Statistical learning with Hawkes processes and new matrix concentration inequalities

Emmanuel Bacry¹, Stéphane Gaïffas¹, Jean-Francois Muzy^{1,2}



Winter 2016

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

¹École Polytechnique and CNRS ²CNRS, Université de Corse

- You have users of a system (a social network, an e-commerce platform, etc.)
- You want to quantify the level of interaction between users
- You don't want to use only *declared* interactions, such as "friendship" or "likes". This information is often deprecated, and not really related to the activity of users
- You want levels of interaction driven by user's actions, using the timestamps' patterns of actions

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Introduction

From:



We want to quantify interactions between users:



▲□▶ ▲□▶ ▲臣▶ ▲臣▶ = 臣 = のへで

- A *d*-dimensional counting process $N = [N_1, \ldots, N_d]^{\top}$
- d is "large"
- Observed on [0, T]. "Asymptotics" in $T \to +\infty$
- N_j has intensity λ_j , namely

 $\mathbb{P}(N_j \text{ has a jump in } [t, t + dt] \mid \mathcal{F}_t) = \lambda_j(t)dt$

for $j = 1, \ldots, d$ where \mathcal{F}_t some filtration

• MHP assumes the following autoregressive structure:

$$\lambda_j(t) = \mu_j(t) + \int_{(0,t)} \sum_{k=1}^d \varphi_{j,k}(t-s) dN_k(s),$$

- $\mu_j(t) \ge 0$ baseline intensity of the *j*-th coordinate
- $\varphi_j : \mathbb{R}^+ \to \mathbb{R}^+$ self-exciting component
- Write this in matrix form

$$\lambda(t) = \mu + \int_{(0,t)} \varphi(t-s) dN(s),$$

with $\boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^\top$ and $\boldsymbol{\varphi}(t) = [\varphi_{j,k}(t)]_{1 \leq j,k \leq d}$.

• Notation:

$$\int_{(0,t)} \varphi(t-s) dN_k(s) = \sum_{i: 0 < T_{i,k} < t} \varphi(t-T_{i,k})$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Introduced by Hawkes in 1971

- Earthquakes and geophysics : Kagan and Knopoff (1981), Zhuang, Harte, Werner, Hainzl and Zhou (2012)
- Genomics : Reynaud-Bouret and Schbath (2010)
- High-frequency Finance : Bacry Delattre Hoffmann and Muzy (2013)
- Terrorist activity : Porter and White (2012)
- Neurobiology : Hansen, Reynaud-Bouret and Rivoirard (2012)
- **Social networks** : Carne and Sornette (2008), Simma and Jordan (2010), Zhou Song and Zha (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,

naa

Hawkes Process, Jonathan Heusser, London, Price, Trading

Parametric estimation (Maximum likelihood)

- First work : Ogata 78
- Simma and Jordan (2010), Zhou Song and Zha (2013)
 - \rightarrow Expected Maximization (EM) algorithms, with priors

Non parametric estimation

Marsan Lengliné (2008), generalized by Lewis, Mohler (2010)
 → EM for penalized likelihood function
 → Monovariate Hawkes processes, Small amount of data, No

- theoretical results
- Reynaud-Bouret and Schbath (2010)
 → Developed for small amount of data (Sparse penalization)
- Bacry and Muzy (2014)
 - \rightarrow Larger amount of data

- Do inference directly from actions of users
- Understand the community structure of users underlying the actions
- Exploit the hidden lower-dimensional structure of the network for inference/prediction

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Dimension *d* is large:

- ullet Need a simple parametric model on μ and arphi
- For inference: we want a **tractable** and **scalable** optimization problem

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

• We want to encode some prior assumptions using (convex) penalization

A simple parametrization of the MHP

Simple parametrization:

• Constant baselines $\mu_j(\cdot) \equiv \mu_j$

Take

$$\varphi_{j,k}(t) = a_{j,k} e^{-\alpha_{j,k}t}$$

• $a_{j,k} =$ level of interaction between nodes j and k

• $\alpha_{j,k} =$ lifetime of instantaneous excitation of node j by node k

The matrix

$$\boldsymbol{A} = [a_{j,k}]_{1 \le j,k \le d}$$

is understood has a weighted adjacency matrix of mutual excitement of nodes $\{1, \ldots, d\}$

• A is non-symmetric: "oriented graph"

We end up with intensities

$$\lambda_{j,\theta}(t) = \mu_j + \int_{(0,t)} \sum_{k=1}^d a_{j,k} e^{-\alpha_{j,k}(t-s)} dN_k(s)$$

for $j \in \{1, \ldots, d\}$ where

$$\theta = [\mu, \mathbf{A}, \boldsymbol{\alpha}]$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

with

- baselines $\mu = [\mu_1, \dots, \mu_d]^\top \in \mathbb{R}^d_+$
- interactions $\boldsymbol{A} = [a_{j,k}]_{1 \leq j,k \leq d} \in \mathbb{R}^{d \times d}_+$
- decays $\boldsymbol{\alpha} = [\alpha_{j,k}]_{1 \leq j,k \leq d} \in \mathbb{R}^{d \times d}_+$

For d = 1, intensity λ_{θ} looks like this:



Minus log-likelihood

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

Least-squares

$$R_{T}(\theta) = \sum_{j=1}^{d} \left\{ \int_{0}^{T} \lambda_{j,\theta}(t)^{2} dt - 2 \int_{0}^{T} \lambda_{j,\theta}(t) dN_{j}(t) \right\}$$

with

$$\lambda_{j,\theta}(t) = \mu_j + \sum_{k=1}^d a_{j,k} \int_{(0,t)} \exp(-\alpha_{j,k}(t-s)) dN_k(s)$$

where $\theta = [\mu, \mathbf{A}, \alpha]$ with $\mu = [\mu_j]$, $\mathbf{A} = [a_{j,k}]$, $\alpha = [\alpha_{j,k}]$

(ロ)、(型)、(E)、(E)、 E) の(の)

A simple framework

Put $\|\lambda_{\theta}\|_{T}^{2} = \langle \lambda_{\theta}, \lambda_{\theta} \rangle_{T}$ with

$$\langle \lambda_ heta, \lambda_{ heta'}
angle_{\mathcal{T}} = rac{1}{\mathcal{T}} \sum_{j=1}^d \int_{[0,\mathcal{T}]} \lambda_{j, heta}(t) \lambda_{j, heta'}(t) dt.$$

so that least-squares writes

$$R_{T}(\theta) = \|\lambda_{\theta}\|_{T}^{2} - \frac{2}{T} \sum_{j=1}^{d} \int_{[0,T]} \lambda_{j,\theta}(t) dN_{j}(t)$$

It is natural: if N has ground truth intensity λ^* then

$$\mathbb{E}[R_{\mathcal{T}}(\theta)] = \mathbb{E}\|\lambda_{\theta}\|_{\mathcal{T}}^2 - 2\mathbb{E}\langle\lambda_{\theta},\lambda^*\rangle_{\mathcal{T}} = \mathbb{E}\|\lambda_{\theta}-\lambda^*\|_{\mathcal{T}}^2 - \|\lambda^*\|_{\mathcal{T}},$$

where we used "signal + noise" decomposition (Doob-Meyer):

$$dN_j(t) = \lambda^*(t)dt + dM_j(t)$$

where M_j martingale

A strong assumption: assume that

$$\varphi_{j,k}(t) = a_{j,k}h_{j,k}(t)$$

for **known** $h_{j,k}$ meaning that

$$\lambda_{j,\theta}(t) = \mu_j + \int_{(0,t)} \sum_{k=1}^d a_{j,k} h_{j,k}(t-s) dN_k(s),$$



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

where
$$\theta = [\mu, \mathbf{A}]$$
 with $\mu = [\mu_1, \dots, \mu_d]^\top$ and $\mathbf{A} = [a_{j,k}]_{1 \le j,k \le d}$

Ways of cheating:

• approximate $h_{j,k}$ using a known dictionary $\{h_1, \ldots, h_L\}$ and learn the coefficients

$$\lambda_{j,\theta}(t) = \mu_j + \int_{(0,t)} \sum_{k=1}^d \sum_{l=1}^L a_{j,k,l} h_l(t-s) dN_k(s),$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

(but then use group Lasso $\sum_{1 \leq j,k \leq d} \|h_{j,k,\bullet}\|_2$)

• Don't estimate the kernels... (ongoing work)

But let's stick with the simple framework

Prior assumptions

• Some users are basically inactive and react only 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

Standard convex relaxations [Tibshirani (01), ..., Srebro et al. (05), Bach (08), Candès & Recht (08), ...]

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

$$\|\boldsymbol{A}\|_1 = \sum_{j,k} |\boldsymbol{A}_{j,k}|$$

• Convex relaxation of rank is trace-norm:

$$\|A\|_* = \sum_j \sigma_j(A) = \|\sigma(A)\|_1$$

where $\sigma_1(A) \geq \cdots \geq \sigma_d(A)$ singular values of **A**

So, we use the following penalizations

- Use ℓ_1 penalization on $oldsymbol{\mu}$
- Use ℓ_1 penalization on **A**
- Use trace-norm penalization on A

[but other choices might be interesting...]

NB1: to induce **sparsity AND low-rank** on **A**, we use the mixed penalization

$$oldsymbol{A}\mapsto w_*\|oldsymbol{A}\|_*+w_1\|oldsymbol{A}\|_1$$

NB2: recent work by Richard et al (2013): much better way to induce sparsity and low-rank than this (but no theory)



 $\{ \pmb{\mathsf{A}} : \| \pmb{\mathsf{A}} \|_* \leq 1 \} \qquad \qquad \{ \pmb{\mathsf{A}} : \| \pmb{\mathsf{A}} \|_1 \leq 1 \} \qquad \{ \pmb{\mathsf{A}} : \| \pmb{\mathsf{A}} \|_1 + \| \pmb{\mathsf{A}} \|_* \leq 1 \}$

The balls are computed on the set of 2×2 symmetric matrices, which is identified with $\mathbb{R}^3.$

▲ロト ▲冊ト ▲ヨト ▲ヨト ヨー の々ぐ

We end up with the problem

$$\hat{\theta} \in \operatorname*{argmin}_{\theta \in \mathbb{R}^{d}_{+} \times \mathbb{R}^{d \times d}_{+}} \big\{ R_{\mathcal{T}}(\theta) + \operatorname{pen}(\theta) \big\},$$

with mixed penalizations

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

But there is the "features scaling" problem

- Features scaling is necessary for "linear approaches" in supervised learning
- No features and labels here!
- \Rightarrow Can be solved here by fine tuning of the penalization terms

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Algorithm

Consider instead

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^{d}_{+} \times \mathbb{R}^{d \times d}_{+}}{\operatorname{argmin}} \{ R_{T}(\theta) + \operatorname{pen}(\theta) \},\$$

where this time

$$\mathsf{pen}(\theta) = \|\mu\|_{1,\hat{w}} + \|\boldsymbol{A}\|_{1,\hat{\boldsymbol{W}}} + \hat{w}_*\|\boldsymbol{A}\|_*$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Penalization tuned by data-driven weights \hat{w} , \hat{W} and \hat{w}_* to solve the "scaling" problem
- Comes from sharp controls of the noise terms, using new probabilistic tools

- Can be solved using first-order routines
- Gradient of $R_T(\theta)$ using a recursion formula [Ogata (1988)]

 \rightarrow When carefully done complexity of one gradient is O(nd) (instead of $O(n^2d)$ for the naive approach), where n = number of events (very large)

 \rightarrow The gradient on each node $j \in \{1, \ldots, d\}$ can be computed in $\ensuremath{\mathsf{parallel}}$

• Computation bootleneck is the heavy use of exp and log [accelerated using some ugly hacking]

• Proximal of trace norm requires many truncated SVD: we use the default's Lanczos's implementation of Python (fast enough)

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ● ● ●

Machine Learning (not only) with Point Processes

NLPP is yet another machine learning library for Python, with a particular emphasis on point processes (Poisson processes, Hawkes processes, Cox regression), but includes also generalized linear models (GLM). It comes with optimization algorithms for inference and provides tools for simulation of datasets. A particular focus is on optimization: an extensive optimization toolbox is proposed, with recent state-of-the-art stochastic solvers.

Main highlights

- · Python 3 only!
- · Fast, most computation are done in C++11 (including multi-threading)
- Contains an optimization toolbox: solvers (optimization algorithms), models (for computing gradient among other things) and prox (penalization) classes can be combined interactively in <u>iPython</u> for instance
- · Features state-of-the-art stochastic optimization algorithms, with parallel implementations
- · Support sparse datasets (sparse features matrices)
- · Interplay between Python and C++ is done using swig

Optimization toolbox

- Models [mlpp.optim.model]
 - Introduction
 - · Contents
 - · Generalized linear models
 - · Hawkes model
 - · What's under the carpet?
- Proximal operators [mlpp.optim.prox]
 - Introduction
 - Available operators
 - Example
- Solvers [mlpp.optim.solver]
 - Introduction
 - · Available solvers

Toy example: take matrix \boldsymbol{A} as



Numerical experiment: dimension 10, 210 parameters



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Numerical experiment: dimension 100, 20100 parameters



◆□▶ ◆□▶ ◆豆▶ ◆豆▶ ̄豆 _ のへで

Numerical experiment: dimension 100, 20100 parameters



ロト 4 聞 ト 4 臣 ト 4 臣 ト 三 三 のへで

Estimation errors of **A** (measured by ℓ_2 norm)



AUC for support selection \boldsymbol{A}





< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Towards a statistical guarantee: first order condition can be written as: for any θ

$$\begin{split} \|\lambda_{\hat{\theta}} - \lambda^*\|_{\mathcal{T}}^2 + \|\lambda_{\hat{\theta}} - \lambda_{\theta}\|_{\mathcal{T}}^2 - \|\lambda_{\theta} - \lambda^*\|_{\mathcal{T}}^2 \\ &\leq -\langle\theta_{\partial}, \hat{\theta} - \theta\rangle + \frac{2}{\mathcal{T}}\langle\hat{\mu} - \mu, \bar{M}_{\mathcal{T}}\rangle + \frac{2}{\mathcal{T}}\langle\hat{\boldsymbol{A}} - \boldsymbol{A}, \boldsymbol{Z}_{\mathcal{T}}\rangle, \end{split}$$

for $\theta_{\partial} \in \partial \operatorname{pen}(\theta)$ and we use $\frac{2}{T} \langle \hat{\boldsymbol{A}} - \boldsymbol{A}, \boldsymbol{Z}_{T} \rangle \leq \frac{2}{T} \| \hat{\boldsymbol{A}} - \boldsymbol{A} \|_{*} \| \boldsymbol{Z}_{T} \|_{\operatorname{op}}$

 $\bar{M}_T = [\int_0^T dM_1(t) \cdots \int_0^T dM_d(t)]^\top$ and \boldsymbol{Z}_t matrix martingale with entries

$$(\boldsymbol{Z}_t)_{j,k} = \int_0^t \int_{(0,s)} h_{j,k}(s-u) dN_k(u) dM_j(s), \qquad (1)$$

or

$$\boldsymbol{Z}_t = \int_0^t \operatorname{diag}[dM_s] \boldsymbol{H}_s,$$

with \boldsymbol{H}_t predictable process with entries

$$(\boldsymbol{H}_t)_{j,j'} = \int_{(0,t)} h_{j,j'}(t-s) dN_{j'}(s)$$

Noise term is a matrix-martingale in continuous time:

$$\frac{1}{T}Z_{7}$$

・ロト ・ 理 ・ ・ ヨ ・ ・ ヨ ・ うへつ

wee need to control $\frac{1}{T} \| \boldsymbol{Z}_T \|_{\text{op}}$

A consequence of our new concentration inequalities (more after):

$$\mathbb{P}\left[\frac{\|\boldsymbol{Z}_t\|_{\text{op}}}{t} \ge \sqrt{\frac{2v(x+\log(2d))}{t}} + \frac{b(x+\log(2d))}{3t}, \\ b_t \le b, \quad \lambda_{\max}(\boldsymbol{V}_t) \le v\right] \le e^{-x},$$

for any v, x, b > 0, where

$$\mathbf{V}_{t} = \frac{1}{t} \int_{0}^{t} \|\mathbf{H}_{s}\|_{2,\infty}^{2} \begin{bmatrix} \operatorname{diag}[\lambda_{s}^{*}] & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_{s}^{\top} \operatorname{diag}[\mathbf{H}_{s}\mathbf{H}_{s}^{\top}]^{-1} \operatorname{diag}[\lambda_{s}^{*}]\mathbf{H}_{s} \end{bmatrix} ds$$

and $b_{t} = \sup_{s \in [0,t]} \|\mathbf{H}_{s}\|_{2,\infty} (\|\cdot\|_{2,\infty} = \operatorname{maximum} \ell_{2} \text{ row norm})$

Useless for statistical learning! Event $\lambda_{\max}(V_t) \le v$ is annoying and V_t is **not observable** (depends on λ^*)!

Theorem [Something better]. For any x > 0, we have

$$\frac{\|\boldsymbol{Z}_t\|_{\mathrm{op}}}{t} \leq 8\sqrt{\frac{(x+\log d+\hat{\ell}_{x,t})\lambda_{\max}(\hat{\boldsymbol{V}}_t)}{t}} + \frac{(x+\log d+\hat{\ell}_{x,t})(10.34+2.65b_t)}{t}$$

with a probability larger than $1 - 84.9e^{-x}$, where

$$\hat{\boldsymbol{V}}_{t} = \frac{1}{t} \int_{0}^{t} \|\boldsymbol{H}_{s}\|_{2,\infty}^{2} \begin{bmatrix} \operatorname{diag}[dN_{s}] & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{H}_{s}^{\top} \operatorname{diag}[\boldsymbol{H}_{s}\boldsymbol{H}_{s}^{\top}]^{-1} \operatorname{diag}[dN_{s}]\boldsymbol{H}_{s} \end{bmatrix} ds$$

and small ugly term:

$$\hat{\ell}_{x,t} = 4\log\log\Big(\frac{2\lambda_{\max}(\hat{\mathbf{V}}_t) + 2(4+b_t^2/3)x}{x} \vee e\Big) + 2\log\log\Big(b_t^2 \vee e\Big).$$

This is a non-commutative deviation inequality with **observable** variance

These concentration inequalities leads to a data-driven tuning of penalization

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

• Solves the "scaling" problem in this context \approx features scaling in supervised learning

Controls on $\|\boldsymbol{Z}_{\mathcal{T}}\|_{\infty} = \max_{j,k} |A_{j,k}|$ and $\|\boldsymbol{Z}_{\mathcal{T}}\|_{\text{op}}$ leads to the following tuning of the penalizations

For ℓ_1 penalization of μ : $\|\mu\|_{1,\hat{w}} = \sum_{j=1}^d \hat{w}_j |\mu_j|$ with

$$\hat{w}_{j} = 6\sqrt{2}\sqrt{\frac{(x + \log d + \hat{\ell}_{x,j,T})N_{j}([0, T])/T}{T}} + 27.93\frac{x + \log d + \hat{\ell}_{x,j,T}}{T}$$

where $N_j([0, T]) = \int_0^T dN_j(t)$, namely

$$\hat{w}_j \approx c \sqrt{\frac{N_j([0,T])/T}{T}}$$

 Each coordinate j of μ is penalized (roughly) by N_j([0, T)]/T: estimated average intensity of events of node j For ℓ_1 penalization of \boldsymbol{A} : $\|\boldsymbol{A}\|_{1,\hat{\boldsymbol{W}}} = \sum_{1 \leq j,k \leq d} \hat{\boldsymbol{W}}_{j,k} |\boldsymbol{A}_{j,k}|$ with

$$\hat{W}_{j,k} = 4\sqrt{2}\sqrt{\frac{(x+2\log d + \hat{\ell}_{x,j,k,T})\hat{V}_{j,k}(T)}{T}} + 18.62\frac{(x+2\log d + \hat{\ell}_{x,j,k,T})B_{j,k}(T)}{T}$$

where

$$B_{j,k}(t) = \sup_{s \in [0,t]} \int_{(0,t)} h_{j,k}(t-s) dN_k(s)$$

$$V_{j,k}(t) = \frac{1}{t} \int_0^t \left(\int_{(0,s)} h_{j,k}(s-u) dN_k(u) \right)^2 dN_j(s)$$

namely

$$\hat{\boldsymbol{W}}_{j,k} pprox c \sqrt{rac{\hat{\boldsymbol{V}}_{j,k}(T)}{T}}$$

 $\hat{V}_{j,k}(t)$ estimates the "variance" of self-excitements between nodes j and k

For trace-norm penalization of \mathbf{A} : $\hat{w}_* \|\mathbf{A}\|_*$ with

$$\hat{w}_{*} = 8\sqrt{\frac{(x + \log d + \hat{\ell}_{x,T})\lambda_{\max}(\hat{V}_{T})}{T}} + \frac{2(x + \log d + \hat{\ell}_{x,T})(10.34 + 2.65b_{t})}{T}}$$

namely

$$\hat{w}_{*} pprox \sqrt{rac{\lambda_{\mathsf{max}}(\hat{oldsymbol{\mathcal{V}}}_{\mathcal{T}})}{\mathcal{T}}}$$

◆□ ▶ < 圖 ▶ < 圖 ▶ < 圖 ▶ < 圖 • 의 Q @</p>

- Data-driven weights that comes from "empirical" Bernstein's inequalities, entrywise and for operator norm of Z_T
- $\hat{\mathbf{V}}_{j,k}(t)$ and $\lambda_{\max}(\hat{\mathbf{V}}_t)$ are estimations (based on optional variation) of the variance terms from Bernstein's inequality
- $B_{j,k}(t)$ and b_t are L^{∞} terms (sub-exponential actually) from these Bernstein's inequalities
- Leads to a data-driven scaling of penalization: deals correctly with the inhomogeneity of information over nodes

A sharp oracle inequality

- Recall $\langle \lambda_1, \lambda_2 \rangle_T = \frac{1}{T} \sum_{j=1}^d \int_0^T \lambda_{1,j}(t) \lambda_{2,j}(t) dt$ and $\|\lambda\|_T^2 = \langle \lambda, \lambda \rangle_T$
- Assume RE in our setting (Restricted Eigenvalues), which is a standard assumption to obtain fast rates for the Lasso (and other convex-relaxation based procedures)

Theorem. We have

$$\begin{aligned} \|\lambda_{\hat{\theta}} - \lambda^*\|_{\mathcal{T}}^2 &\leq \inf_{\theta} \left\{ \|\lambda_{\theta} - \lambda^*\|_{\mathcal{T}}^2 + \kappa(\theta)^2 \Big(\frac{5}{4} \|(\hat{w})_{\mathsf{supp}(\mu)}\|_2^2 \\ &+ \frac{9}{8} \|(\hat{\boldsymbol{W}})_{\mathsf{supp}(\boldsymbol{A})}\|_{F}^2 + \frac{9}{8} \hat{w}_*^2 \operatorname{rank}(\boldsymbol{A}) \Big) \right\} \end{aligned}$$

with a probability larger than $1 - 146e^{-x}$.

• Leading constant 1

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

$$\frac{\|\mu\|_0(x+\log d)}{T} \max_j N_j([0,T])/T \\ + \frac{\|\boldsymbol{A}\|_0(x+2\log d)}{T} \max_{j,k} \hat{\boldsymbol{V}}_{j,k}(T) \\ + \frac{\operatorname{rank}(A)(x+\log d)}{T} \lambda_{\max}(\hat{\boldsymbol{V}}_T)$$

- Complexity measured both by sparsity and rank
- Convergence has shape $(\log d)/T$, where T =length of the observation interval

• These terms are balanced by "empirical variance" terms

Main tool: new concentration inequalities for matrix martingales in continuous time

Introduce

$$\boldsymbol{Z}_t = \int_0^t \boldsymbol{A}_s (\boldsymbol{C}_s \odot d\boldsymbol{M}_s) \boldsymbol{B}_s,$$

where $\{A_t\}$, $\{C_t\}$ and $\{B_t\}$ predictable and where $\{M_t\}_{t\geq 0}$ is a "white" matrix martingale, in the sense that $[\text{vec}M]_t$ is diagonal

NB: entries of Z_t are given by

$$(\boldsymbol{Z}_t)_{i,j} = \sum_{k=1}^p \sum_{l=1}^q \int_0^t (\boldsymbol{A}_s)_{i,k} (\boldsymbol{C}_s)_{k,l} (\boldsymbol{B}_s)_{l,j} (d\boldsymbol{M}_s)_{k,l}.$$

• $\langle \boldsymbol{M} \rangle_t$ = entrywise predictable quadratic variation, so that

 $oldsymbol{M}_t^{\odot 2} - \langle oldsymbol{M}
angle_t$

martingale

- vectorization operator $\mathrm{vec}:\mathbb{R}^{p imes q}\to\mathbb{R}^{pq}$ stacks vertically the columns of \pmb{X}
- $\langle vec \mathbf{M} \rangle_t$ is the $pq \times pq$ matrix with entries that are all pairwise quadratic covariations, so that

$$\operatorname{vec}(\boldsymbol{M}_t)\operatorname{vec}(\boldsymbol{M}_t)^{\top} - \langle \operatorname{vec} \boldsymbol{M} \rangle_t$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

is a martingale.

M_t = *M^c_t* + *M^d_t*, where *M^c_t* is a continuous martingale and *M^d_t* is a purely discountinuous martingale. Its (entrywise) quadratic variation is defined as

$$[\boldsymbol{M}]_t = \langle \boldsymbol{M}^c \rangle_t + \sum_{0 \le s \le t} (\Delta \boldsymbol{M}_t)^2, \qquad (2)$$

and its quadratic covariation by

$$[\operatorname{vec} \boldsymbol{M}]_t = \langle \operatorname{vec} \boldsymbol{M}^c \rangle_t + \sum_{0 \le s \le t} \operatorname{vec}(\Delta \boldsymbol{M}_s) \operatorname{vec}(\Delta \boldsymbol{M}_s)^\top.$$

We say that \boldsymbol{M} is *purely discontinuous* if the process $\langle \operatorname{vec} \boldsymbol{M}^c \rangle_t$ is identically the zero matrix.

Concentration for purely discountinuous matrix martingale:

• **M**_t is purely discountinuous and we have

$$\langle \pmb{M}
angle_t = \int_0^t \pmb{\lambda}_s ds$$

for a non-negative and predictable intensity process $\{\lambda_t\}_{t\geq 0}$. • Standard moment assumptions (subexponential tails) Introduce

$$\boldsymbol{V}_t = \int_0^t \|\boldsymbol{A}_s\|_{\infty,2}^2 \|\boldsymbol{B}_s\|_{2,\infty}^2 \boldsymbol{W}_s ds$$

where

$$\boldsymbol{W}_{t} = \begin{bmatrix} \boldsymbol{W}_{t}^{1} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{W}_{t}^{2} \end{bmatrix}, \qquad (3)$$

$$\begin{split} \boldsymbol{W}_t^1 &= \boldsymbol{A}_t \operatorname{diag}[\boldsymbol{A}_t^{\top} \boldsymbol{A}_t]^{-1} \operatorname{diag}\left[(\boldsymbol{C}_t^{\odot 2} \odot \boldsymbol{\lambda}_t) \mathbb{1}\right] \boldsymbol{A}_t^{\top} \\ \boldsymbol{W}_t^2 &= \boldsymbol{B}_t^{\top} \operatorname{diag}[\boldsymbol{B}_t \boldsymbol{B}_t^{\top}]^{-1} \operatorname{diag}\left[(\boldsymbol{C}_t^{\odot 2} \odot \boldsymbol{\lambda}_t)^{\top} \mathbb{1}\right] \boldsymbol{B}_t \end{split}$$

Introduce also

$$b_t = \sup_{s \in [0,t]} \|\boldsymbol{A}_s\|_{\infty,2} \|\boldsymbol{B}_s\|_{2,\infty} \|\boldsymbol{C}_s\|_{\infty}.$$

Theorem.

$$\mathbb{P}\Big[\|m{Z}_t\|_{\mathrm{op}} \ge \sqrt{2v(x+\log(m+n))} + rac{b(x+\log(m+n))}{3}, \ b_t \le b, \quad \lambda_{\max}(m{V}_t) \le v\Big] \le e^{-x},$$

• First result of this type for matrix-martingale in continuous time

◆□ ▶ < 圖 ▶ < 圖 ▶ < 圖 ▶ < 圖 • 의 Q @</p>

Corollary. { N_t } a $p \times q$ matrix, each (N_t)_{*i*,*j*} is an independent inhomogeneous Poisson processes with intensity $(\lambda_t)_{i,j}$. Consider the martingale $M_t = N_t - \Lambda_t$, where $\Lambda_t = \int_0^t \lambda_s ds$ and let { C_t } be deterministic and bounded. We have

$$\begin{split} \left\| \int_{0}^{t} \boldsymbol{C}_{s} \odot d(\boldsymbol{N}_{t} - \boldsymbol{\Lambda}_{t}) \right\|_{\text{op}} \\ & \leq \sqrt{2 \Big(\left\| \int_{0}^{t} \boldsymbol{C}_{s}^{\odot 2} \odot \boldsymbol{\lambda}_{s} ds \right\|_{1,\infty} \vee \left\| \int_{0}^{t} \boldsymbol{C}_{s}^{\odot 2} \odot \boldsymbol{\lambda}_{s} ds \right\|_{\infty,1} \Big) (x + \log(p + q))} \\ & + \frac{\sup_{s \in [0,t]} \| \boldsymbol{C}_{s} \|_{\infty} (x + \log(p + q))}{3} \end{split}$$

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

Corollary. Even more particular: **N** random matrix where $N_{i,j}$ are independent Poisson variables with intensity $\lambda_{i,j}$. We have

$$egin{aligned} \|oldsymbol{N}-oldsymbol{\lambda}\|_{\mathrm{op}} &\leq \sqrt{2(\|oldsymbol{\lambda}\|_{1,\infty} ee \|oldsymbol{\lambda}\|_{\infty,1})(x+\log(p+q))} \ &+ rac{x+\log(p+q)}{3}. \end{aligned}$$

- Up to our knowledge, not previously stated in literature
- NB: In the Gaussian case: variance depends on maximum l₂ norm of rows and columns (cf. Tropp (2011))

- We have as well a non-commutative Hoeffding's inequality when *M_t* has continuous paths (allowing Itô's formula...), with a similar variance term
- Tools from stochastic calculus, use of the dilation operator and some classical matrix inequalities about the trace exponential and the SDP order.
- A difficult proposition: a control of the quadratic variation of the pure jump process

$$\boldsymbol{U}_t^u = \sum_{0 \le s \le t} \left(e^{u \Delta \mathscr{S}(\boldsymbol{Z}_s)} - u \Delta \mathscr{S}(\boldsymbol{Z}_s) - \boldsymbol{I} \right)$$

given by

$$\langle \boldsymbol{U}^{\boldsymbol{\xi}} \rangle_t \preceq \int_0^t \frac{\varphi\left(\boldsymbol{\xi} \| \boldsymbol{A}_s \|_{\infty,2} \| \boldsymbol{B}_s \|_{2,\infty} \| \boldsymbol{C}_s \|_{\infty}\right)}{\| \boldsymbol{C}_s \|_{\infty}^2} \boldsymbol{W}_s ds,$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

where $\varphi(x) = e^x - x - 1$.

- Theoretical study of learning algorithms for "time-oriented" models needs new probabilistic results
- In our case new concentration results for matrix martingales in continuous time

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

• Leads to a better scaling of penalizations

Going back to maximum-likelihood estimation, with d very large

- For inference, exploit the fact that d is large
- \Rightarrow use a Mean-Field approximation! (from Delattre et al. 2015)



э

Sac

Mean-field inference for Hawkes

We don't understand perfectly why this works yet (proof on a toy example)

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

100
 ≈ 10
1
0.1
0.01
0.001
0.001
0.001
0.001

But it does very well empirically



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

But it does very well empirically



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

And it is faster by several order of magnitude than state-of-the-art solvers



・ロト ・聞ト ・ヨト ・ヨト

э

- Better understanding of the Mean-Field based inference
- Quantifying influence without kernels estimation
- Experiments: ongoing project with Préfecture de l'Oise (car theft), High-frequency Finance
- Non-constant baseline (block-stationarity): block constant with total-variation penalization

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- NMF based learning of the self-excitement matrix
- SBM prior for community detection results
- Even better optimization using recent stochastic gradient algorithms (but needs some tweaking...)