Statistical Learning with Hawkes Processes

5th NUS-USPC Workshop

Stéphane Gaïffas



USPC Université Sorbonne Paris Cité

・ロト ・日・ ・ヨ・ ・ヨ・ クへぐ

- 2 Hawkes processes
- 3 Dimension reduction for MHP
- 4 Random matrix theory
- **5** Causality maps
- 6 Accelerating training time
- 7 Software

Example 1. Social networks. Understand who is influencing **twitter**: based on the timestamps patterns of messages, **web-data**: publication activity of websites/blogs

Example 2. High Frequency Finance.

From zoomed financial signals ($\Delta t \approx 1$ ms, upward / downward price proves and other order book features), build a "causality map"

Example 3. Health-care. Impact of some health events to other health events (all being timestamped, longitudinal data)

▲□▶ ▲□▶ ▲ ≧▶ ▲ ≧▶ ≧ 釣�?

2 Hawkes processes

- 3 Dimension reduction for MHP
- 4 Random matrix theory
- **5** Causality maps
- 6 Accelerating training time
- 7 Software

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ● ● ●

Introduction

From:



Build:







◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶

Introduction

Setting

- For each node $i \in I = \{1, \ldots, d\}$ we have a set Z^i of *events*
- Any $\tau \in Z^i$ is the occurrence time of an event related to i

Counting process

- Put $N_t = [N_t^1 \cdots N_t^d]^\top$
- $N_t^i = \sum_{\tau \in Z^i} \mathbf{1}_{\tau \leq t}$

Intensity

• Stochastic intensities $\lambda_t = [\lambda_t^1 \cdots \lambda_t^d]^\top$, $\lambda_t^i =$ intensity of N_t^i

$$\lambda_t^i = \lim_{dt \to 0} \frac{\mathbb{P}(N_{t+dt}^i - N_t^i = 1 | \mathcal{F}_t)}{dt}$$

- λ_t^i = instantaneous rate of event occurence at time t for node i
- λ_t characterizes the distribution of N_t [Daley et al. 2007]
- Patterns can be captured by *putting structure* on λ_t

▲□▶ ▲□▶ ▲ ≧▶ ▲ ≧▶ ● ④ � �

The Multivariate Hawkes Process (MHP)

Scaling

• We observe N_t on [0, T]. "Asymptotics" in $T \to +\infty$. d is "large"

The Hawkes process

- A particular structure for λ_t : auto-regression
- N_t is called a *Hawkes* process [Hawkes 1971] if

$$\lambda_{t}^{i} = \mu_{i} + \sum_{j=1}^{d} \int_{0}^{t} \varphi^{ij}(t-t') dN_{t'}^{j} = \mu_{i} + \sum_{j=1}^{d} \sum_{t' \in Z_{j}: t' < t} \varphi^{ij}(t-t')$$

- $\mu_i \in \mathbb{R}^+$ exogenous intensity
- φ^{ij} non-negative integrable and causal (support \mathbb{R}_+) functions
- φ^{ij} are called *kernels*. Encodes the impact of an action by node j on the activity of node i
- Captures *auto-excitation* and *cross-excitation* across nodes, a phenomenon observed in social networks [Crane et al. 2008]

◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□▶ ◆□

A simple parametrization of the MHP

For d = 1, K = 1 and $\varphi^{11}(t) = e^{-1}$, intensity $\lambda_{\theta,t}$ looks like:



▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

Stability condition of the MHP

Stability condition

Introduce

$$G^{ij} = \int_0^{+\infty} \varphi^{ij}(t) dt$$

• Spectral norm must satisfy $\| {m G} \| < 1$ to ensure stability and stationarity of the process

Sum of exponentials parametric model:

$$\lambda_{\theta,t}^{i} = \mu_{i} + \int_{(0,t)} \sum_{j=1}^{d} \sum_{k=1}^{K} a_{ij}^{k} \times \alpha_{k} e^{-\alpha_{k}(t-s)} dN_{s}^{j}$$

for $i \in \{1, ..., d\}$ with $\alpha_1, ..., \alpha_K > 0$ given and parameters to infer are $\theta = [\mu, \mathbf{A}]$ with

- baselines $\mu = [\mu_1 \cdots \mu_d]^\top \in \mathbb{R}^d_+$
- interactions $\mathbf{A} = [a_{ij}]_{1 \leq i,j \leq d} \in \mathbb{R}^{d \times d}_+ =$ "adjacency matrix"

<ロト < 母 ト < 臣 ト < 臣 ト 三 の < で</p>

A brief history of MHP

Brief history

- Introduced in Hawkes 1971
- Earthquakes and geophysics [Kagan and Knopoff 1981], [Zhuang et al. 2012]
- Genomics [Reynaud-Bouret and Schbath 2010]
- High-frequency Finance [Bacry et al. 2013]
- Terrorist activity [Mohler et al. 2011, Porter and White 2012]
- Neurobiology [Hansen et al. 2012]
- Social networks [Carne and Sornette 2008], [Zhou et al.2013]
- And even FPGA-based implementation [Guo and Luk 2013]

・・

A brief history of MHP



Home / Bitcoin 201 / Analyzing Trade Clustering To Predict Price Movement In Bitcoin Trading



Analyzing Trade Clustering To Predict Price Movement In Bitcoin Trading

Sep 19, 2013 Posted By Jonathan Heusser In Bitcoin 201, Economics, Featured, News, Trading Tagged Analysis, Bitcoin Trading, Hawkes Process, Jonathan Heusser, London, Price, Trading

SQ (?

MHP in large dimension

What do we want to do?

- Deal with large number of events and large dimension *d* (number of nodes)
- End up with a *tractable* and *scalable* optimization problem

Goodness-of-fit functionals. Two choices: minus log-likelihood

$$-\ell_{T}(\theta) = \frac{1}{T} \sum_{i=1}^{d} \left\{ \int_{0}^{T} \lambda_{\theta,t}^{i} dt - \int_{0}^{T} \log \lambda_{\theta,t}^{i} dN_{t}^{i} \right\}$$

or least-squares

$$R_T(\theta) = \frac{1}{T} \sum_{i=1}^d \left\{ \int_0^T (\lambda_{\theta,t}^i)^2 dt - 2 \int_0^T \lambda_{\theta,t}^i dN_t^i \right\}$$

- 2 Hawkes processes
- Oimension reduction for MHP
- 4 Random matrix theory
- **5** Causality maps
- 6 Accelerating training time
- 7 Software

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ● ● ●

Paper. E. Bacry, S. G., J.-F. Muzy, *A generalization error bound for sparse and low-rank multivariate Hawkes processes*, in revision in JMLR

- Parametric setting $\varphi^{ij}(t) = (\mathbf{A})_{ij} \times h(t)$
- Low-rank and sparsity inducing penalization on A
- Introduces a sharp tuning of the penalizations using data-driven weights
- Leads to optimal error bounds for penalized least-squares (sharp sparse oracle inequality)

Prior assumptions

• Users are basically inactive and react mostly if stimulated:

 μ is sparse

• Everybody does not interact with everybody:

A is sparse

• Interactions have community structure, possibly overlapping, a small number of factors explain interactions:



 \boldsymbol{A} is low-rank

・ロト・日・・日・・日・ うへぐ

Standard convex relaxations

(Tibshirani (01), Srebro et al. (05), Bach (08), Candès & Tao (09), etc.)

• Convex relaxation of $\|\boldsymbol{A}\|_0 = \sum_{ij} \mathbf{1}_{\boldsymbol{A}_{ij} > 0}$ is ℓ_1 -norm:

$$\|oldsymbol{\mathcal{A}}\|_1 = \sum_{ij} |oldsymbol{\mathcal{A}}_{ij}|$$

• Convex relaxation of rank is trace-norm:

$$\|\boldsymbol{A}\|_{*} = \sum_{j} \sigma_{j}(\boldsymbol{A}) = \|\sigma(\boldsymbol{A})\|_{1}$$

where $\sigma_1(\mathbf{A}) \geq \cdots \geq \sigma_d(\mathbf{A})$ singular values of \mathbf{A}

We use the following penalizations

- Use ℓ_1 penalization on μ
- Use ℓ_1 penalization on \boldsymbol{A}
- Use trace-norm penalization on A



Balls are on the set of 2×2 symmetric matrices identified with \mathbb{R}^3 .

▲□▶ ▲□▶ ▲ ≧▶ ▲ ≧▶ ≧ 釣�?

Leads to

$$\hat{\theta} = (\hat{\mu}, \hat{\mathbf{A}}) \in \operatorname*{argmin}_{\theta = (\mu, \mathbf{A}) \in \mathbb{R}^{d}_{+} \times \mathbb{R}^{d \times d}_{+}} \{ R_{T}(\theta) + \operatorname{pen}(\theta) \},\$$

with penalization

$$\mathsf{pen}(\theta) = \tau_1 \|\mu\|_1 + \gamma_1 \|\boldsymbol{A}\|_1 + \gamma_* \|\boldsymbol{A}\|_*$$

The features scaling problem

- Features scaling is necessary for "linear approaches" in supervised learning
- No features and labels here!
- We solve this by sharp data-driven tuning of the penalization terms
- Required a new theory for random matrices with entries that are continuous-time martingales

・ロ・・中・・川・・・日・ ・日・ うくぐ



Left: AUC; Middle: Estimation error; Right: Kandall rank correlation

・ロト・日本・日本・日本・日本・日本

A strong theoretical guarantee

- Recall $\langle \lambda_1, \lambda_2 \rangle_T = \frac{1}{T} \sum_{i=1}^d \int_0^T \lambda_{1,t}^i \lambda_{2,t}^i dt$ and $\|\lambda\|_T^2 = \langle \lambda, \lambda \rangle_T$
- Assume RE in our setting (Restricted Eigenvalues, Compressed Sensing literature)

Theorem. We have

$$\begin{aligned} \|\lambda_{\hat{\theta}} - \lambda^*\|_{\mathcal{T}}^2 &\leq \inf_{\theta} \left\{ \|\lambda_{\theta} - \lambda^*\|_{\mathcal{T}}^2 + c_1 \kappa(\theta)^2 \Big(\|(\hat{w})_{\mathsf{supp}(\mu)}\|_2^2 \\ &+ \|(\hat{W})_{\mathsf{supp}(\mathcal{A})}\|_F^2 + \hat{w}_*^2 \operatorname{rank}(\mathcal{A}) \Big) \right\} \end{aligned}$$

with a probability larger than $1 - c_2 e^{-x}$.

Roughly, $\hat{\theta}$ achieves an optimal tradeoff between approximation and complexity given by

$$\frac{\|\mu\|_0 \log d}{T} \max_i \bar{N}^i([0, T]) + \frac{\|\boldsymbol{A}\|_0 \log d}{T} \max_{ij} \hat{v}_T^{ij} \\ + \frac{\operatorname{rank}(A) \log d}{T} \lambda_{\max}(\hat{\boldsymbol{V}}_T)$$

- Complexity measured both by sparsity and rank
- Convergence has shape $(\log d)/T$, where T =length of the observation interval
- Terms are balanced by "empirical variance" terms

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶

- 2 Hawkes processes
- 3 Dimension reduction for MHP
- 4 Random matrix theory
- **5** Causality maps
- 6 Accelerating training time
- 7 Software

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ▶ ▲□ ▶ ▲□

Contribution 2. Random matrix theory

Paper. E. Bacry, S. G. and J-F Muzy, *Concentration inequalities for matrix martingales in continuous time*, PTRF (2017)

Consider a $m \times n$ matrix-martingale given by

$$\boldsymbol{Z}_t = \int_0^t \mathbb{T}_s \circ d\boldsymbol{M}_s,$$

with $(\mathbf{M}_t)_{t\geq 0}$ "white" Brownian random matrix and $(\mathbb{T}_t)_{t\geq 0}$ rank-4 preditable tensor. Then

$$\mathbb{P}\bigg[\|\boldsymbol{Z}_t\|_{\mathrm{op}} \geq \sqrt{2\boldsymbol{v}(\boldsymbol{x} + \log(\boldsymbol{m} + \boldsymbol{n}))} \ , \ \sigma^2(\boldsymbol{Z}_t) \leq \boldsymbol{v}\bigg] \leq e^{-\boldsymbol{x}}$$

for any v, x > 0 with

$$\sigma^{2}(\boldsymbol{Z}_{t}) = \max\left(\left\|\sum_{j=1}^{n} \langle \boldsymbol{Z}_{\bullet,j} \rangle_{t}\right\|_{\mathrm{op}}, \left\|\sum_{j=1}^{m} \langle \boldsymbol{Z}_{j,\bullet} \rangle_{t}\right\|_{\mathrm{op}}\right).$$

Contribution 2. Random matrix theory

- Strong generalization of previously known inequalities to continuous time (Tropp 2011)
- Very different approach (random matrix tools + stochastic calculus)
- Also the "Poissonian" case: martingale with sub-exponential jumps (counting process, Hawkes processes)

Interesting particular case (previously unknown!). Consider $P = [P_{ij}]$ a $n \times m$ random matrix where P_{ij} is Poisson (λ_{ij}) and put $\lambda = [\lambda_{ij}]$. Then

$$\mathbb{P}igg(\|oldsymbol{N}-oldsymbol{\lambda}\|_{\mathrm{op}}\geq \sqrt{2(\|oldsymbol{\lambda}\|_{1,\infty}ee\|oldsymbol{\lambda}\|_{\infty,1})x}+rac{x}{3}igg)\leq (n+m)e^{-x}$$

for any x > 0, where $\|\lambda\|_{1,\infty}$ (resp. $\|\lambda\|_{\infty,1}$) stands for the maximum ℓ_1 -norm of rows (resp. columns)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ♪ ◇◇◇

- 2 Hawkes processes
- 3 Dimension reduction for MHP
- 4 Random matrix theory
- 5 Causality maps
- 6 Accelerating training time
- 7 Software

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ● ● ●

Paper. Achab et al. Uncovering Causality from Multivariate Hawkes Integrated Cumulants, ICML (2017) and JMLR (2017)

Reminder.

$$\lambda_t^i = \mu_i + \sum_{j=1}^d \int_0^t \varphi^{ij}(t-t') dN_{t'}^j,$$

Idea

• Don't estimate φ^{ij} (parametric, non-parametric), but only

$$g^{ij}=(oldsymbol{G})_{ij}=\int_{0}^{+\infty}arphi^{ij}$$

 Introducing the (unobserved) counting process N^{i←j} = of events from *i* with direct ancestor from *j*, we have

$$\mathbb{E}[dN_t^{i\leftarrow j}] = g^{ij}\mathbb{E}[dN_t^j]$$

<□▶ < @ ▶ < E ▶ < E ▶ 0 < 0</p>

• g^{ij} = average number of events from *i* triggered by one event from *j*

- Actually, if $\varphi^{ij} \ge 0$ then $g^{ij} = 0$ iff N^i does not Granger-cause N^j
- The matrix **G** encodes **causality**

How to estimate **G** directly?

• We know (Jovanovic 2014) how to relate integrated cumulants of N_t to $\boldsymbol{R} = (\boldsymbol{I} - \boldsymbol{G})^{-1}$

$$\Lambda^{i} = \sum_{m=1}^{d} R^{im} \mu^{m}$$

$$C^{ij} = \sum_{m=1}^{d} \Lambda^{m} R^{im} R^{jm}$$

$$K^{ijk} = \sum_{m=1}^{d} (R^{im} R^{jm} C^{km} + R^{im} C^{jm} R^{km} + C^{im} R^{jm} R^{km} - 2\Lambda^{m} R^{im} R^{jm} R^{km}).$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ● ● ●

NPHC. Cumulant matching method for estimation of G

- Compute estimates \widehat{C} and $\widehat{K^c}$ of the third order cumulants of the process
- Find \widehat{R} that matches these empirical cumulants

$$\mathcal{L}(\boldsymbol{R}) = (1-\kappa) \| \boldsymbol{K}^{\boldsymbol{c}}(\boldsymbol{R}) - \widehat{\boldsymbol{K}^{\boldsymbol{c}}} \|_{2}^{2} + \kappa \| \boldsymbol{C}(\boldsymbol{R}) - \widehat{\boldsymbol{C}} \|_{2}^{2},$$

• Put

$$\widehat{oldsymbol{R}} = oldsymbol{I} - \Big(rg\min_{oldsymbol{R}\in\Theta}\mathcal{L}_{\mathcal{T}}(oldsymbol{R})\Big)^{-1}$$

Theorem.

• It is consistent! (under some assumptions... quite technical)

Remarks.

- Highly non-convex problem: polynomial or order 10 with respect to the entries of *R*
- Not so hard, local minima turns out to be good (deep learning literature), we simply use AdaGrad
- Using order **three** is important (two is not enough): integrated covariance contains only symmetric information: unable to provide causal information
- NPHC scales better than state-of-the-art methods, is robust towards the kernel shape and directly outputs the kernel integral
- Simple tensorflow code

<□> <□> <□> <=> <=> <=> <=> <=> <</p>

Experiment with MemeTracker dataset

- keep the 200 most active sites
- contains publication times of articles in many websites/blogs, with hyperlinks
- ullet pprox 8 millions events
- Use hyperlinks to establish an estimated ground truth for the matrix G

Method	ODE	GC	ADM4	NPHC
RelErr	0.162	0.19	0.092	0.071
MRankCorr	0.07	0.053	0.081	0.095
Time (s)	2944	2780	2217	38



Experiment with MemeTracker dataset

◆□▶ ◆□▶ ◆ ≧▶ ◆ ≧ ◆ つへぐ

Order book dynamics.

- Order book: a list of buy and sell orders for a specific financial instrument, the list being updated in real-time throughout the day
- Understand the self and cross-influencing dynamics of all event types in an order book

Introduce

$$N_t = (P_t^{(a)}, P_t^{(b)}, T_t^{(a)}, T_t^{(b)}, L_t^{(a)}, L_t^{(b)}, C_t^{(a)}, C_t^{(b)})$$

where

- $P^{(a)}$ (resp. $P^{(b)}$): upward (resp. downward) price moves;
- T^(a) (resp. T^(b)): market orders at the ask (resp. at the bid) that do not move the price;
- L^(a) (resp. L^(b)): limit orders at the ask (resp. at the bid) that do not move the price;
- C^(a) (resp. C^(b)): cancel orders at the ask (resp. at the bid), that do not move the price.

Data: DAX future contracts between 01/01/2014 and 03/01/2014.

SQA



・ロト・西ト・山下・山下・山下・山下・山下・山下

Interpretable results

- Any 2 × 2 sub-matrix with same kind of inputs (i.e. Prices changes, Trades, Limits or Cancels) is symmetric: ask and bid have symmetric roles;
- Prices are mostly cross-excited: price increase is most likely followed by a price decrease, and conversely;
- Market, limit and cancel orders are strongly self-excited: persistence of order flows, and splitting of meta-orders into sequences of smaller orders.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ♪ ◇◇◇

- 2 Hawkes processes
- 3 Dimension reduction for MHP
- 4 Random matrix theory
- **5** Causality maps



6 Accelerating training time



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ▶ ▲□ ▶ ▲□

Paper. E. Bacry, S. G., J.-F. Muzy, I. Mastromatteo, *Mean-field inference of Hawkes point processes*, Journal of Physics A, 2016

- Dedicated optimization algorithm for the Hawkes MLE with large number of nodes
- Based on a mean-field approximation
- Partially understood (proof on toy cases)
- Improves state-of-the-art by orders of magnitude

Mean-Field approximation (large number of nodes *d* helps!)



< □

590

≣►

4

< E

< 47 ▶

E



Fluctuations $\mathbb{E}^{1/2}[(\lambda_t^1/\Lambda^1 - 1)^2]$



◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□ ◆

No clean proof yet (only on toy example) but works very well empirically



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ▶ ▲□ ▶ ▲□

Faster by several order of magnitude than state-of-the-art solvers



・ロト・日本 エリ・ エリ・ トーロ・

Conclusion

Take-home message

- Hawkes Process for "time-oriented" machine learning
- Surprisingly relevant to fit real-word phenomena (auto-excitation, user influence)
- Very flexible: intensity can depend on features, other processes, etc.

Main contributions

- Sharp theoretical guarantees for low-rank inducing penalization for Hawkes models
- New results about **concentration of matrix-martingales** in continuous time
- Go beyond the parametric approach: **unveil causality using integrated cumulants** matching
- Improved training time of the Hawkes model using a "mean-field" approximation

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□

- 2 Hawkes processes
- 3 Dimension reduction for MHP
- 4 Random matrix theory
- **5** Causality maps
- 6 Accelerating training time
- Software

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ● ●

Software: tick library

- Python 3 et C++11
- Open-source (BSD-3 License)
- pip install tick (on MacOS and Linux...)
- https://x-datainitiative.github.io/tick
- Statistical learning for time-dependent models
- Point processes (Poisson, Hawkes), Survival analysis, GLMs (parallelized, sparse, etc.)
- A strong simulation and optimization toolbox
- Partnership with Intel (use-case for new processors with 256 cores)
- Contributors welcome!

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへで

Software: tick library

tick 0.3 Home Examples API Browse -

Search

Fox ne on Cithus

tick

tick a machine learning library for Python 3. The focus is on statistical learning for time dependent systems, such as point processes. Tick features also tools for generalized linear models, and a generic optimization toolbox.

The core of the library is an optimization module providing model computational classes, solvers and proximal operators for regularization. It comes also with inference and simulation tools intended for end-users.

Show me »

Examples

Examples of how to simulate models, use the optimization toolbox, or use user-friendly inference tools.

Simulation

User-friendly classes for simulation of data

Inference

User-friendly classes for inference of models

Optimization

The core module of the library: an optimization toolbox consisting of models, solvers and prox (penalization) classes. Almost all of them can be combined together.

▲□▶▲□▶▲≡▶▲≡▶ ≡ ∽٩.0

Thank you!

▲□▶▲母▶▲≧▶▲≧▶ ≧ りへぐ