

Machine learning using Hawkes processes and concentration for matrix martingales

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



Winter 2016

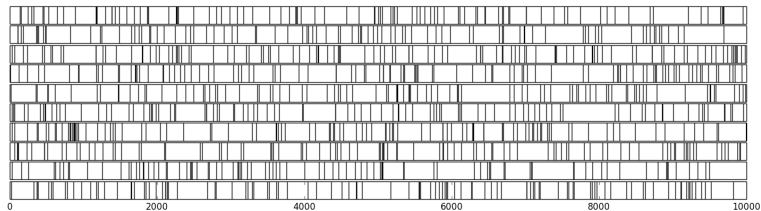
¹École Polytechnique and CNRS

²CNRS, Université de Corse

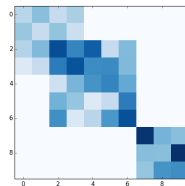
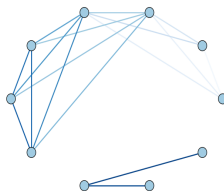
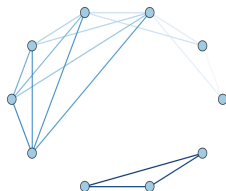
- You have users of a system (a social network, an e-commerce platform, etc.)
- You want to quantify the level of interaction between users
- You don't want to use only *declared* interactions, such as “friendship” or “likes”. This information is often deprecated, and not really related to the activity 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:



Model: Multivariate Hawkes Process (MHP)

- A d -dimensional counting process $N = [N_1, \dots, N_d]^\top$
- d is “large”
- Observed on $[0, T]$. “Asymptotics” in $T \rightarrow +\infty$
- 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 $j = 1, \dots, d$ where \mathcal{F}_t some filtration

Model: Multivariate Hawkes Process (MHP)

- MHP assumes the following autoregressive structure:

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

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

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

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

THE GENESIS BLOCK



Digital currency research and data

[HOME](#)[NEWS](#)[MINING](#)[TRADING](#)[ECONOMICS](#)[REGULATION](#)[BUSINESSES](#)[BITCOIN](#)

[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](#)



Parametric estimation (Maximum likelihood)

- First work : Ogata 78
- Simma and Jordan (2010), Zhou Song and Zha (2013)
 - 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)
 - Larger amount of data

What do we want to do with this?

- 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:

- Need a simple parametric model on μ and φ
- For inference: we want a **tractable** and **scalable** optimization problem
- We want to encode some prior assumptions by penalizing the likelihood

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

$$\mathbf{A} = [a_{j,k}]_{1 \leq j, k \leq d}$$

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

- \mathbf{A} is non-symmetric: oriented graph

A simple parametrization of the MHP

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 $j \in \{1, \dots, d\}$ where

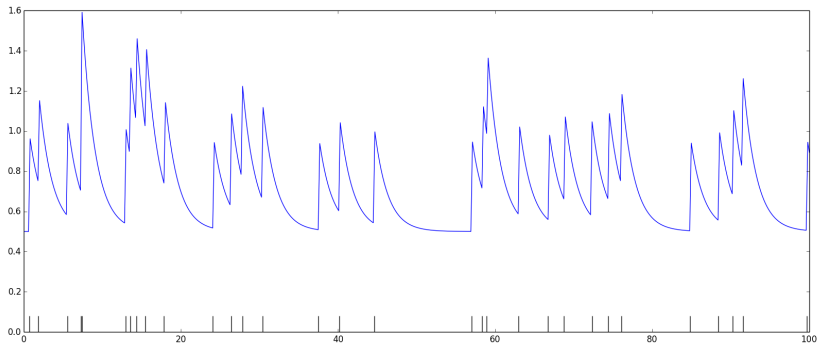
$$\theta = [\mu, \mathbf{A}, \alpha]$$

with

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

A simple parametrization of the MHP

For $d = 1$, intensity λ_θ looks like this:



Goodness-of-fit = $-\log$ -likelihood is given by:

$$-\ell_T(\theta) = \sum_{j=1}^d \left\{ \int_0^T (\lambda_{j,\theta}(t) - 1) dt - \int_0^T \log \lambda_{j,\theta}(t) dN_j(t) \right\}$$

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

Prior assumptions

- Some users are basically inactive and react only if stimulated:

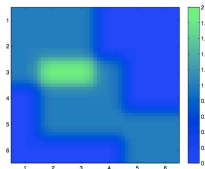
μ is sparse

- Everybody does not interact with everybody:

\mathbf{A} is sparse

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

\mathbf{A} is low-rank



- Decays α are not sparse, but $\alpha_{j,k}$ should be regularized proportionally to $a_{j,k}$

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

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

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

- Convex relaxation of rank is trace-norm:

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

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

So, we use the following penalizations

- Use ℓ_1 penalization on μ
- Use ℓ_1 penalization on \mathbf{A}
- Use trace-norm penalization on \mathbf{A}
- Use ℓ_2^2 penalization on α , weighted by \mathbf{A}

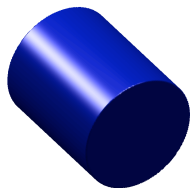
[but other choices might be interesting...]

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

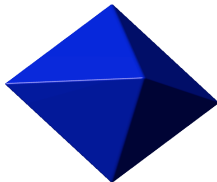
$$\mathbf{A} \mapsto w_* \|\mathbf{A}\|_* + w_1 \|\mathbf{A}\|_1$$

NB2: recent work by Richard et al (2013): better way to induce sparsity and low-rank than the sum

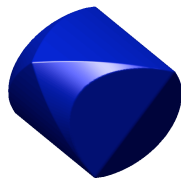
Sparse and low-rank matrices



$$\{\mathbf{A} : \|\mathbf{A}\|_* \leq 1\}$$



$$\{\mathbf{A} : \|\mathbf{A}\|_1 \leq 1\}$$



$$\{\mathbf{A} : \|\mathbf{A}\|_1 + \|\mathbf{A}\|_* \leq 1\}$$

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

Finally, consider

$$\hat{\theta} \in \underset{\theta=(\boldsymbol{\mu}, \mathbf{A}, \boldsymbol{\alpha})}{\operatorname{argmin}} \left\{ -\frac{1}{T} \ell_T(\theta) + \tau \|\boldsymbol{\mu}\|_1 + \gamma_1 \|\mathbf{A}\|_1 \right. \\ \left. + \gamma_* \|\mathbf{A}\|_* + \frac{\kappa}{2} \|\mathbf{A} \odot \boldsymbol{\alpha}\|_F^2 \right\}$$

where we recall

$$-\frac{1}{T} \ell_T(\theta) = \frac{1}{T} \sum_{j=1}^d \left\{ \int_0^T \lambda_{j,\theta}(t) dt - \int_0^T \log \lambda_{j,\theta}(t) dN_j(t) \right\}$$

with

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

Penalized maximum likelihood: a problem

Problem: $\theta \mapsto -\ell_T(\theta)$ not convex! Indeed

$$(a, \alpha) \mapsto ah_\alpha(t)$$

never convex when $\alpha \mapsto h_\alpha(t)$ is convex



We **want** convexity for:

- Convergence to a global optimum
- Plethora of optimization algorithm
- If smooth (Lipschitz gradient): optimal first-order techniques
[first order=mandatory for large scale problems]

Generic in the chosen penalization [if proximal operator easy to compute]

A solution: the **perspective function** trick:

- If $\alpha \mapsto h_\alpha(t)$ is convex, then

$$(a, \alpha) \mapsto ah_{\alpha/a}(t)$$

is **convex**!

- Reparametrization $\beta_{j,k} = a_{j,k}\alpha_{j,k}$, leading to

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

with $\theta = [\mu, \mathbf{A}, \beta]$ for $\beta = [\beta_{j,k}]_{1 \leq j, k \leq d}$

- With this reparametrization

$$\theta \mapsto \lambda_{j,\theta}(t)$$

is **convex**!

The reparametrization $\beta = \mathbf{A} \odot \alpha$ leads to

$$\hat{\theta} \in \underset{\theta=(\mu, \mathbf{A}, \beta)}{\operatorname{argmin}} \left\{ -\frac{1}{T} \ell_T(\theta) + \tau \|\mu\|_1 + \gamma_1 \|\mathbf{A}\|_1 \right. \\ \left. + \gamma_* \|\mathbf{A}\|_* + \frac{\kappa}{2} \|\beta\|_F^2 \right\} \quad (1)$$

where

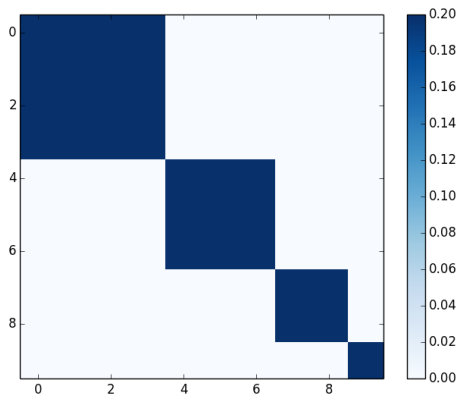
$$-\frac{1}{T} \ell_T(\theta) = \frac{1}{T} \sum_{j=1}^d \left\{ \int_0^T \lambda_{j,\theta}(t) dt - \int_0^T \log \lambda_{j,\theta}(t) dN_j(t) \right\}$$

with

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

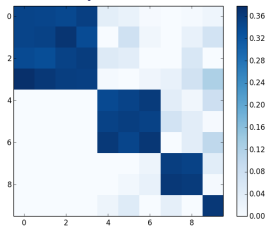
- Can be solved using optimal first-order routines
- Gradient of $-\ell_T(\theta)$ using a recursion formula [Ogata (1988)]
 - When carefully done complexity of one gradient is $O(nd)$ (instead of $O(n^2d)$ for the naive approach), where n = number of events (very large)
 - We have scalable / parallelized code for this: the gradient on each node $j \in \{1, \dots, d\}$ can be computed in a **distributed** fashion
- Computation bottleneck 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 when using incremental truncation

Toy example: take matrix \mathbf{A} as

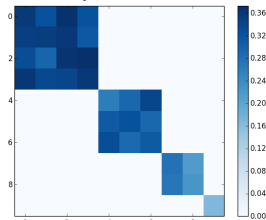


Numerical experiment: dimension 10, 210 parameters

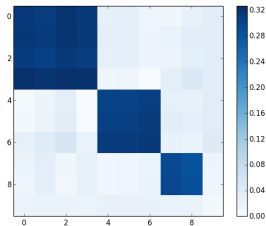
No penalization



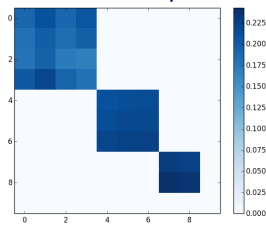
ℓ_1 penalization



trace-norm penalization

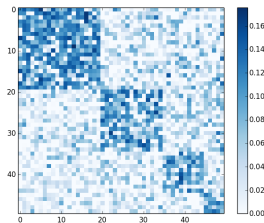


$\ell_1 + \text{trace norm}$ penalization

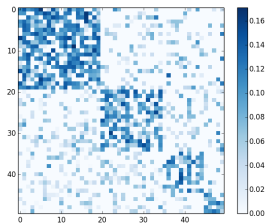


Numerical experiment: dimension 100, 20100 parameters

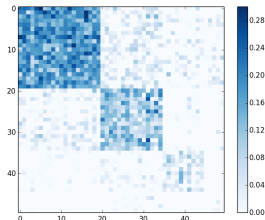
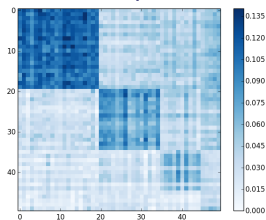
No penalization



ℓ_1 penalization



trace-norm penalization $\ell_1 + \text{trace norm penalization}$



Theoretical guarantees

A simplified framework: fix a set $\{h_{j,k} : 1 \leq j, k \leq d\}$ and intensities

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

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

Instead of $-\log$ likelihood, consider least squares

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

where $\|\lambda_\theta\|_T^2 = \langle \lambda_\theta, \lambda_\theta \rangle_T$ with

$$\langle \lambda_\theta, \lambda_{\theta'} \rangle_T = \frac{1}{T} \sum_{j=1}^d \int_{[0,T]} \lambda_{j,\theta}(t) \lambda_{j,\theta'}(t) dt.$$

Least-squares goodness-of-fit

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

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

$$\mathbb{E}[R_T(\theta)] = \mathbb{E}\|\lambda_\theta\|_T^2 - 2\mathbb{E}\langle \lambda_\theta, \lambda^* \rangle_T = \mathbb{E}\|\lambda_\theta - \lambda^*\|_T^2 - \|\lambda^*\|_T$$

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

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

where M_j martingale

Introduce

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

with

$$\text{pen}(\theta) = \|\mu\|_{1, \hat{w}} + \|\mathbf{A}\|_{1, \hat{w}} + \hat{w}_* \|\mathbf{A}\|_*$$

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

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

$$\begin{aligned} \|\lambda_{\hat{\theta}} - \lambda^*\|_T^2 + \|\lambda_{\hat{\theta}} - \lambda_{\theta}\|_T^2 - \|\lambda_{\theta} - \lambda^*\|_T^2 \\ \leq -\langle \theta_{\partial}, \hat{\theta} - \theta \rangle + \frac{2}{T} \langle \hat{\mu} - \mu, \bar{M}_T \rangle + \frac{2}{T} \langle \hat{\mathbf{A}} - \mathbf{A}, \mathbf{Z}_T \rangle, \end{aligned}$$

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

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

$$(\mathbf{Z}_t)_{j,k} = \int_0^t \int_{(0,s)} h_{j,k}(s-u) dN_k(u) dM_j(s), \quad (2)$$

or

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

with \mathbf{H}_t predictable process with entries

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

Noise term is a matrix-martingale in continuous time:

$$\frac{1}{T} \mathbf{Z}_T$$

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

A consequence of our new concentration inequalities:

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

for any $v, x, b > 0$, where

$$\mathbf{V}_t = \frac{1}{t} \int_0^t \|\mathbf{H}_s\|_{2,\infty}^2 \begin{bmatrix} \text{diag}[\lambda_s^*] & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_s^\top \text{diag}[\mathbf{H}_s \mathbf{H}_s^\top]^{-1} \text{diag}[\lambda_s^*] \mathbf{H}_s \end{bmatrix} ds$$

and $b_t = \sup_{s \in [0,t]} \|\mathbf{H}_s\|_{2,\infty}$ ($\|\cdot\|_{2,\infty}$ = maximum ℓ_2 row norm)

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

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

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

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

$$\hat{\mathbf{V}}_t = \frac{1}{t} \int_0^t \|\mathbf{H}_s\|_{2,\infty}^2 \begin{bmatrix} \text{diag}[dN_s] & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_s^\top \text{diag}[\mathbf{H}_s \mathbf{H}_s^\top]^{-1} \text{diag}[dN_s] \mathbf{H}_s \end{bmatrix} ds$$

and small ugly term:

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

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 $\|\mathbf{Z}_T\|_\infty = \max_{j,k} |A_{j,k}|$ and $\|\mathbf{Z}_T\|_{\text{op}}$ 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 \mathbf{A} : $\|\mathbf{A}\|_{1,\hat{\mathbf{W}}} = \sum_{1 \leq j,k \leq d} \hat{\mathbf{W}}_{j,k} |\mathbf{A}_{j,k}|$ with

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

where

$$\mathbf{B}_{j,k}(t) = \sup_{s \in [0,t]} \int_{(0,t)} h_{j,k}(t-s) dN_k(s) \\ \hat{\mathbf{V}}_{j,k}(t) = \frac{1}{t} \int_0^t \left(\int_{(0,s)} h_{j,k}(s-u) dN_k(u) \right)^2 dN_j(s)$$

namely

$$\hat{\mathbf{W}}_{j,k} \approx c \sqrt{\frac{\hat{\mathbf{V}}_{j,k}(T)}{T}}$$

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

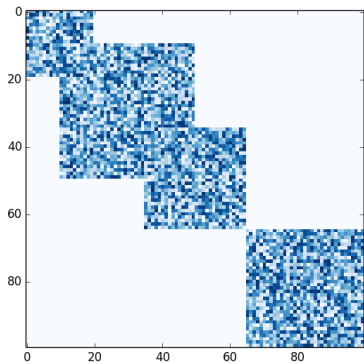
For trace-norm penalization of \mathbf{A} : $\hat{w}_* \|\mathbf{A}\|_*$ with

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

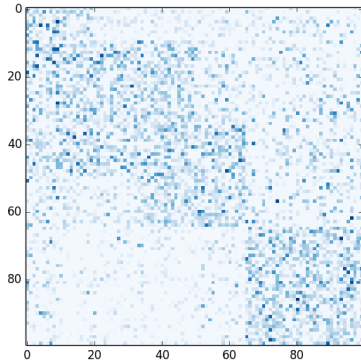
namely

$$\hat{w}_* \approx \sqrt{\frac{\lambda_{\max}(\hat{\mathbf{V}}_T)}{T}}$$

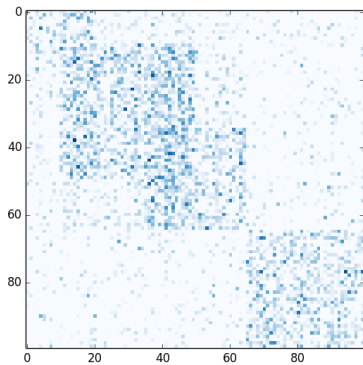
- Data-driven weights that comes from “empirical” Bernstein's inequalities, entrywise and for operator norm of \mathbf{Z}_T
- $\hat{\mathbf{V}}_{j,k}(t)$ and $\lambda_{\max}(\hat{\mathbf{V}}_t)$ are estimations (based on optional variation) of the variance terms from Bernstein's inequality
- $\mathbf{B}_{j,k}(t)$ 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



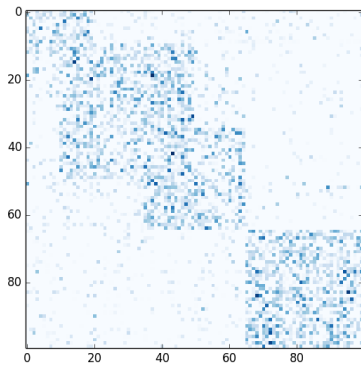
Truth



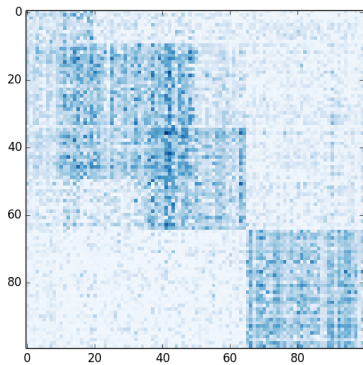
No pen



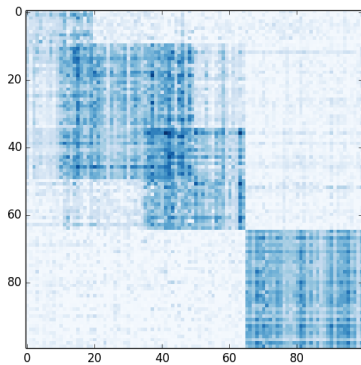
L1



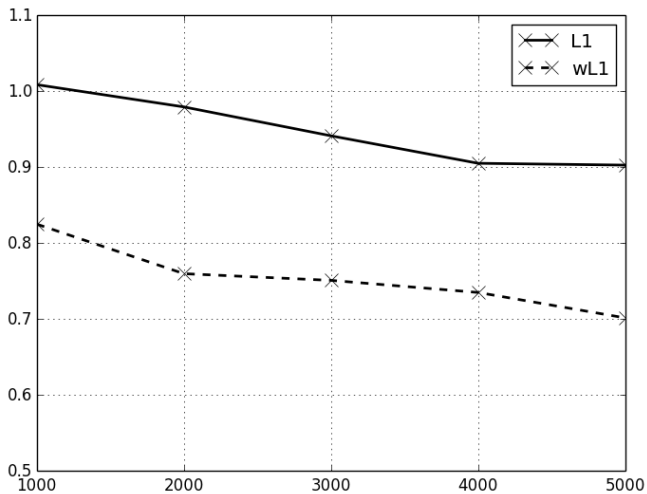
wL1



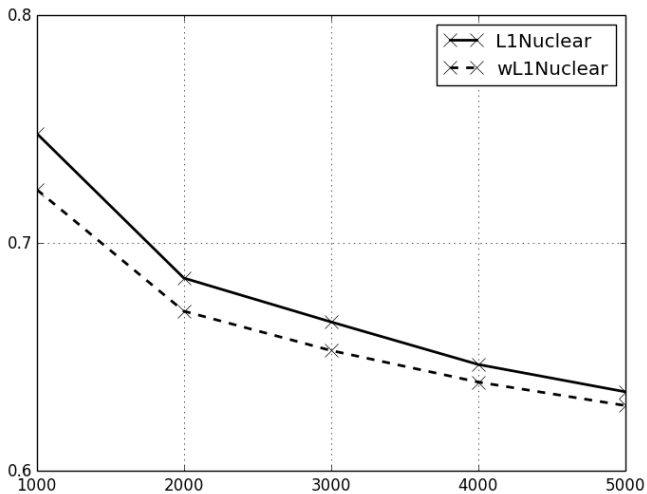
Nuclear



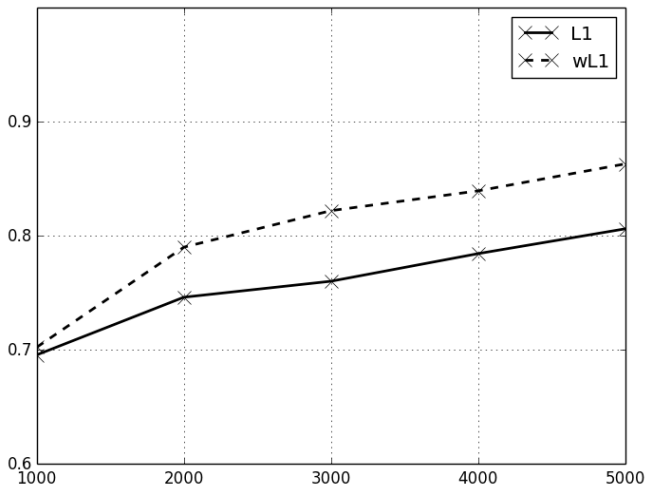
wNuclear



Error for L1 and wL1



Error for L1Nuclear and wL1Nuclear



AUC (support selection) for L1Nuclear and wL1Nuclear

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

$$\begin{aligned} \|\lambda_{\hat{\theta}} - \lambda^*\|_T^2 \leq \inf_{\theta} \left\{ \|\lambda_{\theta} - \lambda^*\|_T^2 + \kappa(\theta)^2 \left(\frac{5}{4} \|(\hat{\mathbf{W}})_{\text{supp}(\mu)}\|_2^2 \right. \right. \\ \left. \left. + \frac{9}{8} \|(\hat{\mathbf{W}})_{\text{supp}(\mathbf{A})}\|_F^2 + \frac{9}{8} \hat{w}_*^2 \text{rank}(\mathbf{A}) \right) \right\} \end{aligned}$$

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

- Leading constant 1

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

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

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

Concentration inequalities for matrix martingales in continuous time

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

Introduce

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

where $\{\mathbf{A}_t\}$, $\{\mathbf{C}_t\}$ and $\{\mathbf{B}_t\}$ predictable and where $\{\mathbf{M}_t\}_{t \geq 0}$ is a “white” matrix martingale, in the sense that $[\text{vec} \mathbf{M}]_t$ is diagonal

NB: entries of \mathbf{Z}_t are given by

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

- $\langle \mathbf{M} \rangle_t$ = entrywise predictable quadratic variation, so that

$$\mathbf{M}_t^{\odot 2} - \langle \mathbf{M} \rangle_t$$

martingale

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

$$\text{vec}(\mathbf{M}_t) \text{vec}(\mathbf{M}_t)^\top - \langle \text{vec} \mathbf{M} \rangle_t$$

is a martingale.

- $\mathbf{M}_t = \mathbf{M}_t^c + \mathbf{M}_t^d$, where \mathbf{M}_t^c is a continuous martingale and \mathbf{M}_t^d is a purely discontinuous martingale. Its (entrywise) quadratic variation is defined as

$$[\mathbf{M}]_t = \langle \mathbf{M}^c \rangle_t + \sum_{0 \leq s \leq t} (\Delta \mathbf{M}_s)^2, \quad (3)$$

and its quadratic covariation by

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

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

Concentration for purely discontinuous matrix martingale:

- \mathbf{M}_t is purely discontinuous and we have

$$\langle \mathbf{M} \rangle_t = \int_0^t \lambda_s ds$$

for a non-negative and predictable intensity process $\{\lambda_t\}_{t \geq 0}$.

- Standard moment assumptions (subexponential tails)

Introduce

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

where

$$\mathbf{W}_t = \begin{bmatrix} \mathbf{W}_t^1 & \mathbf{0} \\ \mathbf{0} & \mathbf{W}_t^2 \end{bmatrix}, \quad (4)$$

$$\mathbf{W}_t^1 = \mathbf{A}_t \operatorname{diag}[\mathbf{A}_t^\top \mathbf{A}_t]^{-1} \operatorname{diag}[(\mathbf{C}_t^{\odot 2} \odot \lambda_t) \mathbf{1}] \mathbf{A}_t^\top$$

$$\mathbf{W}_t^2 = \mathbf{B}_t^\top \operatorname{diag}[\mathbf{B}_t \mathbf{B}_t^\top]^{-1} \operatorname{diag}[(\mathbf{C}_t^{\odot 2} \odot \lambda_t)^\top \mathbf{1}] \mathbf{B}_t$$

Introduce also

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

Theorem.

$$\mathbb{P} \left[\|\mathbf{Z}_t\|_{\text{op}} \geq \sqrt{2v(x + \log(m+n))} + \frac{b(x + \log(m+n))}{3}, \right. \\ \left. b_t \leq b, \quad \lambda_{\max}(\mathbf{V}_t) \leq v \right] \leq e^{-x},$$

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

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

$$\begin{aligned} & \left\| \int_0^t \mathbf{C}_s \odot d(\mathbf{N}_t - \mathbf{\Lambda}_t) \right\|_{\text{op}} \\ & \leq \sqrt{2 \left(\left\| \int_0^t \mathbf{C}_s^{\odot 2} \odot \lambda_s ds \right\|_{1,\infty} \vee \left\| \int_0^t \mathbf{C}_s^{\odot 2} \odot \lambda_s ds \right\|_{\infty,1} \right) (x + \log(p+q))} \\ & \quad + \frac{\sup_{s \in [0,t]} \|\mathbf{C}_s\|_{\infty} (x + \log(p+q))}{3} \end{aligned}$$

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

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

$$\|\mathbf{N} - \boldsymbol{\lambda}\|_{\text{op}} \leq \sqrt{2(\|\boldsymbol{\lambda}\|_{1,\infty} \vee \|\boldsymbol{\lambda}\|_{\infty,1})(x + \log(p + q))} \\ + \frac{x + \log(p + q)}{3}.$$

- Up to our knowledge, not previously stated in literature
- NB: In the Gaussian case: variance depends on maximum ℓ_2 norm of rows and columns (cf. Tropp (2011))

- We have as well a non-commutative Hoeffding's inequality when \mathbf{M}_t has continuous paths, 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

$$\mathbf{U}_t^u = \sum_{0 \leq s \leq t} \left(e^{u \Delta \mathcal{J}(\mathbf{Z}_s)} - u \Delta \mathcal{J}(\mathbf{Z}_s) - I \right)$$

given by

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

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

- Theoretical study of learning algorithms for “time-oriented” models need new probabilistic results
- In our case new concentration results for matrix martingales in continuous time
- Solves the scaling problem of penalizations

Perspectives:

- Experiments on Twitter, BlogoSphere and High-frequency Finance (ongoing)
- Superposition of Hawkes for viral diffusion of contents
- Better solvers using stochastic gradient based algorithms