- [13] X. CHEN AND X. WANG, *Perfect sampling of multivariate hawkes processes*, in Proceedings of the Winter Simulation Conference (WSC '20), IEEE Press, 2021, pp. 469–480. (Cited on p. 368)
- [14] F. CHEYSSON AND G. LANG, *Strong Mixing Condition for Hawkes Processes and Application to Whittle Estimation from Count Data*, preprint, <https://arxiv.org/abs/2003.04314>, 2020. (Cited on p. 368)
- [15] W.-H. CHIANG AND G. MOHLER, *Hawkes Process Multi-armed Bandits for Disaster Search and Rescue*, preprint, <https://arxiv.org/abs/2004.01580>, 2020. (Cited on p. 368)
- [16] L. CUI, A. HAWKES, AND H. YI, *An elementary derivation of moments of Hawkes processes*, Adv. Appl. Probab., 52 (2020), pp. 102–137, <https://doi.org/10.1017/apr.2019.53>. (Cited on p. 368)
- [17] D. DALEY AND D. VERE-JONES, *An Introduction to the Theory of Point Processes, Volume I: Elementary Theory and Methods*, Springer, 2003. (Cited on pp. 332, 337)
- [18] S. DONNET, V. RIVOIRARD, AND J. ROUSSEAU, *Nonparametric Bayesian estimation of multivariate Hawkes processes*, Ann. Statist., 48 (2020), pp. 2698–2727. (Cited on pp. 350, 352)
- [19] N. DU, H. DAI, R. TRIVEDI, U. UPADHYAY, M. GOMEZ-RODRIGUEZ, AND L. SONG, *Recurrent marked temporal point processes: Embedding event history to vector*, in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, 2016, pp. 1555–1564. (Cited on pp. 352, 359)
- [20] N. DU, M. FARAJTABAR, A. AHMED, A. J. SMOLA, AND L. SONG, *Dirichlet-Hawkes processes with applications to clustering continuous-time document streams*, in Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, Australia, 2015, pp. 219–228. (Cited on pp. 333, 350, 367)
- [21] H. S. DUTTA, V. R. DUTTA, A. ADHIKARY, AND T. CHAKRABORTY, *HawkesEye: Detecting fake retweeters using Hawkes process and topic modeling*, IEEE Trans. Inform. Forensics Security, 15 (2020), pp. 2667–2678, <https://doi.org/10.1109/TIFS.2020.2970601>. (Cited on p. 365)
- [22] F. FERRIANI AND P. ZOI, *The dynamics of price jumps in the stock market: An empirical study on Europe and U.S.*, European J. Finance, 28 (2022), pp. 718–742, <https://doi.org/10.1080/1351847X.2020.1740288>. (Cited on p. 367)
- [23] F. FIGUEIREDO, G. R. BORGES, P. O. S. V. DE MELO, AND R. ASSUNÇÃO, *Fast estimation of causal interactions using Wold processes*, in Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, Montreal, Canada, 2018, pp. 2975–2986. (Cited on p. 367)
- [24] S. FLAXMAN, M. CHIRICO, P. PEREIRA, AND C. LOEFFLER, *Scalable high-resolution forecasting of sparse spatiotemporal events with kernel methods: A winning solution to the NIJ “Real-Time Crime Forecasting Challenge,”* Ann. Appl. Stat., 13 (2019), pp. 2564–2585, <https://doi.org/10.1214/19-AOAS1284>. (Cited on pp. 336, 338)
- [25] I. J. GOODFELLOW, J. POUGET-ABADIE, M. MIRZA, B. XU, D. WARDE-FARLEY, S. OZAIR, A. C. COURVILLE, AND Y. BENGIO, *Generative adversarial nets*, in Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, NeurIPS 2014, Montreal, Canada, 2014, pp. 2672–2680. (Cited on p. 346)
- [26] X. GUO, A. HU, R. XU, AND J. ZHANG, *Consistency and Computation of Regularized MLEs for Multivariate Hawkes Processes*, preprint, <https://arxiv.org/abs/1810.02955>, 2018.

- (Cited on p. 368)
- [27] Z. HAN, Y. MA, Y. WANG, S. GÜNNEMANN, AND V. TRESP, *Graph Hawkes neural network for forecasting on temporal knowledge graphs*, in 8th Automated Knowledge Base Construction (AKBC), 2020. (Cited on p. 368)
- [28] A. G. HAWKES, *Point spectra of some mutually exciting point processes*, J. Roy. Statist. Soc. Ser. B, 33 (1971), pp. 438–443, <http://www.jstor.org/stable/2984686>. (Cited on p. 333)
- [29] A. G. HAWKES, *Spectra of some self-exciting and mutually exciting point processes*, Biometrika, 58 (1971), pp. 83–90. (Cited on pp. 332, 333, 337)
- [30] A. G. HAWKES, *Hawkes processes and their applications to finance: A review*, Quant. Finance, 18 (2018), pp. 193–198. (Cited on pp. 333, 338, 367)
- [31] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, 2016, pp. 770–778. (Cited on p. 358)
- [32] A. HELMSTETTER AND D. SORNETTE, *Diffusion of epicenters of earthquake aftershocks, Omori’s law, and generalized continuous-time random walk models*, Phys. Rev. E, 66 (2002), art. 061104. (Cited on p. 332)
- [33] S. HOCHREITER AND J. SCHMIDHUBER, *Long short-term memory*, Neural Comput., 9 (1997), pp. 1735–1780. (Cited on p. 353)
- [34] A. J. HOLBROOK, X. JI, AND M. A. SUCHARD, *Bayesian mitigation of spatial coarsening for a Hawkes model applied to gunfire, wildfire and viral contagion*, Ann. Appl. Statist., 16 (2022), pp. 573–595. (Cited on p. 368)
- [35] F. ILHAN AND S. S. KOZAT, *Modeling of spatio-temporal Hawkes processes with randomized kernels*, IEEE Trans. Signal Process., 68 (2020), pp. 4946–4958, <https://doi.org/10.1109/tsp.2020.3019329>. (Cited on p. 368)
- [36] A. E. JOHNSON, T. J. POLLARD, L. SHEN, L. H. LEHMAN, M. FENG, M. GHASSEMI, B. MOODY, P. SZOLOVITS, L. A. CELI, AND R. G. MARK, *MIMIC-III, a freely accessible critical care database*, Sci. Data, 3 (2016), art. 160035. (Cited on p. 359)
- [37] S. JOSEPH, L. D. KASHYAP, AND S. JAIN, *Shallow Neural Hawkes: Non-parametric Kernel Estimation for Hawkes Processes*, preprint, <https://arxiv.org/abs/2006.02460>, 2020. (Cited on p. 358)
- [38] A. T. KALAI AND R. SASTRY, *The isotron algorithm: High-dimensional isotonic regression*, in the 22nd Conference on Learning Theory, COLT 2009, Montreal, Canada, 2009. (Cited on p. 341)
- [39] K. KANAZAWA AND D. SORNETTE, *Field master equation theory of the self-excited Hawkes process*, Phys. Rev. Res., 2 (2020), art. 138301, <https://doi.org/10.1103/physrevresearch.2.033442>. (Cited on p. 368)
- [40] K. KANAZAWA AND D. SORNETTE, *Nonuniversal power law distribution of intensities of the self-excited Hawkes process: A field-theoretical approach*, Phys. Rev. Lett., 125 (2020), art. 138301, <https://doi.org/10.1103/PhysRevLett.125.138301>. (Cited on p. 368)
- [41] T. N. KIPF AND M. WELLING, *Semi-supervised classification with graph convolutional networks*, in 5th International Conference on Learning Representations (ICLR 2017), 2017. (Cited on p. 358)
- [42] R. KOBAYASHI AND R. LAMBIOTTE, *TiDeH: Time-dependent Hawkes process for predicting retweet dynamics*, in Proceedings of the Tenth International Conference on Web and Social Media, Cologne, Germany, 2016. (Cited on pp. 333, 339, 340)
- [43] Y. LEE, K. W. LIM, AND C. S. ONG, *Hawkes processes with stochastic excitations*, in Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York, 2016, pp. 79–88. (Cited on p. 360)
- [44] R. LEMMONIER, K. SCAMAN, AND A. KALOGERATOS, *Multivariate Hawkes processes for large-scale inference*, in Proceedings of the Conference on Artificial Intelligence, 2017, pp. 2168–2174. (Cited on pp. 333, 343)
- [45] J. LESKOVEC AND A. KREVL, *SNAP Datasets: Stanford Large Network Dataset Collection*, 2014, <http://snap.stanford.edu/data>. (Cited on pp. 350, 359)
- [46] E. LEWIS AND G. MOHLER, *A nonparametric EM algorithm for multiscale Hawkes processes*, J. Nonparametric Statist., 1 (2011), pp. 1–20. (Cited on p. 347)
- [47] S. LI, S. XIAO, S. ZHU, N. DU, Y. XIE, AND L. SONG, *Learning temporal point processes via reinforcement learning*, in Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, Montreal, Canada, 2018, pp. 10804–10814. (Cited on pp. 333, 364)
- [48] T. LI AND Y. KE, *Thinning for accelerating the learning of point processes*, in Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, Vancouver, Canada, 2019, pp. 4093–4103. (Cited



- on p. 368)
- [49] S. W. LINDERMAN AND R. P. ADAMS, *Discovering latent network structure in point process data*, in Proceedings of the International Conference on Machine Learning, 2014, pp. 1413–1421. (Cited on p. 360)
- [50] T. LINIGER, *Multivariate Hawkes Processes*, Ph.D. thesis, ETH Zurich, 2009. (Cited on p. 337)
- [51] Y. LIU, T. YAN, AND H. CHEN, *Exploiting graph regularized multi-dimensional Hawkes processes for modeling events with spatio-temporal characteristics*, in Proceedings of the 27th International Joint Conference on Artificial Intelligence, AAAI Press, 2018, pp. 2475–2482. (Cited on pp. 336, 360)
- [52] C. E. LOEFFLER AND S. R. FLAXMAN, *Is gun violence contagious? A spatiotemporal test*, *J. Quant. Criminology*, 34 (2018), pp. 999–1017, <https://doi.org/10.1007/s10940-017-9363-8>. (Cited on pp. 336, 338)
- [53] D. LUO, H. XU, AND L. CARIN, *Fused Gromov-Wasserstein alignment for Hawkes processes*, in Learning with Temporal Point Processes, NeurIPS 2019 Workshop, 2019. (Cited on p. 360)
- [54] M. MAGRIS, *On the Simulation of the Hawkes Process via Lambert-W Functions*, preprint, <https://arxiv.org/abs/1907.09162>, 2019. (Cited on p. 368)
- [55] M. MARK AND T. A. WEBER, *Robust identification of controlled Hawkes processes*, *Phys. Rev. E*, 101 (2020), art. 043305, <https://doi.org/10.1103/PhysRevE.101.043305>. (Cited on p. 368)
- [56] H. MEI AND J. EISNER, *The neural Hawkes process: A neurally self-modulating multivariate point process*, in Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, NeurIPS 2017, Long Beach, CA, 2017, pp. 6757–6767. (Cited on pp. 353, 359, 365)
- [57] H. MEI, G. QIN, AND J. EISNER, *Imputing missing events in continuous-time event streams*, in Proceedings of the 36th International Conference on Machine Learning, ICML 2019, Long Beach, CA, Proc. Mach. Learn. Res. 97, PMLR, 2019, pp. 4475–4485. (Cited on p. 365)
- [58] S. MISHRA, M. RIZOIU, AND L. XIE, *Feature driven and point process approaches for popularity prediction*, in Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, ACM, 2016, pp. 1069–1078. (Cited on p. 360)
- [59] G. O. MOHLER, M. B. SHORT, P. J. BRANTINGHAM, F. P. SCHOENBERG, AND G. E. TITA, *Self-exciting point process modelling of crime*, *J. Amer. Statist. Assoc.*, 106 (2012), pp. 100–108. (Cited on pp. 333, 342)
- [60] J. MOLLER AND J. G. RASMUSSEN, *Perfect simulation of Hawkes processes*, *Adv. Appl. Probab.*, 37 (2010), pp. 629–646. (Cited on p. 339)
- [61] K. MUANDET, K. FUKUMIZU, B. SRIPERUMBUDUR, AND B. SCHÖLKOPF, *Kernel mean embedding of distributions: A review and beyond*, *Found. Trends Mach. Learn.*, 10 (2017), pp. 1–141, <https://doi.org/10.1561/22000000060>. (Cited on p. 361)
- [62] Y. OGATA, *On Lewis' simulation method for point processes*, *IEEE Trans. Inform. Theory*, 27 (1981), pp. 23–31. (Cited on pp. 339, 355)
- [63] Y. OGATA, *Seismicity analysis through point-process modelling: A review*, *Pure Appl. Geophys.*, 155 (1999), pp. 471–507. (Cited on pp. 332, 366)
- [64] T. OMI, N. UEDA, AND K. AIHARA, *Fully neural network based model for general temporal point processes*, in Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, Vancouver, Canada, 2019, pp. 2120–2129. (Cited on p. 358)
- [65] N. PRIVAULT, *Recursive computation of the Hawkes cumulants*, *Statist. Probab. Lett.*, 177 (2021), art. 109161. (Cited on p. 368)
- [66] C. E. RASMUSSEN AND C. K. I. WILLIAMS, *Gaussian Processes for Machine Learning*, *Adapt. Comput. Mach. Learn.*, MIT Press, 2006. (Cited on p. 352)
- [67] A. REINHART, *A review of self-exciting spatio-temporal point processes and their applications*, *Statist. Sci.*, 33 (2018), pp. 299–318, <https://doi.org/10.1214/17-sts629>. (Cited on p. 333)
- [68] D. J. REZENDE AND S. MOHAMED, *Variational inference with normalizing flows*, in Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, F. R. Bach and D. M. Blei, eds., JMLR Workshop Conf. Proc. 37, 2015, pp. 1530–1538. (Cited on p. 367)
- [69] M. RIZOIU, Y. LEE, AND S. MISHRA, *Hawkes processes for events in social media*, in *Frontiers of Multimedia Research*, ACM, Morgan & Claypool, 2018, pp. 191–218. (Cited on p. 333)
- [70] M. RIZOIU, S. MISHRA, Q. KONG, M. J. CARMAN, AND L. XIE, *SIR-Hawkes: Linking epidemic models and Hawkes processes to model diffusions in finite populations*, in Proceedings of

- the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, ACM, 2018, pp. 419–428. (Cited on p. 360)
- [71] M. G. RODRIGUEZ AND I. VALERA, *Learning with Temporal Point Processes*, 2018, <http://learning.mpi-sws.org/tpp-icml18/>. (Cited on p. 333)
- [72] F. SALEHI, W. TROULEAU, M. GROSSGLAUSER, AND P. THIRAN, *Learning Hawkes processes from a handful of events*, in Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, Vancouver, Canada, 2019, pp. 12694–12704. (Cited on p. 365)
- [73] J. SHANG AND M. SUN, *Geometric Hawkes processes with graph convolutional recurrent neural networks*, in the Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, Honolulu, AAAI Press, 2019, pp. 4878–4885. (Cited on p. 358)
- [74] O. SHCHUR, M. BILOS, AND S. GÜNNEMANN, *Intensity-free learning of temporal point processes*, in 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, OpenReview.net, 2020. (Cited on p. 367)
- [75] O. SHCHUR, N. GAO, M. BILOS, AND S. GÜNNEMANN, *Fast and flexible temporal point processes with triangular maps*, in Advances in Neural Information Processing Systems (NeurIPS), 2020. (Cited on p. 367)
- [76] C. SHELTON, Z. QIN, AND C. SHETTY, *Hawkes process inference with missing data*, in Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, 2018, pp. 6425–6432. (Cited on p. 365)
- [77] I. SUTSKEVER, O. VINYALS, AND Q. V. LE, *Sequence to sequence learning with neural networks*, in Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, NeurIPS 2014, Montreal, Canada, 2014, pp. 3104–3112. (Cited on p. 357)
- [78] W. TROULEAU, J. ETESAMI, M. GROSSGLAUSER, N. KIYAVASH, AND P. THIRAN, *Learning Hawkes processes under synchronization noise*, in Proceedings of the 36th International Conference on Machine Learning, ICML 2019, Long Beach, CA, Proc. Mach. Learn. Res. 97, PMLR, 2019, pp. 6325–6334. (Cited on p. 364)
- [79] H. J. T. UNWIN, I. ROUTLEDGE, S. FLAXMAN, M.-A. RIZOU, S. LAI, J. COHEN, D. J. WEISS, S. MISHRA, AND S. BHATT, *Using Hawkes processes to model imported and local malaria cases in near-elimination settings*, PLOS Comput. Biol., 17 (2021), pp. 1–18, <https://doi.org/10.1371/journal.pcbi.1008830>. (Cited on pp. 336, 338, 366)
- [80] U. UPADHYAY, A. DE, AND M. GOMEZ-RODRIGUEZ, *Deep reinforcement learning of marked temporal point processes*, in Proceedings of the 32nd International Conference on Neural Information Processing Systems (NIPS’18), 2018, pp. 3172–3182. (Cited on pp. 333, 363)
- [81] A. VASWANI, N. SHAZEER, N. PARMAR, J. USZKOREIT, L. JONES, A. N. GOMEZ, L. KAISER, AND I. POLOSUKHIN, *Attention is all you need*, in Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, NeurIPS 2017, Long Beach, CA, 2017, pp. 5998–6008. (Cited on p. 358)
- [82] H. WANG, L. XIE, A. CUOZZO, S. MAK, AND Y. XIE, *Uncertainty quantification for inferring Hawkes networks*, in Proceedings of the 34th International Conference on Neural Information Processing Systems (NIPS’20), 2020, pp. 7125–7134. (Cited on p. 368)
- [83] Y. WANG, G. WILLIAMS, E. THEODOROU, AND L. SONG, *Variational policy for guiding point processes*, in Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, Australia, 2017, pp. 3684–3693. (Cited on p. 333)
- [84] Y. WANG, B. XIE, N. DU, AND L. SONG, *Isotonic Hawkes processes*, in Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York, 2016, pp. 2226–2234. (Cited on pp. 333, 339, 341)
- [85] S. WEI, S. ZHU, M. ZHANG, AND Y. XIE, *Goodness-of-fit test for mismatched self-exciting processes*, in Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS), Proc. Mach. Learn. Res. 130, PMLR, 2021, pp. 1243–1251. (Cited on p. 368)
- [86] S. XIAO, M. FARAJTABAR, X. YE, J. YAN, X. YANG, L. SONG, AND H. ZHA, *Wasserstein learning of deep generative point process models*, in Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, NeurIPS 2017, Long Beach, CA, 2017, pp. 3250–3259. (Cited on pp. 355, 356)
- [87] S. XIAO, H. XU, J. YAN, M. FARAJTABAR, X. YANG, L. SONG, AND H. ZHA, *Learning conditional generative models for temporal point processes*, in Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18), the 30th Innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, LA, 2018, pp. 6302–6310. (Cited on p. 357)

- [88] S. XIAO, J. YAN, M. FARAJTABAR, L. SONG, X. YANG, AND H. ZHA, *Joint Modeling of Event Sequence and Time Series with Attentional Twin Recurrent Neural Networks*, preprint, <https://arxiv.org/abs/1703.08524>, 2017. (Cited on p. 353)
- [89] S. XIAO, J. YAN, X. YANG, H. ZHA, AND S. M. CHU, *Modeling the intensity function of point process via recurrent neural networks*, in Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, San Francisco, CA, 2017, pp. 1597–1603. (Cited on pp. 354, 359)
- [90] H. XU, L. CARIN, AND H. ZHA, *Learning registered point processes from idiosyncratic observations*, in Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, 2018, pp. 5439–5448. (Cited on pp. 333, 368)
- [91] H. XU, M. FARAJTABAR, AND H. ZHA, *Learning Granger causality for Hawkes processes*, in Proceedings of the International Conference on Machine Learning, 2016, pp. 1717–1726. (Cited on pp. 333, 336, 341, 350, 365)
- [92] H. XU, D. LUO, AND H. ZHA, *Learning Hawkes processes from short doubly-censored event sequences*, in Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, Australia, 2017, pp. 3831–3840. (Cited on pp. 333, 339, 344)
- [93] J. YAN, *Recent Advance in Temporal Point Process: From Machine Learning Perspective*, 2019, [http://thinklab.sjtu.edu.cn/src/pp\\_survey.pdf](http://thinklab.sjtu.edu.cn/src/pp_survey.pdf). (Cited on p. 333)
- [94] J. YAN, X. LIU, L. SHI, C. LI, AND H. ZHA, *Improving maximum likelihood estimation of temporal point process via discriminative and adversarial learning*, in Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, Stockholm, Sweden, 2018, pp. 2948–2954. (Cited on pp. 333, 346)
- [95] G. YANG, Y. CAI, AND C. K. REDDY, *Recurrent spatio-temporal point process for check-in time prediction*, in Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM 2018, Torino, Italy, 2018, ACM, 2018, pp. 2203–2211. (Cited on pp. 334, 359)
- [96] Y. YANG, J. ETESAMI, N. HE, AND N. KIYAVASH, *Online learning for multivariate Hawkes processes*, in Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, NeurIPS 2017, Long Beach, CA, 2017, pp. 4944–4953. (Cited on pp. 333, 349, 350)
- [97] B. YUAN, H. LI, A. L. BERTOZZI, P. J. BRANTINGHAM, AND M. A. PORTER, *Multivariate spatiotemporal Hawkes processes and network reconstruction*, SIAM J. Math. Data Sci., 1 (2019), pp. 356–382, <https://doi.org/10.1137/18M1226993>. (Cited on p. 333)
- [98] A. ZAREZADE, A. DE, H. R. RABIEE, AND M. GOMEZ-RODRIGUEZ, *Cheshire: An online algorithm for activity maximization in social networks*, in the 55th Annual Allerton Conference on Communication, Control, and Computing, 2017. (Cited on pp. 333, 362)
- [99] A. ZAREZADE, U. UPADHYAY, H. R. RABIEE, AND M. GOMEZ-RODRIGUEZ, *RedQueen: An online algorithm for smart broadcasting in social networks*, in Proceedings of the Tenth ACM International Conference on Web Search and Data Mining, WSDM 2017, Cambridge, UK, 2017, pp. 51–60. (Cited on pp. 333, 361)
- [100] Q. ZHANG, A. LIPANI, Ö. KIRNAP, AND E. YILMAZ, *Self-attentive Hawkes process*, in Proceedings of the 37th International Conference on Machine Learning (ICML'20), Proc. Mach. Learn. Res. 119, PMLR, 2020, pp. 11183–11193. (Cited on pp. 358, 359)
- [101] R. ZHANG, C. J. WALDER, AND M.-A. RIZOIU, *Variational inference for sparse Gaussian process modulated Hawkes process*, in the 34th AAAI Conference on Artificial Intelligence (AAAI-20), 2020, pp. 6803–6810. (Cited on p. 359)
- [102] R. ZHANG, C. J. WALDER, M. RIZOIU, AND L. XIE, *Efficient non-parametric Bayesian Hawkes processes*, in Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, S. Kraus, ed., 2019, pp. 4299–4305. (Cited on pp. 350, 352, 359)
- [103] Q. ZHAO, M. A. ERDOGDU, H. Y. HE, A. RAJARAMAN, AND J. LESKOVEC, *Seismic: A self-exciting point process model for predicting tweet popularity*, in Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2015, pp. 1513–1522. (Cited on pp. 339, 341, 359)
- [104] F. ZHOU, Z. LI, X. FAN, Y. WANG, A. SOWMYA, AND F. CHEN, *Fast multi-resolution segmentation for nonstationary Hawkes process using cumulants*, Internat. J. Data Sci. Anal., 10 (2020), pp. 321–330, <https://doi.org/10.1007/s41060-020-00223-3>. (Cited on p. 368)
- [105] K. ZHOU, H. ZHA, AND L. SONG, *Learning social infectivity in sparse low-rank networks using multi-dimensional Hawkes processes*, in Proceedings of the Sixteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2013, Scottsdale, AZ, 2013, pp. 641–649. (Cited on p. 350)

- [106] K. ZHOU, H. ZHA, AND L. SONG, *Learning triggering kernels for multi-dimensional Hawkes processes*, in Proceedings of the International Conference on Machine Learning, 2013, pp. 1301–1309. (Cited on pp. 342, 350)
- [107] S. ZUO, H. JIANG, Z. LI, T. ZHAO, AND H. ZHA, *Transformer Hawkes process*, in Proceedings of the 37th International Conference on Machine Learning (ICML'20), Proc. Mach. Learn. Res. 119, PMLR, 2020, pp. 11692–11702. (Cited on pp. 358, 359)