

# Machine learning and time: “time accounting” learning

Stéphane Gaïffas<sup>1</sup>

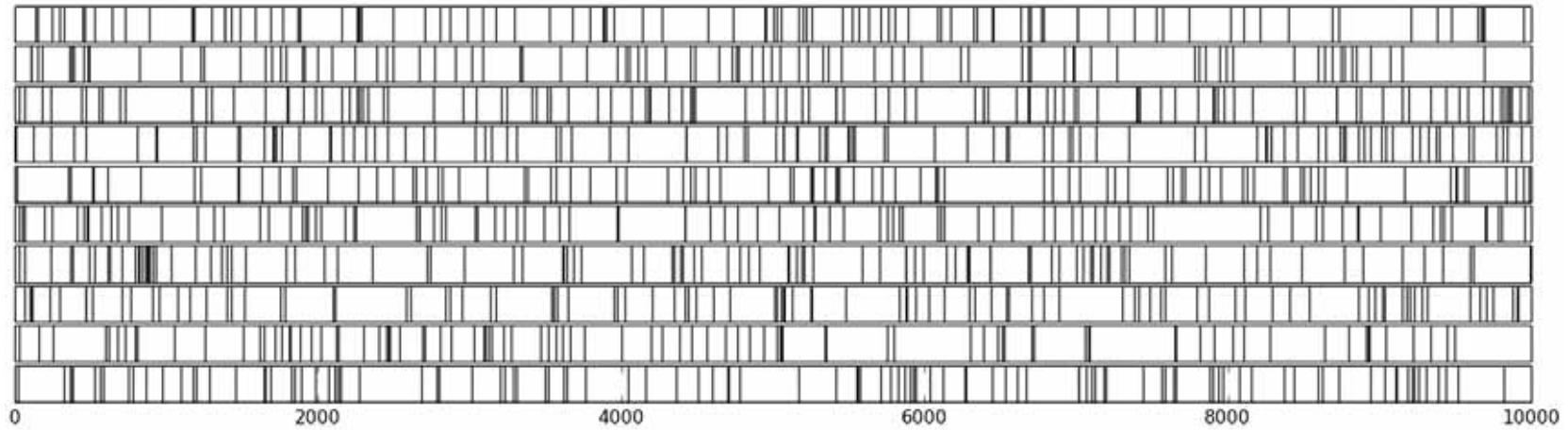


14 mars 2014

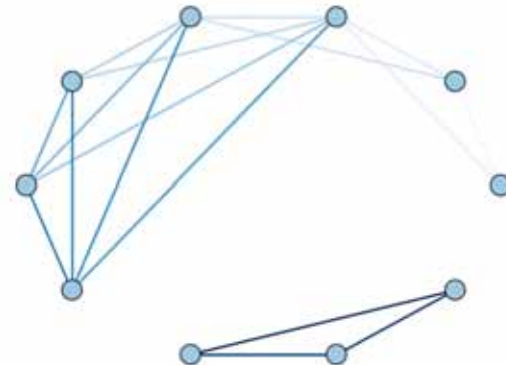
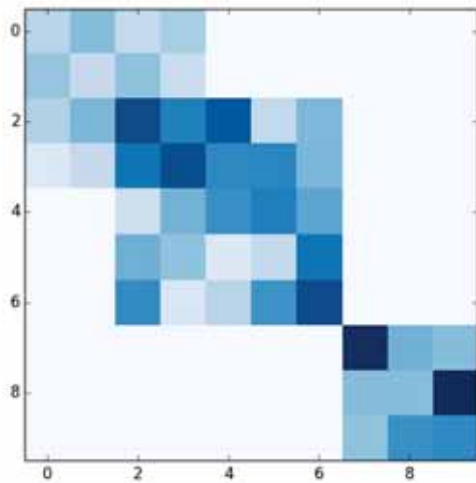
- Finite network with nodes  $\{1, \dots, d\}$ : users of a social network, of an e-commerce platform, etc.
- For each node  $j \in \{1, \dots, d\}$  we observe the timestamps  $\{t_{j,1}, t_{j,2}, \dots\}$  of nodes' actions
- Goal: recover **levels of interactions** between users based on the timestamps patterns

# Introduction

From



Quantify interactions between users



- 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

# Model: Multivariate Hawkes Process (MHP)

- Counting process  $N_j(t) = \sum_{i \geq 1} \mathbf{1}_{t_j, i \leq t}$
- Data: 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  $\lambda_j$ , namely

$$\begin{aligned} & \mathbb{P}(j \text{ does something at time } t \text{ knowing the past}) \\ &= \mathbb{P}(N_j \text{ has a jump in } [t, t + dt] \mid \mathcal{F}_t) = \lambda_j(t)dt \end{aligned}$$

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

# Model: Multivariate Hawkes Process (MHP)

- MHP assumes an autoregressive structure on the intensities:

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

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

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

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

- Notation:

$$\int_{(0,t)} \varphi(t-s) dN_j(s) = \sum_{i: t_{j,i} < t} \varphi(t - t_{j,i})$$

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)

## 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

Dimension  $d$  is large:

- Need a simple parametric model on  $\mu$  and  $\varphi$
- We want a **convex** optimization problem with smooth loss
- We want to encode some prior assumptions by penalizing this loss

# 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

# 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

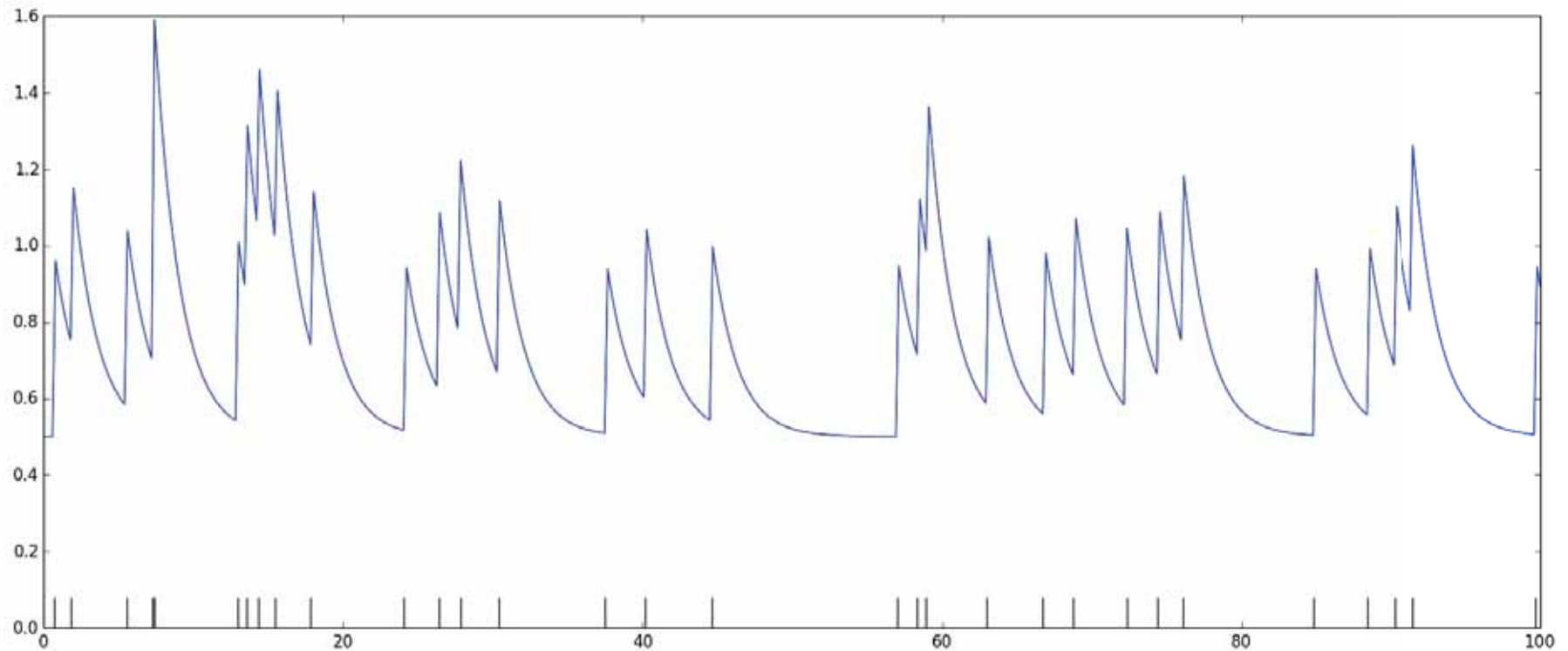
$$\theta = [\mu, \mathbf{A}, \alpha]$$

with

- baselines  $\mu = [\mu_1, \dots, \mu_d]^\top \in \mathbb{R}_+^d$
- adjacencies  $\mathbf{A} = [a_{j,k}]_{1 \leq j, k \leq d} \in \mathbb{R}_+^{d \times d}$
- 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  $\lambda_\theta$  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  $\theta = (\boldsymbol{\mu}, \mathbf{A}, \boldsymbol{\alpha})$  with  $\boldsymbol{\mu} = [\mu_j]$ ,  $\mathbf{A} = [A_{j,k}]$ ,  $\boldsymbol{\alpha} = [\alpha_{j,k}]$

## Prior assumptions

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

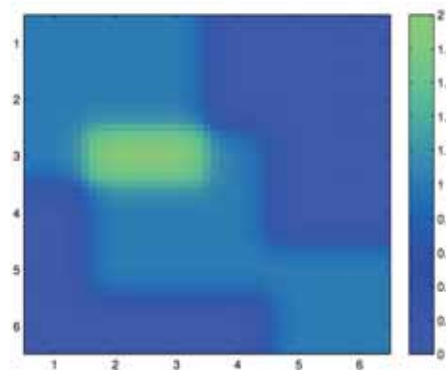
$\mu$  is sparse

- 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  $\alpha$  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, ...]

- Tightest convex relaxation of  $\|\mathbf{A}\|_0 = \sum_{j,k} \mathbf{1}_{\mathbf{A}_{j,k} > 0}$  is  $\ell_1$ -norm:

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

- Tightest convex relaxation of rank is trace-norm:

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

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



So, we use the following penalizations

- Use  $\ell_1$  penalization on  $\mu$
- Use  $\ell_1$  penalization on  $\mathbf{A}$
- Use trace-norm penalization on  $\mathbf{A}$
- Use  $\ell_2^2$  penalization on  $\alpha$ , weighted by  $\mathbf{A}$

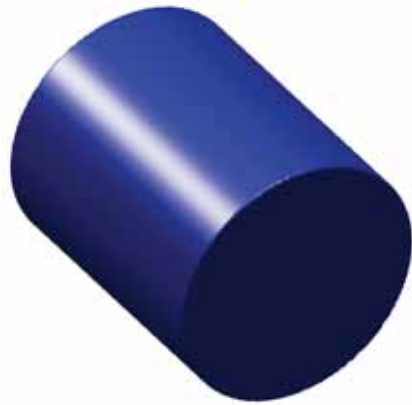
[but other choices might be interesting...]

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

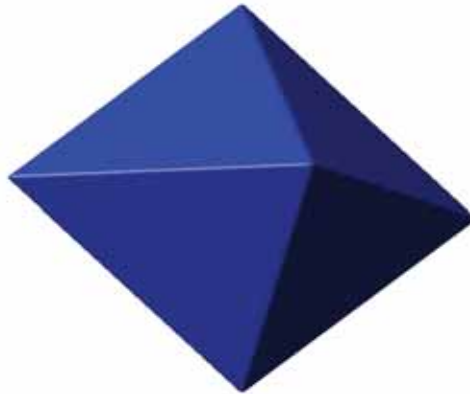
$$\mathbf{A} \mapsto w_* \|\mathbf{A}\|_* + w_1 \|\mathbf{A}\|_1$$

NB2: recent works by Richard et al (2013, 2014): better way to induce sparsity and low-rank than the sum, but not-scalable / non-convex

# 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  $2 \times 2$  symmetric matrices, which is identified with  $\mathbb{R}^3$ .

[show video]

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

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

# Penalized maximum likelihood: reparametrization

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  $\beta = \mathbf{A} \circ \alpha$ , 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

# Penalized maximum likelihood: reparametrization

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

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

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

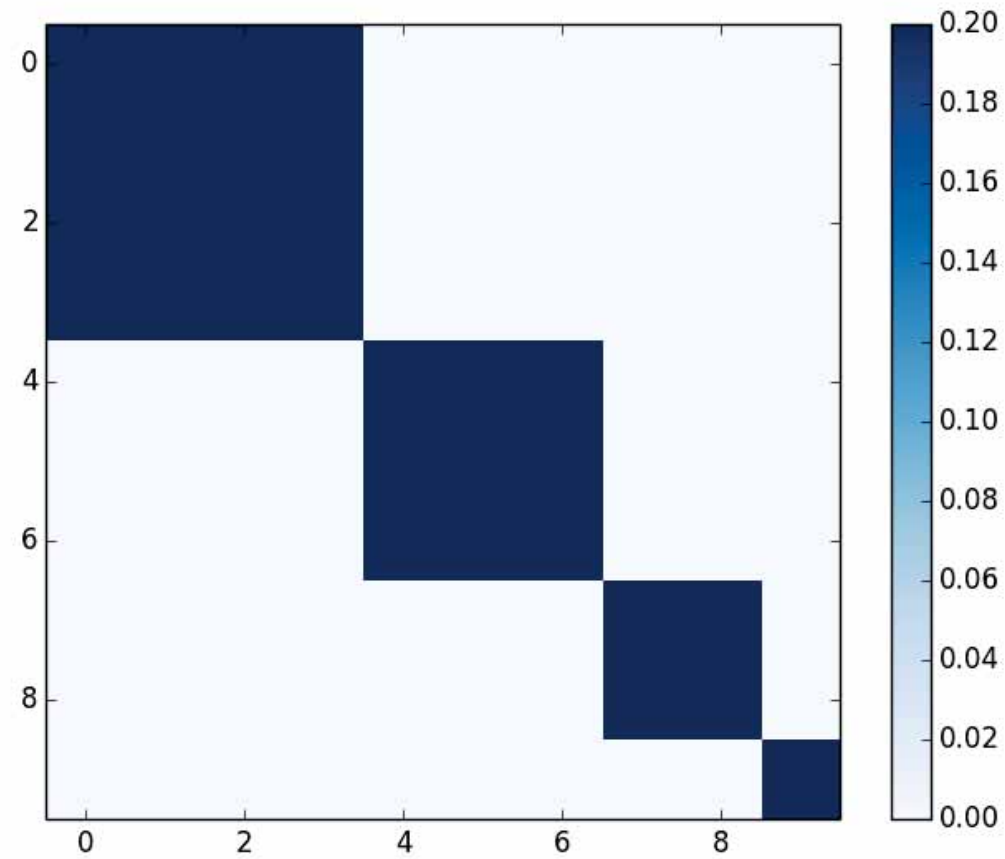
$$\lambda_{j,\theta}(t) = \mu_j + \sum_{k=1}^d a_{j,k} \int_{(0,t)} \exp \left( -\frac{\beta_{j,k}}{a_{j,k}} (t-s) \right) dN_k(s)$$

# Convex optimization – numerical aspects

- Can be solved using first-order routines:  
Fista [Beck Teboulle (2009)], Prisma [Orabona et al (2012)], GFB [Peyre et al. (2011)], Primal-Dual [Chambolle et al. (2009), Condat et al. (2013)], ADMM [Boyd (2012)], etc...
- Gradient of  $-\ell_{\mathcal{T}}(\theta)$  using a recursion formula  
→ Naively  $O(n^2d)$  with  $n =$  number of events (very large) but  $O(nd)$  when careful (using recursion formulas)  
→ Parallelized code for this: gradient of each node  $j \in \{1, \dots, d\}$  computed **in parallel**
- Computation bottleneck: exp and log, accelerated using ugly hacking
- Trace norm penalization, truncated SVD: default's Lanczos's implementation of Python is fast enough for  $d \approx 1K$ , use a non-convex factorized formulation  $\mathbf{A} = \mathbf{U}\mathbf{V}^{\top}$  for  $d \gg 1K$

# Numerical experiment

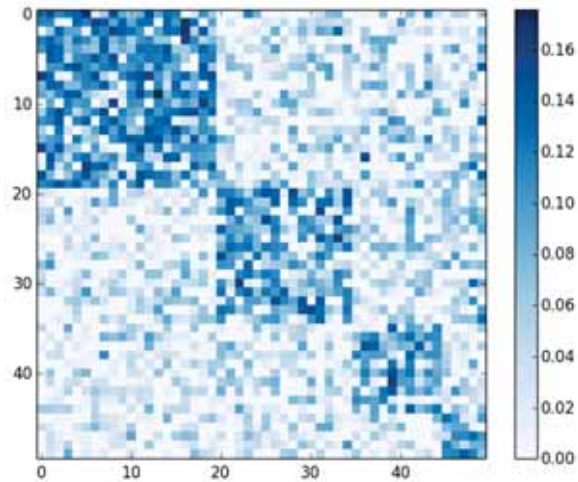
Toy example: take matrix  $\mathbf{A}$  as



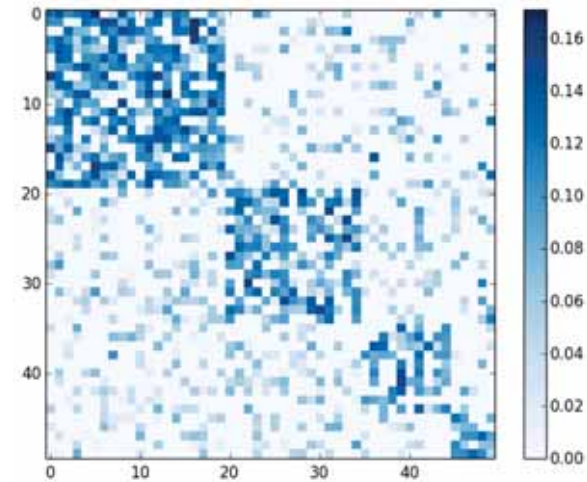


# Numerical experiment: dimension 100, 20100 parameters

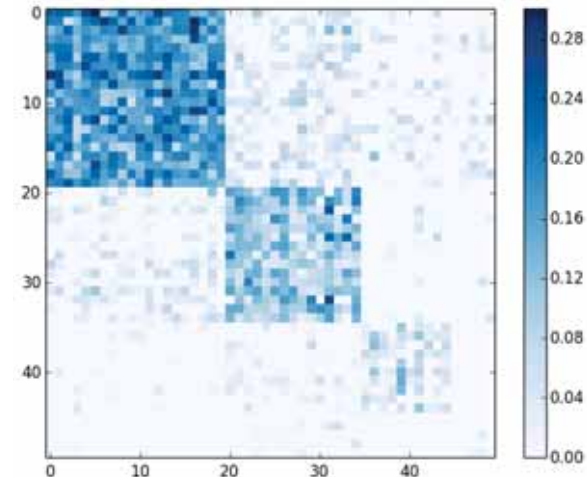
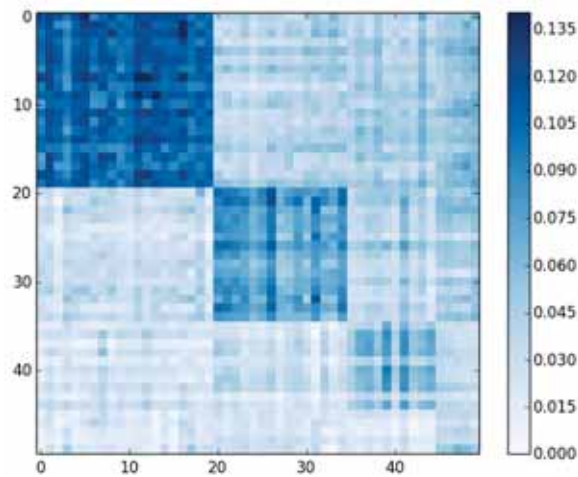
No penalization



$\ell_1$  penalization



trace-norm penalization  $\ell_1 +$  trace norm penalization



# Some theory: sharp oracle inequalities

We consider 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

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

# Some theory: sharp oracle inequalities

Introduce

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

with

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

- Penalization tuned by data-driven weights  $\hat{w}$ ,  $\hat{\mathbf{W}}$  and  $\hat{w}_*$
- Comes from sharp controls of the noise terms
- Solves the **scaling** problem for this model (e.g. feature scaling)

# Solving the “feature scaling” problem

$\ell_1$ -penalization of  $\mu$

$$\|\mu\|_{1,\hat{w}} = \sum_{j=1}^d \hat{w}_j |\mu_j|$$

with

$$\hat{w}_j \approx \sqrt{\frac{(x + \log d) N_j([0, T]) / T}{T}}$$

where  $N_j([0, T]) = \#$  events for node  $j$

- Each  $\mu_j$  penalized by its average events intensity

# Solving the “feature scaling” problem

$\ell_1$ -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

$$\hat{\mathbf{W}}_{j,k} \approx \sqrt{\frac{(x + \log d) \hat{\mathbf{V}}_{j,k}(T)}{T}}$$

where

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

= variance estimation of the self-excitement for  $k \rightarrow j$

# Solving the “feature scaling” problem

Trace-norm penalization of  $\mathbf{A}$  [difficult]

$$\hat{w}_* \|\mathbf{A}\|_* = \hat{w}_* \sum_{j=1}^d \sigma_j(\mathbf{A})$$

with

$$\hat{w}_* \approx \sqrt{\frac{(x + \log d)(\|\hat{\mathbf{V}}_1(T)\|_{\text{op}} \vee \|\hat{\mathbf{V}}_2(T)\|_{\text{op}})}{T}}$$

where  $\|\cdot\|_{\text{op}} = \text{operator norm}$

# Solving the “feature scaling” problem

and where  $\hat{\mathbf{V}}_1(t)$  diagonal matrix with entries

$$(\hat{\mathbf{V}}_1(t))_{j,j} = \frac{1}{t} \int_0^t \|\mathbf{H}(s)\|_{2,\infty}^2 dN_j(s),$$

$\hat{\mathbf{V}}_2(t)$  matrix with entries

$$(\hat{\mathbf{V}}_2(t))_{j,k} = \frac{1}{t} \int_0^t \|\mathbf{H}(s)\|_{2,\infty}^2 \sum_{l=1}^d \frac{H_{j,l}(s)H_{k,l}(s)}{\|\mathbf{H}_{l,\bullet}(s)\|_2^2} dN_l(s),$$

with  $\|\cdot\|_{2,\infty} = \text{maximum } \ell_2 \text{ row norm}$  and  $\mathbf{H}(t)$  matrix with entries

$$\mathbf{H}_{j,k}(t) = \int_{(0,t)} h_{j,k}(t-s) dN_k(s)$$

# Solving the “feature scaling” problem

- $\hat{\mathbf{V}}_{j,k}(t)$ ,  $\|\hat{\mathbf{V}}_1(t)\|_{\text{op}}$  and  $\|\hat{\mathbf{V}}_2(t)\|_{\text{op}}$  are estimations (based on optional variation) of non-observable variance terms
- It comes from **new Bernstein’s concentration inequalities**, used on the noise term
- We develop a new probabilistic tool: **non-commutative concentration inequality for random matrix martingales in continuous time** (theory given by Tropp (2011) applies to discrete time only, and depend on unobserved variance terms)

These tools give a **sharp data-driven tuning** of the penalizations, solving the scaling problem



# Some theory: sharp oracle inequalities

Define

$$\langle \lambda_1, \lambda_2 \rangle_T = \frac{1}{T} \sum_{j=1}^d \int_0^T \lambda_{1,j}(t) \lambda_{2,j}(t) dt$$

and  $\|\lambda\|_T^2 = \langle \lambda, \lambda \rangle_T$

- We use a standard assumption to obtain fast rates for the Lasso: the RE (Restricted Eigenvalue) Assumption [Bickel et al. (2009), Koltchinskii (2011), ...]

## Theorem 1

We have

$$\|\lambda_{\hat{\theta}} - \lambda_0\|_T^2 \leq \inf_{\theta} \left\{ \|\lambda_{\theta} - \lambda_0\|_T^2 + \kappa(\theta)^2 \left( \frac{5}{4} \|(\hat{\mathbf{W}})_{\text{supp}(\mu)}\|_2^2 + \frac{9}{8} \|(\hat{\mathbf{W}})_{\text{supp}(\mathbf{A})}\|_F^2 + \frac{9}{8} \hat{w}_*^2 \text{rank}(\mathbf{A}) \right) \right\}$$

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

- $\kappa(\theta)$ : RE constant
- Sharp: leading constant 1

## Some theory: sharp oracle inequalities

Take-home message:  $\hat{\theta}$  achieves an optimal tradeoff between approximation and complexity given by

$$\begin{aligned} & \frac{\|\mu\|_0(x + \log d)}{T} \max_j N_j([0, T]) / T \\ & + \frac{\|\mathbf{A}\|_0(x + \log d)}{T} \max_{j,k} \hat{\mathbf{V}}_{j,k}(T) \\ & + \frac{\text{rank}(A)(x + \log d)}{T} \|\hat{\mathbf{V}}_1(T)\|_{\text{op}} \vee \|\hat{\mathbf{V}}_2(T)\|_{\text{op}}. \end{aligned}$$

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

# A new concentration inequality

- New Bernstein's empirical concentration inequality for continuous-time matrix martingale
- Consider the random matrix  $\mathbf{Z}(t)$  with entries

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

where  $M_j(t) = N_j(t) - \int_0^t \lambda_j(s) ds$  are martingales obtained by compensation

- This is the noise term in our problem

# A new concentration inequality

A classical concentration inequality for  $\mathbf{Z}_{j,k}$  [Lipster Shiriyayev 1986] is

$$\frac{1}{t}(\mathbf{Z}(t))_{j,k} \leq \sqrt{\frac{2vx}{t}} + \frac{bx}{3t}$$

for any  $x > 0$ , with a probability  $\geq 1 - e^{-x}$  whenever

$$\frac{1}{t}\langle \mathbf{Z}_{j,k} \rangle_t = \frac{1}{t} \int_0^t \left( \int_{(0,s)} h_{j,k}(s-u) dN_k(u) \right)^2 \lambda_j(s) ds \leq v$$

and

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

# A new concentration inequality

- Predictable variation  $\langle \mathbf{Z}_{j,k} \rangle_t$  depends on non-observed  $\lambda_j$ : this concentration is **useless** for statistics
- Need an **empirical Bernstein's inequality**, with a variance term using the **optional variation**

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

- We need also to remove the event  $\{\langle \mathbf{Z}_{j,k} \rangle_t \leq tv\}$  from this inequality

We provide:

- A control of all the entries  $\mathbf{Z}_{j,k}$  of  $\mathbf{Z}$
- A control of  $\|\mathbf{Z}_t\|_{\text{op}}$

## Theorem 2

We have

$$\frac{1}{t} |\mathbf{Z}_{j,k}(t)| \leq 2\sqrt{2} \sqrt{\frac{(x + 2 \log d + \hat{\mathbf{L}}_{j,k}(t)) \hat{\mathbf{V}}_{j,k}(t)}{t}} + 9.31 \frac{(x + 2 \log d + \hat{\mathbf{L}}_{j,k}(t)) \mathbf{B}_{j,k}(t)}{t}$$

for any  $1 \leq j, k \leq d$ , with a probability larger than  $1 - 30.55e^{-x}$ .

- Based on a previous result by G. and Guilloux (2011), see also Hansen et al (2012)
- Reminiscent of previous works by Audibert (2008)

## Theorem 3

For any  $x > 0$ , we have

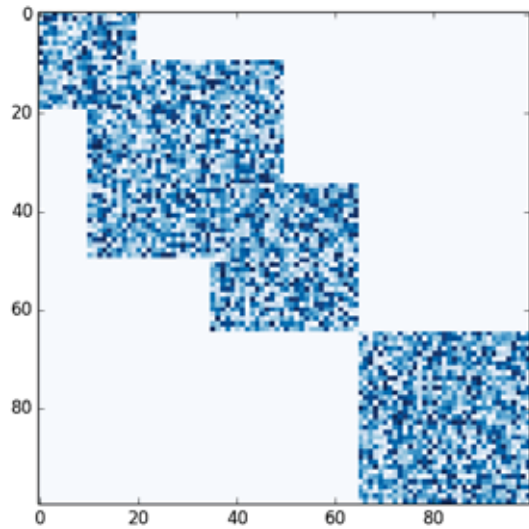
$$\begin{aligned} & \frac{\|\mathbf{Z}(t)\|_{\text{op}}}{t} \\ & \leq 4\sqrt{\frac{(x + \log d + \hat{\ell}_x(t))\|\hat{\mathbf{V}}_1(t)\|_{\text{op}} \vee \|\hat{\mathbf{V}}_2(t)\|_{\text{op}}}{t}} \\ & \quad + \frac{(x + \log d + \hat{\ell}_x(t))(10.34 + 2.65 \sup_{t \in [0, T]} \|\mathbf{H}(t)\|_{2, \infty})}{t} \end{aligned}$$

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

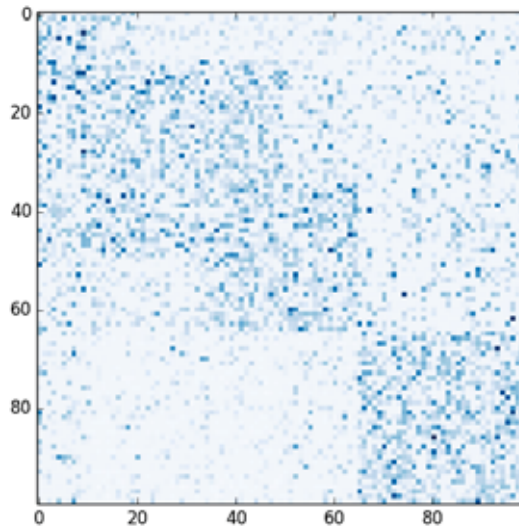
- First non-commutative Bernstein's inequality for continuous time martingales
- Can be extended to a wider class of martingales
- Extension of [Tropp (2012)] results



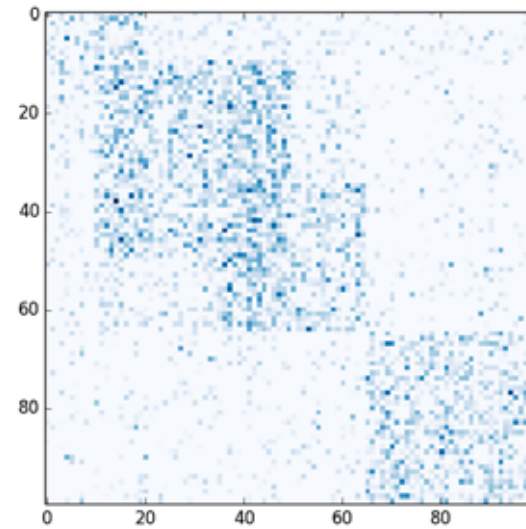
# Consequence: a sharp scaling of penalizations



Ground Truth



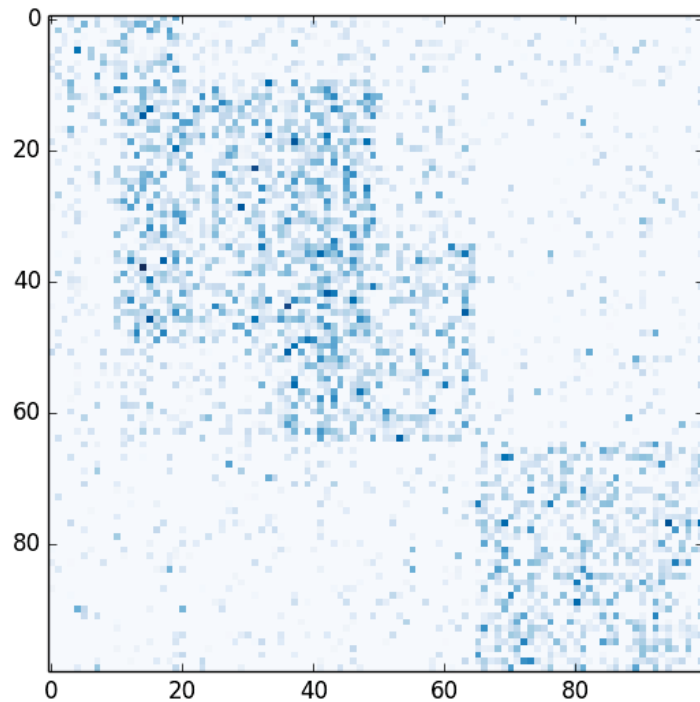
NoPen



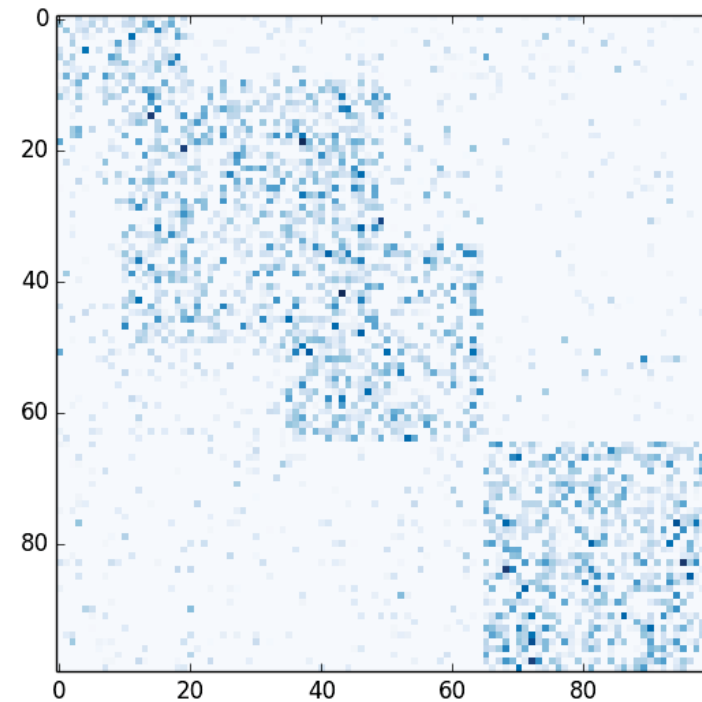
L1

# Consequence: a sharp scaling of penalizations

## L1 vs wL1



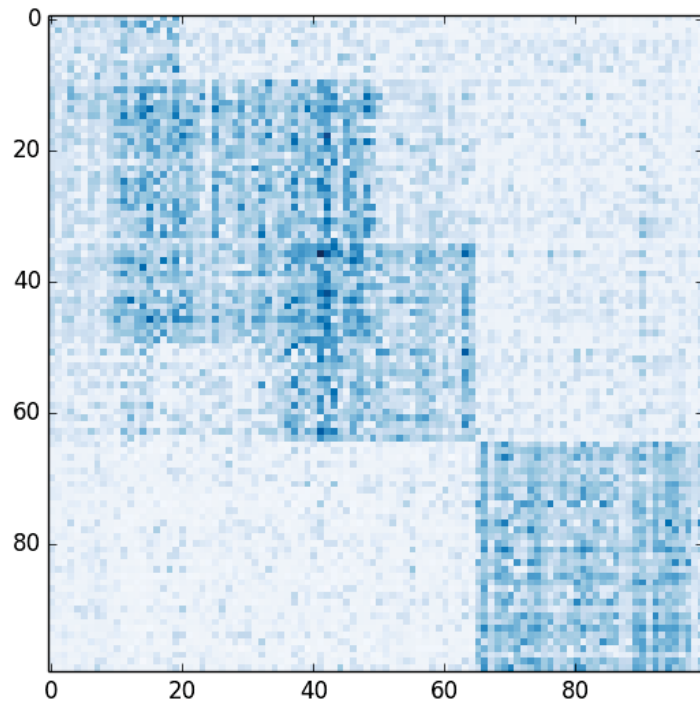
L1



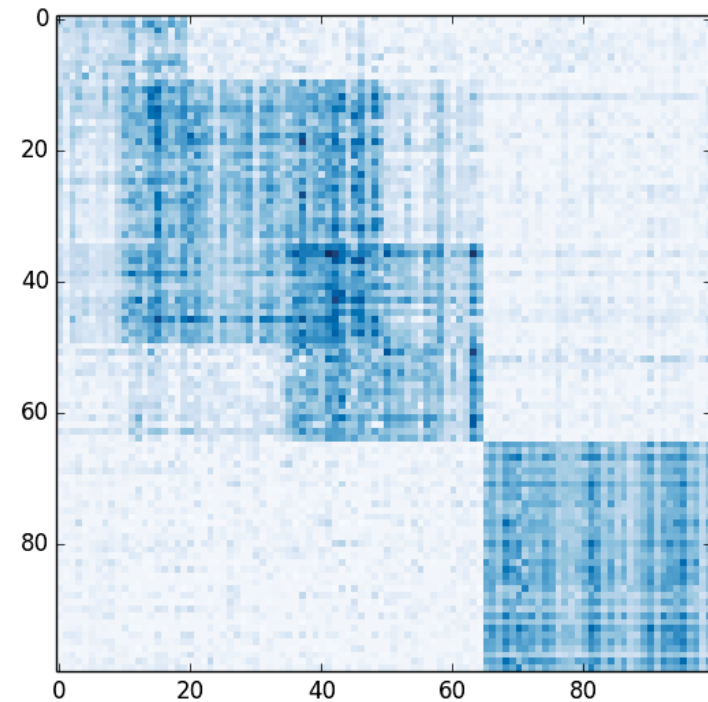
wL1

# Consequence: a sharp scaling of penalizations

## Nuclear vs wNuclear

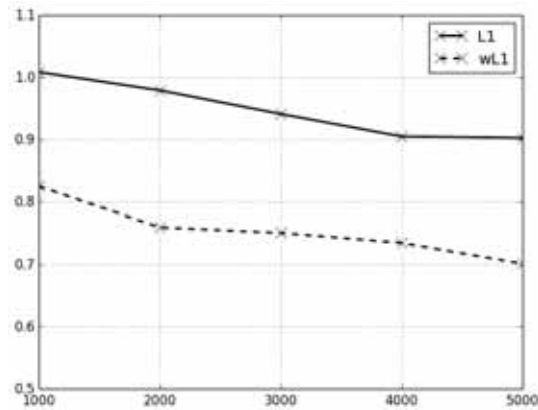


L1Nuclear

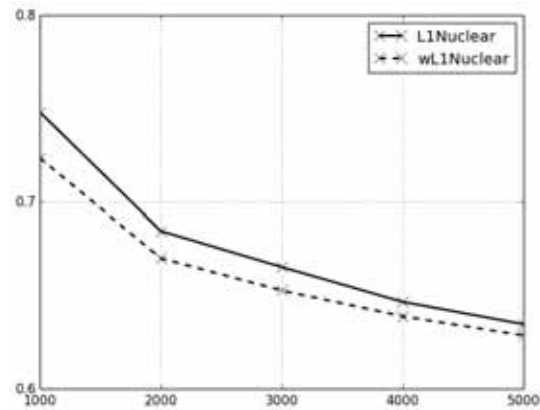


wL1Nuclear

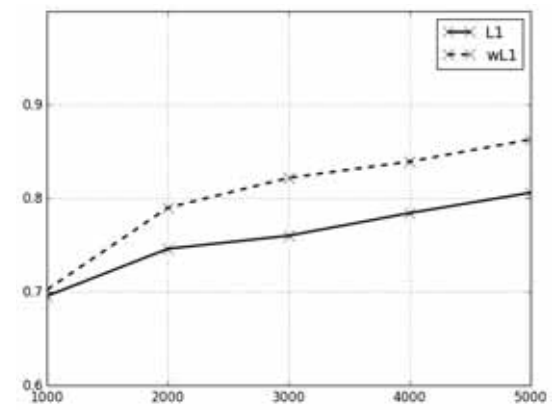
# Consequence: a sharp scaling of penalizations



Error for L1 and wL1

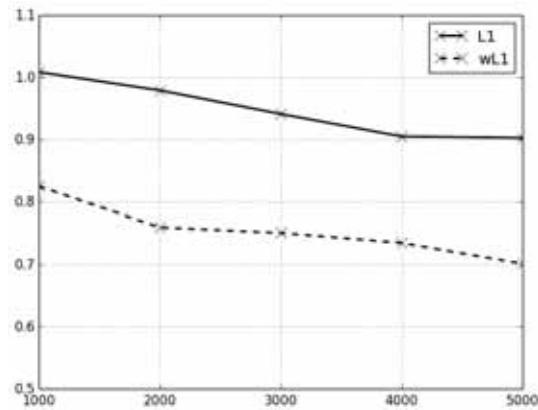


Error for L1Nuclear and wL1Nuclear

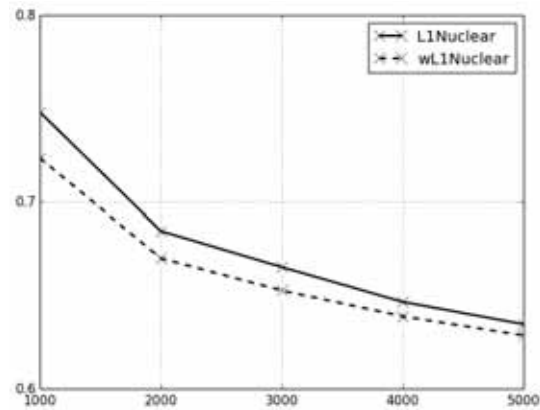


AUC for L1 and wL1

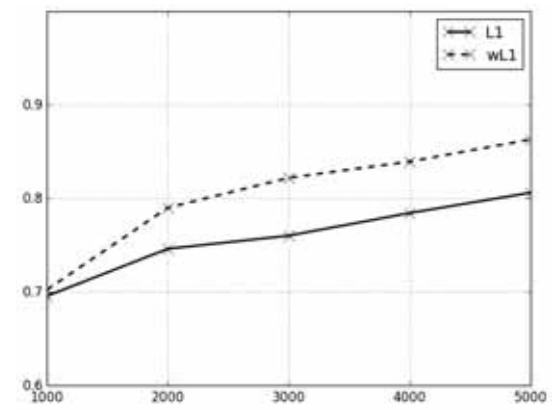
# Consequence: a sharp scaling of penalizations



Error for L1 and wL1



Error for L1Nuclear and wL1Nuclear



AUC for L1 and wL1

# Take-home message

- Reparametrization of the problem
- Theoretical analysis gives insight to choose the correct scaling of the penalizations
- First oracle inequality for this problem
- This required new probabilistic tools for matrix martingales in continuous time

- Larger scale: factorized form  $\mathbf{A} = \mathbf{UV}^\top$
- Incorporation of features (text, time-varying graph-features, etc.)
- Time varying baseline  $\mu(t)$  for non-stationarity

Thank you!