

Abstract

We investigate the problem of causal discovery in temporal point processes named Hawkes processes. This is a special class of temporal point processes exhibiting a natural notion of causality, as occurrence of events in the past may increase the probability of events in the future. Our approach to Granger causal inference is based on minimum description length (MDL) principle for statistical inference. Our contribution is two fold: 1. We propose a novel general routine for estimation of MDL objective function based on Monte-Carlo simulations. 2. We employ our method to perform causal discovery in Hawkes processes. Experiments on real and synthetic datasets demonstrate significant superiority of our method in comparison with state-of-the-art algorithms.

Temporal Point Processes

A multivariate-dimensional temporal point process (TPP) is a set of random processes which is used to model occurrence of multi-type events in time. A p-variate TPP can be represented by a collection of counting processes $\{U_i\}$ where $U_i(t)$ for $t \in [0, T]$ is the number of events happened prior to time t . For each dimension i the conditional intensity function is defined as

$$\lambda_i(t) = \mathbb{E}[dU_i(t)|\mathcal{H}_t] = \lim_{\Delta t \rightarrow 0} \mathbb{E}[U_i(t + \Delta t) - U_i(t)|\mathcal{H}_t], \quad (1)$$

where \mathcal{H}_t is called the filtration, i.e the history of the process prior to time t .

Granger Causality in Temporal Point Processes

In p-dimensional time-series variable \mathbf{x}_j Granger-causes \mathbf{x}_i variable when the future of \mathbf{x}_i is better predicted when taking into account the past of variable \mathbf{x}_j .

The i-th dimension of TPP Granger-causes the j-th dimension if

$$\lambda_i(t) = \mathbb{E}[dU_i(t)|\mathcal{H}_t] \neq \mathbb{E}[dU_i(t)|\mathcal{H}_t^{-j}] \quad (2)$$

where \mathcal{H}_t^{-j} denotes the history of the process excluding events in the j-th dimension.

Hawkes 1971: A p-dimensional Hawkes process (MHP) is a p-dim TPP with conditional intensity

$$\lambda_i(t) = \mu_i + \sum_{j=1}^p \int_{-\infty}^t \phi_{ij}(t - \tau) dU_j(\tau). \quad (3)$$

μ is a vector of exogenous baseline intensities. Functions $\phi_{ij} : \mathbb{R}_0^+ \rightarrow \mathbb{R}_0^+$ are kernel functions specifying how the previous events in the j-th dim influences the intensity of the i-th dim.

Theorem 1. (Eichler et al., 2016). In a Hawkes process, the events in the j-th dimension do not Granger-cause events in the i-th dimension if and only if $\phi_{ij} \equiv 0$.

Parameter learning in a parametric class

We consider p-dim Hawkes processes with exp-decay kernels (exp-MHP) $\phi_{ij} = \alpha_{ij} \exp(-\beta_{ij}t)$ where the decay matrix β is constant and known, while the influence coefficient matrix α determines the magnitude of the effect, and it is the parameter of the model. Let $\theta_i = (\mu_i, \alpha_i^T)^T \in (\mathbb{R}_0^+)^{p+1}$ be the parameters associated to i-th dimension. We define the parameter vector of exp-MHP as $\theta = [\theta_1^T, \theta_2^T, \dots, \theta_p^T]^T \in (\mathbb{R}_0^+)^{p+p^2}$.

Facts:

1. Likelihood for exp-MHP is convex in θ and has a closed form (Ozaki 1979).
2. The objective is convex in θ (Ogata 1981).
3. Solving MLE for exp-MHP can be done efficiently (Bacry et al. 2015).

Normalized Maximum Likelihood

Denote a parametric statistical model by a family of probability distributions $M = \{p(\cdot|\theta) : \theta \in \Theta\}$ where Θ is the parameter space. We partition Θ into disjoint subsets $\{\Theta_\gamma : \gamma \in \Gamma\}$ to define family of models $\{M_\gamma : \gamma \in \Gamma\}$. For data $\mathbf{x} \in \mathcal{X}$, the **normalized maximum likelihood** (NML) distribution (**Shtarkov 1978**) is a general technique applying the MDL principle to statistical model selection. It depends on the luckiness function $v : \Theta \rightarrow \mathbb{R}_0^+$ and for each submodel M_γ is

$$p_{v|\gamma}^{NML}(\mathbf{x}) = \frac{\max_{\theta \in \Theta_\gamma} p(\mathbf{x}|\theta)v(\theta)}{\int_{\mathcal{X}} \max_{\theta \in \Theta_\gamma} p(\mathbf{s}|\theta)v(\theta) ds}. \quad (4)$$

The logarithm of the integral is called the **model complexity** $COMP(M_\gamma; v)$.

MDL-Based Statistical Inference

Denote $R_v(\theta, \mathbf{x}) = p(\mathbf{x}|\theta)v(\theta)$ and let $r_v(\theta; \mathbf{x}) := \log R_v(\theta, \mathbf{x}) = \log p(\mathbf{x}|\theta) + \log v(\theta)$ for each $\mathbf{x} \in \mathcal{X}$ and $\theta \in \Theta$. We define the MDL estimator based on v for a specific model $M_\gamma \subset M$ as $\hat{\theta}_{v|\gamma}(\mathbf{x}) := \arg \min_{\theta \in \Theta_\gamma} -r_v(\theta; \mathbf{x})$. The model selection on NML picks over Γ the model minimizing

$$L_v(\gamma; \mathbf{x}) = -\log \pi(\gamma) - \log p_{v|\gamma}^{NML}(\mathbf{x}) = -\log \pi(\gamma) - r_v(\hat{\theta}_{v|\gamma}(\mathbf{x}); \mathbf{x}) + COMP(M_\gamma; v) \quad (5)$$

where $\pi(\cdot)$ is a given distribution on Γ .

The MDL function incorporates the trade-off between goodness-of-fit measured by $r_v(\hat{\theta}_{v|\gamma}(\mathbf{x}); \mathbf{x})$ and the model complexity $COMP(M_\gamma; v)$. The model selection based on MDL selects

$$\hat{\gamma}^{MDL} = \arg \min_{\gamma \in \Gamma} L_v(\gamma; \mathbf{x}). \quad (6)$$

Our Method for Estimation of MDL Objective Function

Computation of the goodness-of-fit

To compute the goodness-of-fit term $r_v(\hat{\theta}_{v|\gamma}(\mathbf{x}); \mathbf{x})$ we need to efficiently compute the MDL estimator $\hat{\theta}_{v|\gamma}$ which is the result of a minimization problem over the restricted parameter space Θ_γ . If the minimization problem is convex, the solution would be unique and efficiently computable. Therefore we require two (mild) conditions:

1. Θ_γ is a convex set, e.g., a Euclidean space.
2. $-r_v(\theta; \mathbf{x}) = -\log(p(\mathbf{x}|\theta)) - \log(v(\theta))$ is a convex function.

Estimation of model complexity

Denote $COMP(M_\gamma; v) := \log \int_{\mathcal{X}} R_v(\hat{\theta}_{v|\gamma}(\mathbf{s}); \mathbf{s}) ds$. We define

$$Q_{v|\gamma}(\mathbf{s}, \mathbf{z}) = \frac{R_v(\hat{\theta}_{v|\gamma}(\mathbf{s}); \mathbf{s})}{p(\mathbf{s}|\mathbf{z})} \quad (7)$$

and rewrite model complexity as $COMP(M_\gamma; v) = \log \mathbb{E}_{\mathbf{X} \sim p}[Q_{v|\gamma}(\mathbf{X}, \theta)|\theta = \mathbf{z}]$.

We randomize \mathbf{z} to any full support distribution, we can rewrite model complexity

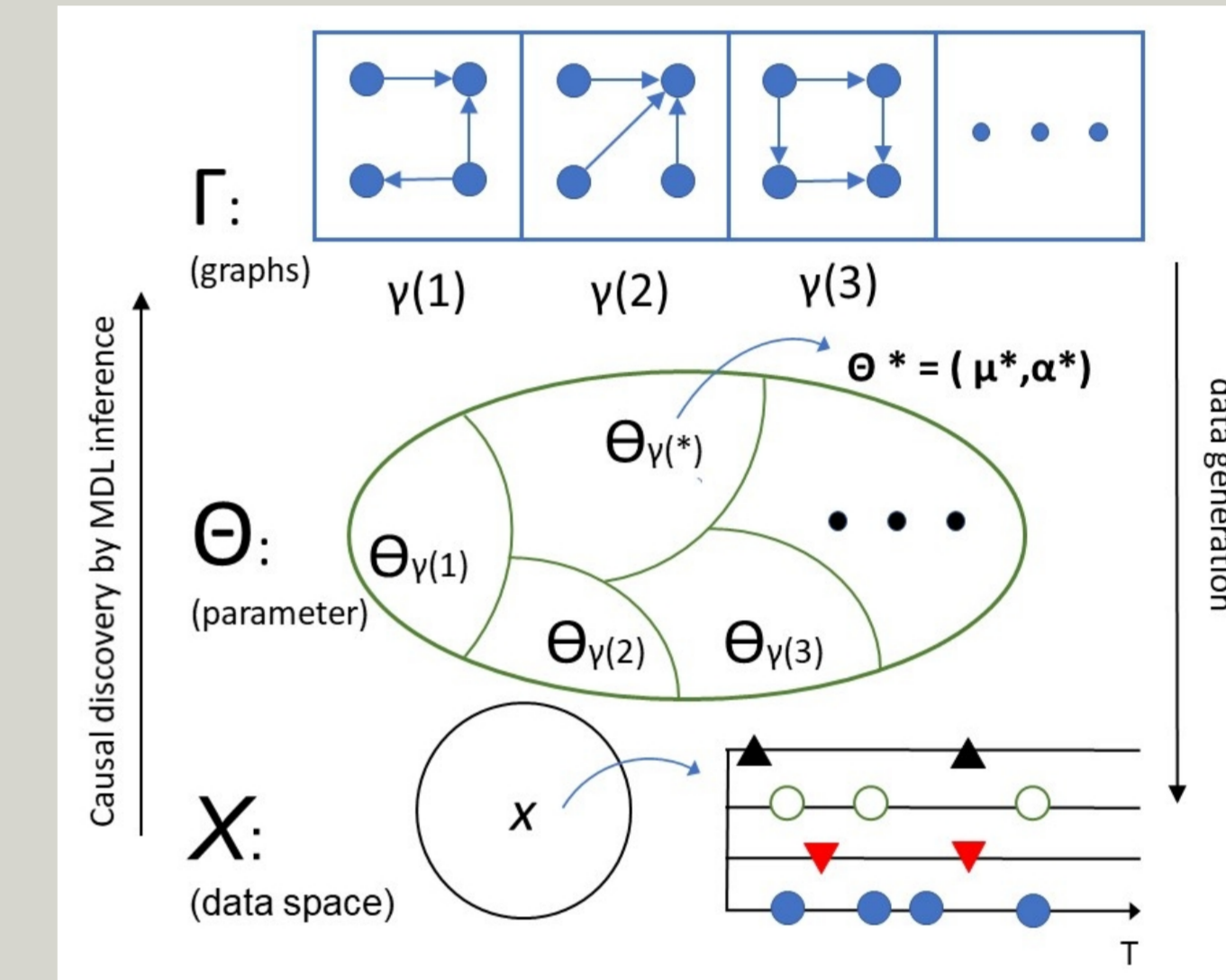
$$COMP(M_\gamma; v) = \log \mathbb{E}_{\mathbf{X}, \theta \sim p}[Q_{v|\gamma}(\mathbf{X}, \theta)]. \quad (8)$$

$Q_{v|\gamma}$ is efficiently computable at all points. By taking multiple joint samples of \mathbf{X} and using them in $Q_{v|\gamma}$ we may get random draws of $Q_{v|\gamma}(\mathbf{X}, \theta)$. Taking the average of these values gives us an unbiased estimator for the above expectation.

Causal Discovery in Hawkes Processes

Causal discovery as a model selection problem

Let Γ be the set of all $p \times p$ binary matrices and let M be a set of all p-dim exp-MHP models. For each $\gamma \in \Gamma$, denote $M_\gamma \subset M$ a set of all p-dim exp-MHP models with **Granger causal graph** γ (adjacency matrix). According to Theorem 1, $\alpha_{ij}^* = 0 \iff \gamma_{ij}^* = 0$. Hence, the true parameter θ^* lies in exactly one sub-model: M_{γ^*} . Therefore, **finding the true causal graph γ^* is equivalent to finding the true sub-bmodel M_{γ^*} which contains the true parameter θ^* .**



Algorithm MDLH:

Solving the independent optimization problems $\hat{\gamma}_i^{MDL} = \arg \min_{\gamma_i \in \{0,1\}^p} L_v^i(\gamma_i; \mathbf{x})$, we get the MDL estimator - causal graph $\hat{\gamma}^{MDL} = [\hat{\gamma}_1^{MDL} | \hat{\gamma}_2^{MDL} | \dots | \hat{\gamma}_p^{MDL}]^T$ with $\mathcal{O}(p2^p)$ complexity.

For sparse graphs with max degree $\delta \ll p$ we get polynomial complexity $\mathcal{O}(Np^\delta)$ for N Monte Carlo simulations. The algorithm can be also amortized and parallelized.

Synthetic and real experiments

p	7			20		
	200	400	700	500	1300	2000
MDLH	77.4	84.7	89.3	79.4	82.8	84.4
ADM4	68.4	72.6	78.5	26.8	29.9	31.5
NPHC	49.3	58.8	61.3	27.3	34.5	40.0
ML	68.7	74.6	80.4	25.8	28.2	29.4
LS	68.3	74.4	76.9	26.4	29.8	31.3
IC	NA	NA	NA	NA	NA	NA
Random	30.0	30.0	30.0	7.5	7.5	7.5

MDLH and state-of-art methods in F1 measure: ADM4 from Zhou et al. 2013, NPHC from Achab et al. 2017; ML is max. likelihood, LS least squares, IC the best F1 of AIC, BIC and HQ.

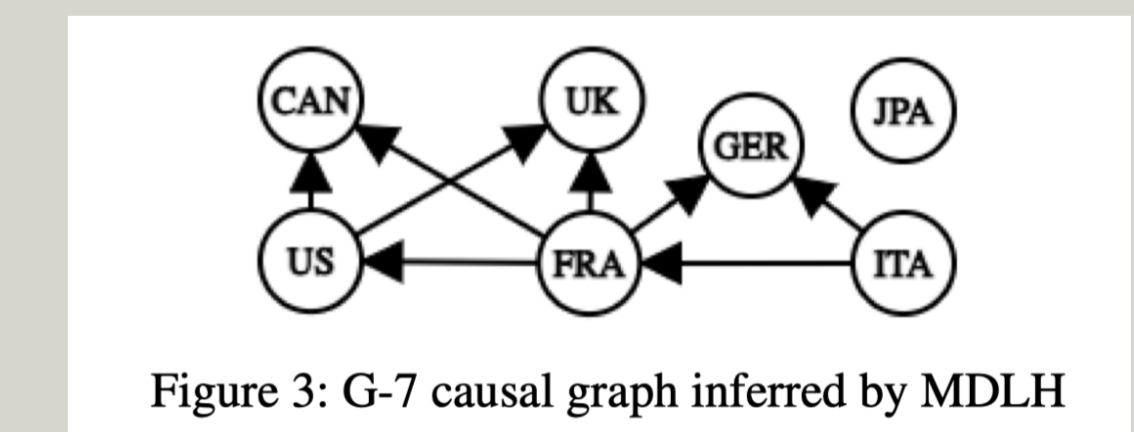


Figure 3: G-7 causal graph inferred by MDLH

Influence network among sovereign bonds of 7 large and developed economies called G-7 from daily return volatility 2003-2014 (data from Demirev et al. 2018). More details in our paper.