

博士論文

論文題目

**Penalized Least Squares Approximation
Methods and Their Applications to
Stochastic Processes**

(罰則付き最小二乗近似法とその確率過程への応用)

氏名 鈴木 拓海

**Penalized Least Squares
Approximation Methods and Their
Applications to Stochastic Processes**

SUZUKI Takumi

January 8, 2020

Acknowledgements

I would like to express my sincere thanks to Professor Nakahiro Yoshida for giving me many valuable instructions and advices. It was a great pleasure in my life that I had received his solicitous education and generous support.

Contents

1	Introduction	4
1.1	A brief review of sparse estimation	4
1.2	A brief review of statistical inference for stochastic processes	7
1.3	Organization of this thesis	8
2	Penalized LSA estimator	10
2.1	Definition of penalized LSA estimator	10
2.2	Main theorem	12
2.3	P-O estimator	14
2.4	Proofs of main theorems	16
2.4.1	Proof of Theorem 2.1	16
2.4.2	Proof of Theorem 2.2	17
2.4.3	Proof of Theorem 2.3	18
2.4.4	Proof of Theorem 2.5	19
2.4.5	Proof of Theorem 2.6	20
3	Applications	21
3.1	Point process	21
3.1.1	Cox type of process with ergodic covariates	24
3.1.2	Hawkes process	28
3.2	Diffusion process	31
3.2.1	Ergodic case	31
3.2.2	Non-ergodic case	34
4	Simulations	37
4.1	Simulation for the Cox model	37
4.2	Simulations for the Hawkes process	44
4.3	Simulation for the diffusion type process	53
	Bibliography	56

<i>CONTENTS</i>	3
A. Appendix	61

Chapter 1

Introduction

In this thesis, we are interested in two topics: one is the sparse estimation and the other is statistical inference for stochastic processes. As a prologue, we shall begin to introduce the general history of these two topics briefly.

1.1 A brief review of sparse estimation

For several decades, the sparse modeling has received attention from various fields. The most commonly used sparse modeling in statistics is L^1 regularization, and typical one is the LASSO (least absolute shrinkage and selection operator), which is proposed by Tibshirani [35]. LASSO is a useful and widely studied approach to the problem of variable selection. Compared with other estimation methods, LASSO's major advantage is simultaneous execution of both parameter estimation and variable selection ([16], [35]). Originally, LASSO was introduced for linear regression problems. Suppose that $\mathbf{y} = [y_1, \dots, y_T]'$ is a response vector and $\mathbf{x}_j = [x_{1j}, \dots, x_{Tj}]'$, $j = 1, \dots, d$, are the linearly independent predictors.¹ Then the LASSO estimator is defined by

$$\hat{\theta}_{\text{LASSO}} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \left\| \mathbf{y} - \sum_{j=1}^d \mathbf{x}_j \theta_j \right\|^2 + \lambda \sum_{j=1}^d |\theta_j| \right\}, \quad (1.1)$$

where λ is a nonnegative regularization parameter. The second term on the right-hand side of (1.1) is the so-called L^1 penalty. Thanks to the singularity of the L^1 penalty at the origin, LASSO can perform automatic variable selection.

¹The prime denotes the matrix transpose.

LASSO has evolved in various directions over the last two decades, and the directions of its development are roughly categorized into the following three.

- (1) Extension of models and penalties.
- (2) Investigations to the properties of estimators.
- (3) Developments of algorithms to calculate estimators.

First, regarding (1), LASSO is applied to not only the linear regression model, but also the generalized linear regression model ([20], [33]), the graphical model ([29], [48]), the multivariate analysis ([44], [53]), the (quasi-) likelihood analysis ([14]), etc. Regarding the penalty terms, there are various extensions according to the analysis. One extension is to replace the L^1 -penalty with another penalty like the adaptive LASSO ([52]), Bridge ([18]), SCAD (smoothly clipped absolute deviation [16]), MCP (minimax concave penalty [49]), etc. Another extension is using the relation among parameters like the Group LASSO ([47]), the fused LASSO ([36]), etc.

In relation to (2), the concept of the LASSO is that “*we sacrifice a little bias to reduce the variance of the predicted values and hence may improve the overall prediction accuracy*” ([36]), but it is preferable that the estimator of interest is consistent or asymptotically normal as the sample size is large enough. Thus Knight and Fu ([26]) studied the asymptotic properties of the LASSO estimator and showed that the LASSO type estimator is consistent and asymptotically normal. Recently, the convergence rate and selection consistency of the estimators have been well studied in the case where, not only the sample size, the number of variables is also large ([5], [42], [50]). However, when the number of variables is quite large, strong conditions are required to derive the good properties of the LASSO estimator. Thus, the LASSO cannot be used in the practical situation with high dimensional data.

With regards to (3), the LARS algorithm (Least Angle Regression [15]) was the first noticed algorithm for the LASSO. The estimator obtained by the LARS algorithm is quite similar to the LASSO estimator and, moreover, those two estimators coincide by using the slightly modified version of the LARS algorithm. In recent years, fast and versatile algorithms such as the coordinate decent ([19]) or the ADMM algorithm (Alternating Direction Method and Multipliers [8]) have been proposed.

Based on the above, in this thesis, we will discuss the following.

Regarding (1), by replacing the first term on the right-hand side of (1.1) with a general loss function \mathcal{L}_T , we can easily apply the LASSO to various models. However, the asymptotic and numerical theories are established in a

case-by-case manner. One of the solutions to this problem is the least squares approximation (LSA) method proposed by Wang and Leng ([43]). The LSA is defined by a simple approximation to the original loss function:

$$\frac{1}{T}\mathcal{L}_T(\theta) \approx (\theta - \tilde{\theta})'\hat{G}(\theta - \tilde{\theta})$$

where \hat{G} is a non-singular matrix depending on the data and $\tilde{\theta}$ is the minimizer of the loss function \mathcal{L}_T . Using the LSA method, we can deal with many different models in a unified frame. Choice of the penalty term is also a crucial issue in regularization techniques. In this thesis, we adopt the Bridge and the adaptive lasso type penalty, i.e., a weighted L^q penalty because the Bridge estimator with $0 < q < 1$ and the adaptive LASSO estimator have the ‘‘oracle properties’’. Oracle properties were recognized by Fan and Li ([16]) and a good estimator with variable selection should have these properties. Let $\theta^* = [\theta_j^*]_j$ is the true value of θ and $\mathcal{A} = \{j; \theta_j^* \neq 0\}$. An estimator $\hat{\theta}$ has oracle properties if $\hat{\theta}$ satisfies

- selection consistency: $P[\hat{\theta}_{\mathcal{A}^c} = 0] \rightarrow 1$, and
- asymptotic normality: $\sqrt{T}(\hat{\theta}_{\mathcal{A}} - \theta_{\mathcal{A}}^*) \rightarrow^d N(0, G^{-1})$, for some $|\mathcal{A}| \times |\mathcal{A}|$ positive definite symmetric matrix G .

Then our objective function $Q_T^{(q)}(\theta)$ consists of an LSA term and an weighted L^q penalty term:

$$Q_T^{(q)}(\theta) = (\theta - \tilde{\theta})'\hat{G}(\theta - \tilde{\theta}) + \lambda_T \sum_{j=1}^d \hat{w}_j |\theta_j|^q. \quad (1.2)$$

Moreover, for $0 < q \leq 1$, we define a penalized least squares approximation (penalized LSA, pLSA) estimator by $\hat{\theta}^{(q)} = \operatorname{argmin}_{\theta} Q_T^{(q)}(\theta)$.

Next, respecting (2), we first show that the pLSA estimator defined above has the oracle properties. Moreover, we investigate the properties of the pLSA estimator in more detail. Specifically, we consider the L^p -boundedness of the initial estimator. This concept is often used and plays an important role in the statistical inference for stochastic processes. By assuming the L^p -boundedness of the initial estimator, we have the L^p -boundedness of the pLSA estimator and evaluate the convergence rate of selection consistency.

Finally, with respect to (3), since the pLSA objective function is non-convex when $q < 1$, the pLSA methods have a disadvantage in optimization in comparison with the L^1 regularization methods. However, by simplifying the objective function (1.2) by replacing the coefficient matrix \hat{G} with the

identity matrix, optimization of the pLSA objective function comes down to the one-dimensional optimization. Due to that simplification, the pLSA estimator loses the efficiency, but we can obtain the efficient estimator by estimating the parameter again under the model which is selected by the pLSA estimator.

1.2 A brief review of statistical inference for stochastic processes

A stochastic process is the mathematical concept which describes random phenomenon depending on time or space. Studies on the statistics of stochastic processes began in the 1970s, and it grew rapidly, especially with the development of the martingale limit theorem in the 1980s. In recent years, research has been conducted using not only martingale theory, but also the Malliavin calculus and the theory of empirical processes, and the basic theory has become wider and deeper. On the other hand, its applications to various fields such as actually finance, actuarial science, biostatistics, and survival analysis have been actively studied, and its importance has been widely recognized.

This thesis treats some specific examples of continuous-time stochastic processes, thus, we will mention this a little. The asymptotic inference for continuous-time stochastic processes has been studied based on the likelihood theory by many authors. There are many examples of this: Markov processes with the general state space (Billingsley [6]), Markov branching processes (Athreya and Keiding [2], Feigin [17]), point processes (Brown [11], Brillinger [10]) and general Levy processes (Akritas and Johnson [1]), for instance. In the 1990s, the asymptotic inference for stochastic processes began to be discussed in the general framework of semimartingales that was a wide class of stochastic processes. The LAMN property discussed by Jeganathan [24] and Basawa and Scott [4] in general framework was applied to the class of semimartingales by Luschgy [27], and he introduced the new concept of the local asymptotic quadraticity (LAQ). Yoshida [46] gave a polynomial type large deviation inequality in this LAQ setting to carry out the Ibragimov-Has'minskii-Kutoyants scheme for stochastic processes. The polynomial type large deviation inequality works in various settings, in particular, it is applicable to our examples. As mentioned in [46], L^p -boundedness is derived from the PLDI.

1.3 Organization of this thesis

This thesis is mostly based on [34] and consists of three chapters.

In Chapter 2, we first define the penalized least squares approximation estimator and derive the oracle properties of the pLSA estimator and the convergence rate of selection consistency. This is the main part of this thesis. Next, we define P-O estimator, which does not have the oracle properties but is consist. “P-O” is the abbreviation of “penalized method to ordinary method”, i.e., P-O estimator is obtained by the following two steps : (i) obtain an estimator which satisfies the selection consistency (not necessarily having the oracle properties), (ii) obtain an estimator by using the ordinary method under the model which is selected by the estimator in Step (i). Specifically, instead of using the pLSA objective function $Q_T^{(q)}(\theta)$, we use the objective function

$$Q_{T,I}^{(q)}(\theta) = |\theta - \tilde{\theta}|^2 + \lambda_T \sum_{j=1}^d \hat{w}_j |\theta_j|^q.$$

Thus, the optimization of the objective function is reduced the one-dimensional one. This gives a computational advantage. We sill prove the theorems at the end of Chapter 2.

In Chapter 3, we apply the pLSA methods to stochastic processes. We focus on two types of processes: point process and diffusion type process. For the point process, we first introduce the general theory of the point process. Next, a Cox model and a Hawkes model will be considered as the specific example. For the Cox model, we consider the intensity with following:

$$\lambda(t, \theta) = \exp \left(\sum_{j \in \mathbf{J}} \theta_j X_t^j \right),$$

where $X^j = (X_t^j)_t$ are ergodic covariate processes. For the Hawkes model, in particular, we treat the exponential Hawkes process. In this model, since the parameter space is a subset of $\mathbb{R}_+^d \times \mathbb{R}_+^{d \times d} \times \mathbb{R}_+^{d \times d}$, the true value of the parameter is on the boundary of the parameter space under the sparse situation. We will derive the selection consistency of the pLSA estimator even in this case. There are applications to various research field on the Hawkes process and recently we often consider the high dimensional situation. Therefore, it is worth considering the variable selection of the Hawkes model. For the diffusion type process, we treat the both ergodic and non-ergodic cases and also use QLA. In the ergodic case, drift parameter and volatility parameter are simultaneously estimated by QLA and we show that the variable selection is also executed simultaneously by using the pLSA methods. In regard

to the non-ergodic case, we consider only the volatility parameter. Since the QMLE (or QBE) has the asymptotic mixed normality, this is an example where the limit G of coefficient matrix \hat{G} is random matrix.

In Chapter 4, we report three simulations: (i) Cox model, (ii) Hawkes model and (iii) non-ergodic diffusion type process. For the Cox model, we take 20 Ornstein-Uhlenbeck processes as the covariates. We use the P-O estimator to avoid optimization involving the 20-dimensional parameter. For the Hawkes process, we first show the results that the initial estimator (i.e. QMLE) performs well even if the case where the true value of parameter is on the boundary of the parameter space. It will be shown that the pLSA estimator performs well. In the cases (i) and (iii), we also use the unified LASSO type estimator and the Bridge type estimator for comparison.

Chapter 2

Penalized LSA estimator

In this chapter, we will discuss the theory of penalized LSA estimation. As mentioned in the previous chapter, we can deal with various kinds of loss functions in a unified way by using the penalized least squares approximation (pLSA) methods defined in this chapter. Penalized methods in the general case has often been discussed. For example, [28] gives some results for penalized methods in M-estimation. Since optimization is usually not easy for a high dimensional parameter, it is worth considering LSA type estimation.

For the pLSA estimator, we will show the oracle property, and derive the L^p -boundedness of pLSA estimator. Moreover we will obtain the convergence rate of variable selection consistency from the L^p -boundedness of the initial estimator. We also construct an objective function in a simpler form than the conventional LSA type estimation by replacing the coefficient matrix by the identity matrix. Optimization for such objective function is easy even if the case $q < 1$.

2.1 Definition of penalized LSA estimator

Suppose that $\theta = [\theta_1, \dots, \theta_p]' \in \Theta \subset \mathbb{R}^p$ is a parameter of interest and $\tilde{\theta} = [\tilde{\theta}_1, \dots, \tilde{\theta}_p]' \in \mathbb{R}^p$ is an estimator of θ , where the parameter space Θ is an open bounded subset of \mathbb{R}^p . In many cases, $\tilde{\theta}$ minimizes some loss function $\mathcal{L}_T(\theta)$, but we will not assume the existence of the loss function. Here, T is a time index and we often consider the case where $T \in \mathbb{N}$ with discrete time observation or $T \in \mathbb{R}_+$ with continuous time observation. $\tilde{\theta}$ depends on T , however, we omit T for the sake of notational simplicity : $\tilde{\theta} = \tilde{\theta}_T$.

Example. Consider a linear regression model $y_t = \mathbf{x}'_t \theta + \epsilon_t$, ($t = 1, \dots, T, T \in \mathbb{N}$), where $\{\epsilon_t\}_t$ are independent and identically distributed random variables

with mean 0 and covariance σ^2 and $\{\mathbf{x}_t\}_t$ is independent of $\{\epsilon_t\}_t$. Then we can take $\tilde{\theta}$ as the least square estimator for $\mathcal{L}_T(\theta) = \sum_t |y_t - \mathbf{x}_t' \theta|^2$.

Example. If we consider a negative log-likelihood function as the loss function, then $\tilde{\theta}$ is the maximum likelihood estimator (MLE) of θ .

Hereafter, we assume that there exists a true value $\theta^* = [\theta_1^*, \dots, \theta_p^*]' \in \mathbb{R}^p$ of θ and that \mathbf{p}^0 components of θ^* do not equal to 0, $\mathbf{p}^0 = \#\{j; \theta_j^* \neq 0\}$. Here, for convenience of explanation, we consider a loss function $\mathcal{L}_T(\theta)$. In order to carry out parameter estimation and variable selection simultaneously, we add a penalty term to the loss function $\mathcal{L}_T(\theta)$. For example, we can take a penalized loss function as the adaptive lasso objective function by Zou (2006 [52]):

$$\frac{1}{T} \mathcal{L}_T(\theta) + \sum_{j=1}^p \kappa_T^j |\theta_j|, \quad (2.1)$$

where $\kappa_T^j = \alpha_T |\tilde{\theta}_j|^{-\gamma}$ for a deterministic sequence $(\alpha_T)_T$ and a \sqrt{T} -consistent estimator $\tilde{\theta}$.

We consider quadratic approximation of the loss function instead of the first term of (2.1). Thanks to this approximation, we can discuss the various cases into a unified methodology, and because the behavior as $T \rightarrow \infty$ is simply described, we can have a more in-depth discussion like large deviation. Moreover, we replace L^1 penalty with L^q penalty ($0 < q \leq 1$) instead. Under this setting, we will show that the parameter estimation and the variable selection can be executed simultaneously in this case. More precisely, for a $\mathbf{p} \times \mathbf{p}$ almost surely positive definite symmetric random matrix \hat{G} depending on T , we use the objective function

$$Q_T^{(q)}(\theta) = \hat{G}[(\theta - \tilde{\theta})^{\otimes 2}] + \sum_{j=1}^p \kappa_T^j |\theta_j|^q,$$

where κ_T^j are nonnegative random variables, $A^{\otimes 2} = AA'$ for a matrix or a vector A , and $A[B] = \text{Tr}(AB')$ for matrices A and B of the same size.

For twice differentiable $\mathcal{L}_T(\theta)$, $T^{-1} \mathcal{L}_T(\theta)$ is approximated as

$$\frac{1}{T} \mathcal{L}_T(\theta) \approx \frac{1}{T} \mathcal{L}_T(\tilde{\theta}) + \frac{1}{T} (\theta - \tilde{\theta})' \partial_{\theta} \mathcal{L}_T(\tilde{\theta}) + \frac{1}{2} \left\{ \frac{1}{T} \partial_{\theta}^2 \mathcal{L}_T(\tilde{\theta}) \right\} [(\theta - \tilde{\theta})^{\otimes 2}].$$

Here, the first term on the right hand side is constant with respect to θ and the second term vanishes by the definition of $\tilde{\theta}$. Thus, instead of minimizing $T^{-1} \mathcal{L}_T(\theta)$, we may minimize $\left\{ T^{-1} \partial_{\theta}^2 \mathcal{L}_T(\tilde{\theta}) \right\} [(\theta - \tilde{\theta})^{\otimes 2}]$ and in this case we can take $\hat{G} = T^{-1} \partial_{\theta}^2 \mathcal{L}_T(\tilde{\theta})$ for example.

Let $\hat{\theta}^{(q)} = [\hat{\theta}_1^{(q)}, \dots, \hat{\theta}_p^{(q)}]'$ be a minimizer of this objective function $Q_T^{(q)}(\theta)$:

$$\hat{\theta}^{(q)} \in \operatorname{argmin}_{\theta \in \bar{\Theta}} Q_T^{(q)}(\theta)$$

We call $\hat{\theta}^{(q)}$ the penalized least squares approximation (penalized LSA) estimator.

2.2 Main theorem

In this section, we will show asymptotic properties of the penalized LSA estimator $\hat{\theta}^{(q)}$ based on $Q_T^{(q)}(\theta)$. Suppose that the statistics are realized on a probability space (Ω, \mathcal{F}, P) . To describe the results, we may suppose that $\theta_1^* \neq 0, \dots, \theta_{p^0}^* \neq 0$ and $\theta_{p^0+1}^* = \dots = \theta_p^* = 0$ without loss of generality. Let

$$a_T = \max\{\kappa_T^j; j \leq p^0\} \quad \text{and} \quad b_T = \min\{\kappa_T^j; j > p^0\}.$$

For a vector $v = [v_1, \dots, v_p]' \in \mathbb{R}^p$, we denote subvectors $[v_1, \dots, v_{p^0}]'$ and $[v_{p^0+1}, \dots, v_p]'$ by $v_{\mathcal{J}^1}$ and $v_{\mathcal{J}^0}$ respectively.

We consider the following conditions with respect to $\tilde{\theta}$ and \hat{G} . Let r_T be a sequence of positive numbers tending to 0 as $T \rightarrow \infty$. We often consider the case that $r_T = T^{-1/2}$.

Assumption 1. There exists a positive definite symmetric random matrix G such that $\hat{G} \rightarrow^p G$.

Assumption 2. $\tilde{\theta}$ is r_T^{-1} -consistent, i.e., $r_T^{-1}(\tilde{\theta} - \theta^*) = O_p(1)$.

Assumption 3. $r_T^{-1}(\tilde{\theta} - \theta^*) \rightarrow^{d_s} \Gamma^{-\frac{1}{2}}\zeta$ holds, where Γ is a $p \times p$ positive definite random symmetric matrix, ζ is a p -dimensional standard Gaussian random vector defined on an extended probability space of (Ω, \mathcal{F}, P) and independent of \mathcal{G} , and d_s denotes the \mathcal{G} -stable convergence for some σ -field \mathcal{G} such that $\sigma(\Gamma) \subset \mathcal{G} \subset \mathcal{F}$.

Of course, Assumption 3 is stronger than Assumption 2, but r_T^{-1} -consistency and selection consistency of the penalized LSA estimator $\hat{\theta}^{(q)}$ are derived from Assumptions 1 and 2. We need Assumption 3 to show asymptotic normality of penalized LSA estimator $\hat{\theta}^{(q)}$.

For a $p \times p$ matrix $M = [m_{ij}]_{1 \leq i \leq p, 1 \leq j \leq p}$, we denote the $p^0 \times p^0$ matrix $[m_{ij}]_{1 \leq i \leq p^0, 1 \leq j \leq p^0}$, $p^0 \times (p - p^0)$ matrix $[m_{ij}]_{1 \leq i \leq p^0, p^0 < j \leq p}$, $(p - p^0) \times p^0$ matrix $[m_{ij}]_{p^0 < i \leq p, 1 \leq j \leq p^0}$ and $(p - p^0) \times (p - p^0)$ matrix $[m_{ij}]_{p^0 < i \leq p, p^0 < j \leq p}$ by $M_{\mathcal{J}^{11}}, M_{\mathcal{J}^{10}}, M_{\mathcal{J}^{01}}$ and $M_{\mathcal{J}^{00}}$ respectively:

$$M = \begin{bmatrix} M_{\mathcal{J}^{11}} & M_{\mathcal{J}^{10}} \\ M_{\mathcal{J}^{01}} & M_{\mathcal{J}^{00}} \end{bmatrix}.$$

Theorem 2.1 (r_T^{-1} -consistency). Under Assumptions 1 and 2, if $r_T^{-1}a_T = O_p(1)$, then

$$r_T^{-1}(\hat{\theta}^{(q)} - \theta^*) = O_p(1).$$

Let $\hat{\mathcal{J}}^1 = \{j = 1, \dots, \mathbf{p}; \hat{\theta}_j^{(q)} \neq 0\}$.

Theorem 2.2 (Selection consistency). Under Assumptions 1 and 2, if $r_T^{-1}a_T = O_p(1)$ and $r_T^{-(2-q)}b_T \rightarrow^p \infty$, then

$$P[\hat{\mathcal{J}}^1 = \{1, \dots, \mathbf{p}^0\}] \rightarrow 1. \quad (2.2)$$

Theorem 2.3 (Asymptotic normality). Let $\mathfrak{G} = [I_{\mathbf{p}^0} \quad (G_{\mathcal{J}^{11}})^{-1}G_{\mathcal{J}^{10}}]$ for $\mathbf{p}^0 \times \mathbf{p}^0$ identity matrix $I_{\mathbf{p}^0}$. Under Assumptions 1 and 2, if $r_T^{-1}a_T = o_p(1)$ and $r_T^{-(2-q)}b_T \rightarrow^p \infty$, then

$$r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)_{\mathcal{J}^1} - \mathfrak{G}\{r_T^{-1}(\tilde{\theta} - \theta^*)\} \rightarrow^p 0.$$

In particular, under Assumption 3 and $G = \Gamma$, we have

$$r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)_{\mathcal{J}^1} \rightarrow^{d_s} \mathfrak{G}\Gamma^{-\frac{1}{2}}\zeta \sim \text{MN}_{\mathbf{p}^0}(0, (\Gamma_{\mathcal{J}^{11}})^{-1}).$$

Hereafter, we consider $\kappa_T^j = \alpha_T |\tilde{\theta}_j|^{-\gamma}$, where γ is a constant satisfying $\gamma > -(1 - q)$ and $(\alpha_T)_T$ is a deterministic sequence. If $(\alpha_T)_T$ satisfies the conditions

$$r_T^{-(2-q+\gamma)}\alpha_T \rightarrow \infty \quad \text{and} \quad r_T^{-1}\alpha_T = o(1), \quad (2.3)$$

then the conditions in Theorems 2.1-2.3 are fulfilled. Moreover we will show that the probability $P[\hat{\theta}_{\mathcal{J}^0}^{(q)} = 0]$ can be evaluated by any power of r_T .

Let $\tilde{u} = r_T^{-1}(\tilde{\theta} - \theta^*)$ and $\hat{u} = r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)$.

Definition 2.4. For a stochastic process $X = \{X_T\}_T$ is $L^{\infty-}$ -bounded if and only if $\sup_T E[|X_T|^p] < \infty$ holds for all $p \geq 1$.

Additionally, we consider the following conditions:

Assumption 4. $\{\hat{G}\}_T$, $\{\hat{G}^{-1}\}_T$ and $\{\tilde{u}\}_T$ are $L^{\infty-}$ -bounded.

Remark. The L^p -boundedness of a sequence of estimators can be obtained by the quasi-likelihood analysis with a polynomial type large deviation inequality for an associated statistical random field. See [46] for details.

Theorem 2.5. Let $\epsilon \in (-1 + q, \gamma)$. Assume that $r_T^{1+\gamma-\epsilon}\alpha_T^{-1} = O(1)$ and $r_T^{-1}\alpha_T = O(1)$. Then under Assumptions 1 and 4, $\{\hat{u}\}_T$ is $L^{\infty-}$ -bounded. Moreover, for all $L > 0$, there exists a constant C_L such that

$$P[\hat{\mathcal{J}}^1 = \{1, \dots, \mathbf{p}^0\}] \geq 1 - C_L r_T^{2L}. \quad (2.4)$$

for all $T > 0$.

2.3 P-O estimator

We now discuss the coefficient matrix \hat{G} . In the above theorems, we assume convergence of \hat{G} to G or L^∞ -boundedness of $\{\hat{G}\}$ and $\{\hat{G}^{-1}\}$ but we should not necessarily find such coefficient matrix \hat{G} . In fact, if we take $\hat{G} = I_p$, then we can apply Theorems 2.1-2.5 except that the conditional asymptotic variance in Theorem 3 becomes $(\Gamma^{-1})_{\mathcal{J}^{11}}$. Since $(\Gamma_{\mathcal{J}^{11}})^{-1} = (\Gamma^{-1})_{\mathcal{J}^{11}} - (\Gamma^{-1})_{\mathcal{J}^{10}}((\Gamma^{-1})_{\mathcal{J}^{00}})^{-1}(\Gamma^{-1})_{\mathcal{J}^{01}}$, this estimator is not efficient. However, the objective function has the following simple form

$$Q_T^{(q)}(\theta) = \sum_{j=1}^p \left((\theta_j - \tilde{\theta}_j)^2 + \kappa_T^j |\theta_j|^q \right).$$

From a computational point of view, this fact is useful because it is difficult to optimize the non-convex function in the high-dimensional case. Then we obtain a new estimator under the model selected by the penalized LSA estimator with coefficient matrix I_p . We call this estimator the P-O (penalized method to ordinary method) estimator and denote it by $\check{\theta}$. More precisely, we define the P-O estimator as follows.

Let Θ be a bounded open subset of \mathbb{R}^p . First, we assume the r_T^{-1} -consistency of the initial estimator $\tilde{\theta}$. Second, we get the penalized LSA estimator $\hat{\theta}_{I_p}^{(q)}$ with coefficient matrix I_p defined by

$$\hat{\theta}_{I_p}^{(q)} \in \operatorname{argmin}_{\theta \in \tilde{\Theta}} Q_T^{(q)}(\theta),$$

where $\kappa_T^j = \alpha_T |\tilde{\theta}_j|^{-\gamma}$, $\gamma > -(1 - q)$ and α_T satisfies (2.3). Let $\hat{\mathcal{J}}_p^0 = \{j = 1, \dots, p; \hat{\theta}_{I_p, j}^{(q)} = 0\}$ and $\tilde{\Theta} = \{\theta \in \Theta; \theta_j = 0, j \in \hat{\mathcal{J}}_p^0\}$. Here, we consider another loss function $\mathbb{L}_T(\theta)$. Then, we define the P-O estimator $\check{\theta}$ by

$$\check{\theta} \in \operatorname{argmin}_{\theta \in \tilde{\Theta}} \mathbb{L}_T(\theta).$$

Before we turn to the statement of the results for the P-O estimator $\check{\theta}$, we consider some conditions. We denote a parameter $\theta = \begin{bmatrix} \phi \\ \psi \end{bmatrix} \in \mathbb{R}^{p^0 + (p - p^0)}$ and its true value $\theta^* = \begin{bmatrix} \phi^* \\ \psi^* \end{bmatrix} = \begin{bmatrix} \phi^* \\ 0 \end{bmatrix}$. Let $\bar{\mathbb{L}}_T(\phi) = \mathbb{L}_T\left(\begin{bmatrix} \phi \\ 0 \end{bmatrix}\right)$ and

$$\bar{\phi} \in \operatorname{argmin}_{\phi} \bar{\mathbb{L}}_T(\phi). \quad (2.5)$$

Assumption 5. For any sequence of estimations $\{\bar{\phi}\}$ satisfying (2.5),

- (i) $\{\tilde{u}\}_T = \{r_T^{-1}(\tilde{\theta} - \theta^*)\}_T$ is $L^{\infty-}$ -bounded.
- (ii) $r_T^{-1}(\bar{\phi} - \phi^*) \rightarrow^{d_s} \Lambda^{-\frac{1}{2}}\eta$, where Λ is a $\mathbf{p}^0 \times \mathbf{p}^0$ positive definite symmetric random matrix, η is a \mathbf{p}^0 -dimensional standard Gaussian random vector independent of Λ .
- (iii) $\{r_T^{-1}(\bar{\phi} - \phi^*)\}_T$ is $L^{\infty-}$ -bounded.

Remark. In many cases, we take $\mathbb{L}_T(\theta) = \mathcal{L}_T(\theta)$ and $\Lambda = \Gamma_{\mathcal{J}^{11}}$. Then, the sufficient condition for Assumption 5 is as follows. We define the random field $\mathbb{Z}_T : \mathbb{U}_T \rightarrow \mathbb{R}_+$ by $\mathbb{Z}_T(u) = \exp\{-\mathcal{L}_T(\theta^* + r_T u) + \mathcal{L}_T(\theta^*)\}$, where $\mathbb{U}_T = \{u \in \mathbb{R}^{\mathbf{p}}; \theta^* + r_T u \in \Theta\}$. We denote $B(R) = \{u \in \mathbb{R}^{\mathbf{p}}; |u| \leq R\}$. If $\mathbb{Z}_T(u) \rightarrow^{d_s} \mathbb{Z}(u)$ in $C(B(R))$ for every $R > 0$ as $T \rightarrow \infty$ and $\{\bar{\phi}\}$ is tight, then Assumption 5 (ii) holds. Here, \mathbb{Z} is a random field defined by $\mathbb{Z}(u) = \exp\left(u' \Gamma^{\frac{1}{2}} \zeta - \frac{1}{2} u' \Gamma u\right)$. Moreover, if the random field \mathbb{Z}_T satisfies the polynomial type large deviation inequality (Theorem 1 in [46]), then by using Proposition 1 in [46], Assumptions 5(i) and (iii) hold.

Theorem 2.6. (a) Under Assumption 2,

$$P\left[\check{\theta}_{\mathcal{J}^1} \in \operatorname{argmin}_{\phi} \bar{\mathbb{L}}_T(\phi)\right] \rightarrow 1. \quad (2.6)$$

Additionally, under Assumption 5(ii),

$$r_T^{-1}(\check{\theta} - \theta^*)_{\mathcal{J}^1} \rightarrow^{d_s} \Lambda^{-\frac{1}{2}}\eta \sim \operatorname{MN}_{\mathbf{p}^0}(0, \Lambda^{-1}).$$

- (b) Let $\epsilon \in (-1 + q, \gamma)$. Assume Assumptions 5(i) and (iii), $r_T^{1+\gamma-\epsilon} \alpha_T^{-1} = O(1)$ and $r_T^{-1} \alpha_T = O(1)$. Then $\{r_T^{-1}(\check{\theta} - \theta^*)\}_T$ is $L^{\infty-}$ -bounded. Moreover, for all $L > 0$ there exists a constant C_L such that

$$P[\check{\mathcal{J}}^1 = \{1, \dots, \mathbf{p}^0\}] \geq 1 - C_L r_T^{2L}. \quad (2.7)$$

for all $T > 0$ where $\check{\mathcal{J}}^1 = \{j = 1, \dots, \mathbf{p}; \check{\theta}_j^{(q)} \neq 0\}$.

2.4 Proofs of main theorems

2.4.1 Proof of Theorem 2.1

Since $\hat{\theta}^{(q)}$ minimizes $Q_T^{(q)}(\theta)$, we obtain

$$\begin{aligned}
0 &\geq Q_T^{(q)}(\hat{\theta}^{(q)}) - Q_T^{(q)}(\theta^*) \\
&= \hat{G}[(\hat{\theta}^{(q)} - \tilde{\theta})^{\otimes 2}] + \sum_{j=1}^p \kappa_T^j |\hat{\theta}_j^{(q)}|^q - \hat{G}[(\theta^* - \tilde{\theta})^{\otimes 2}] - \sum_{j=1}^p \kappa_T^j |\theta_j^*|^q \\
&= \hat{G}[(\hat{\theta}^{(q)} - \theta^*)^{\otimes 2}] + 2(\hat{\theta}^{(q)} - \theta^*)' \hat{G}(\theta^* - \tilde{\theta}) + \sum_{j=1}^p \kappa_T^j |\hat{\theta}_j^{(q)}|^q - \sum_{j=1}^p \kappa_T^j |\theta_j^*|^q.
\end{aligned} \tag{2.8}$$

Since $0 \leq |\hat{\theta}_j^{(q)}| < |\theta_j^*|$ implies $(|\theta_j^*|^q - |\hat{\theta}_j^{(q)}|^q)/(|\theta_j^*| - |\hat{\theta}_j^{(q)}|) \leq |\theta_j^*|^q/|\theta_j^*| = |\theta_j^*|^{q-1}$, we obtain $|\hat{\theta}_j^{(q)}|^q - |\theta_j^*|^q \geq -K^*|\hat{\theta}_j^{(q)} - \theta_j^*|$ where $K^* = \max_{1 \leq j \leq p^0} |\theta_j^*|^{q-1}$. Thus

$$\begin{aligned}
\sum_{j=1}^p \kappa_T^j |\hat{\theta}_j^{(q)}|^q - \sum_{j=1}^p \kappa_T^j |\theta_j^*|^q &\geq \sum_{j=1}^{p^0} \kappa_T^j (|\hat{\theta}_j^{(q)}|^q - |\theta_j^*|^q) \\
&\geq - \sum_{j=1}^{p^0} K^* \kappa_T^j |\hat{\theta}_j^{(q)} - \theta_j^*| \\
&\geq -p^0 K^* a_T |\hat{\theta}^{(q)} - \theta^*|.
\end{aligned}$$

Therefore, by multiplying both sides of (2.8) by r_T^{-2} , we obtain

$$\begin{aligned}
0 &\geq \hat{G}[\{r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)\}^{\otimes 2}] + 2\{r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)\}' \hat{G}\{r_T^{-1}(\theta^* - \tilde{\theta})\} - \\
&\quad p^0 K^* r_T^{-1} a_T |r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)| \\
&\geq \tau_{\min}(\hat{G}) |r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)|^2 - 2\tau_{\max}(\hat{G}) |r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)| |r_T^{-1}(\tilde{\theta} - \theta^*)| - \\
&\quad p^0 K^* r_T^{-1} a_T |r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)|,
\end{aligned}$$

where $\tau_{\min}(\hat{G})$ and $\tau_{\max}(\hat{G})$ are the minimum and maximum eigenvalues of matrix \hat{G} , respectively. After all,

$$|r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)| \leq \left\{ \frac{1}{\tau_{\min}(\hat{G})} \left(2\tau_{\max}(\hat{G}) |r_T^{-1}(\tilde{\theta} - \theta^*)| + p^0 K^* r_T^{-1} a_T \right) \right\}. \tag{2.9}$$

Since the right hand side is $O_p(1)$ by the assumption, we obtain $r_T^{-1}(\hat{\theta}^{(q)} - \theta^*) = O_p(1)$. \square

2.4.2 Proof of Theorem 2.2

First, we assume that $\hat{\theta}^{(a)} \in \partial\Theta$. Then, since $|\hat{\theta}^{(a)} - \theta^*| \geq \epsilon_0$ where $\epsilon_0 = \inf\{|\theta - \theta^*|; \theta \in \partial\Theta\} > 0$, we obtain

$$P[\hat{\theta}^{(a)} \in \partial\Theta] \leq P[r_T^{-1}|\hat{\theta}^{(a)} - \theta^*| \geq r_T^{-1}\epsilon_0] \rightarrow 0. \quad (2.10)$$

Next, we assume $\hat{\theta}^{(a)} \notin \partial\Theta$ and $\hat{\theta}_j^{(a)} \neq 0$ for some $j(\mathbf{p}^0 < j \leq \mathbf{p})$. Since $Q_T^{(a)}(\theta)$ is differentiable at $\theta = \hat{\theta}^{(a)}$ with respect to the j -th component and $\hat{\theta}^{(a)}$ minimizes $Q_T^{(a)}(\theta)$,

$$\begin{aligned} 0 &= r_T^{-1} \frac{\partial Q_T^{(a)}(\theta)}{\partial \theta_j} \Big|_{\theta=\hat{\theta}^{(a)}} \\ &= 2\hat{G}^{(j)} \{r_T^{-1}(\hat{\theta}^{(a)} - \tilde{\theta})\} + r_T^{-1} \kappa_T^j q |\hat{\theta}_j^{(a)}|^{q-1} \text{sgn}(\hat{\theta}_j^{(a)}), \end{aligned}$$

where $\hat{G}^{(j)}$ means the j -th row vector of \hat{G} . Therefore, we have

$$\begin{aligned} 2|\hat{G}^{(j)} \{r_T^{-1}(\hat{\theta}^{(a)} - \tilde{\theta})\}| |r_T^{-1} \hat{\theta}_j^{(a)}|^{1-q} &= q r_T^{-(2-q)} \kappa_T^j \\ &\geq q r_T^{-(2-q)} b_T. \end{aligned} \quad (2.11)$$

Since, by Theorem 2.1 and the assumption, the left hand side of above equation is $O_p(1)$ and $r_T^{-(2-q)} b_T \rightarrow^p \infty$, we obtain

$$\begin{aligned} &P[\hat{\theta}_j^{(a)} \neq 0, \hat{\theta}^{(a)} \notin \partial\Theta] \\ &\leq P[2|\hat{G}^{(j)} \{r_T^{-1}(\hat{\theta}^{(a)} - \tilde{\theta})\}| |r_T^{-1} \hat{\theta}_j^{(a)}|^{1-q} \geq q r_T^{-(2-q)} b_T] \rightarrow 0 \end{aligned} \quad (2.12)$$

for $j = \mathbf{p}^0 + 1, \dots, \mathbf{p}$.

Thus we have

$$P[\hat{\theta}_{\mathcal{J}^0}^{(a)} \neq 0] \leq P[\hat{\theta}^{(a)} \in \partial\Theta] + \sum_{j=\mathbf{p}^0+1, \dots, \mathbf{p}} P[\hat{\theta}_j^{(a)} \neq 0, \hat{\theta}^{(a)} \notin \partial\Theta] \rightarrow 0.$$

In particular, since $\hat{\mathcal{J}}^1 \neq \{1, \dots, \mathbf{p}^0\}$ implies

$$|(\hat{\theta}^{(a)} - \theta^*)_{\mathcal{J}^1}| \geq \min_{j=1, \dots, \mathbf{p}^0} |\theta_j^*| > 0$$

or $\hat{\theta}_{\mathcal{J}^0}^{(a)} \neq 0$, we have (2.2) by Theorem 2.1. \square

2.4.3 Proof of Theorem 2.3

For $\theta = \begin{bmatrix} \theta_{\mathcal{J}^1} \\ \theta_{\mathcal{J}^0} \end{bmatrix} \in \mathbb{R}^p$,

$$\begin{aligned} Q_T^{(q)}(\theta) &= \hat{G}[(\theta - \tilde{\theta})^{\otimes 2}] + \sum_{j=1}^p \kappa_T^j |\theta_j|^q \\ &= \hat{G}_{\mathcal{J}^{11}}[(\theta - \tilde{\theta})_{\mathcal{J}^1}^{\otimes 2}] + 2(\theta - \tilde{\theta})'_{\mathcal{J}^1} \hat{G}_{\mathcal{J}^{10}}(\theta - \tilde{\theta})_{\mathcal{J}^0} + \hat{G}_{\mathcal{J}^{00}}[(\theta - \tilde{\theta})_{\mathcal{J}^0}^{\otimes 2}] \\ &\quad + \sum_{j=1}^{p^0} \kappa_T^j |\theta_j|^q + \sum_{j=p^0+1}^p \kappa_T^j |\theta_j|^q. \end{aligned}$$

In particular, for $\theta^\ddagger = \begin{bmatrix} \theta_{\mathcal{J}^1} \\ 0 \end{bmatrix} \in \mathbb{R}^p$,

$$Q_T^{(q)}(\theta^\ddagger) = \hat{G}_{\mathcal{J}^{11}}[(\theta - \tilde{\theta})_{\mathcal{J}^1}^{\otimes 2}] - 2(\theta - \tilde{\theta})'_{\mathcal{J}^1} \hat{G}_{\mathcal{J}^{10}} \tilde{\theta}_{\mathcal{J}^0} + \hat{G}_{\mathcal{J}^{00}}[\tilde{\theta}_{\mathcal{J}^0}^{\otimes 2}] + \sum_{j=1}^{p^0} \kappa_T^j |\theta_j|^q.$$

Let

$$A_T = \left\{ \min_{1 \leq j \leq p^0} |\hat{\theta}_j^{(q)}| > 0, \hat{\theta}_{\mathcal{J}^0}^{(q)} = 0, \det(\hat{G}_{\mathcal{J}^{11}}) \neq 0 \right\}.$$

Then Theorems 2.1 and 2.2 imply $P[A_T] \rightarrow 1$. Let $\mathbb{R}_0^p = \{\theta \in \mathbb{R}^p; \theta_{\mathcal{J}^0} = 0\}$.

Since $Q_T^{(q)}(\hat{\theta}^{(q)}) = \min_{\theta^\ddagger \in \mathbb{R}_0^p} Q_T^{(q)}(\theta^\ddagger)$ on A_T ,

$$\begin{aligned} 0 &= \frac{1}{2} \frac{\partial Q_T^{(q)}(\theta)}{\partial \theta_{\mathcal{J}^1}} \Big|_{\theta = \hat{\theta}^{(q)}} \\ &= \hat{G}_{\mathcal{J}^{11}}(\hat{\theta}^{(q)} - \tilde{\theta})_{\mathcal{J}^1} - \hat{G}_{\mathcal{J}^{10}} \tilde{\theta}_{\mathcal{J}^0} + V(\hat{\theta}_{\mathcal{J}^1}^{(q)}) \end{aligned}$$

holds on A_T , where $V(\hat{\theta}_{\mathcal{J}^1}^{(q)}) = [2^{-1} q \kappa_T^j |\hat{\theta}_j^{(q)}|^{q-1} \text{sgn}(\hat{\theta}_j^{(q)})]_{j=1, \dots, p^0} \in \mathbb{R}^{p^0}$. Let $\hat{\mathfrak{G}} = [I_{p^0} \quad (\hat{G}_{\mathcal{J}^{11}})^{-1} \hat{G}_{\mathcal{J}^{10}}]$. Since $\hat{\mathfrak{G}} \rightarrow^p \mathfrak{G}$ and $1_{A_T} \{r_T^{-1} (\hat{G}_{\mathcal{J}^{11}})^{-1} V(\hat{\theta}_{\mathcal{J}^1}^{(q)})\} \rightarrow^p$

0, we have

$$\begin{aligned}
& r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)_{\mathcal{J}^1} - \mathfrak{G}\{r_T^{-1}(\tilde{\theta} - \theta^*)\} \\
&= 1_{A_T} \left\{ r_T^{-1}(\tilde{\theta} - \theta^*)_{\mathcal{J}^1} + r_T^{-1}(\hat{G}_{\mathcal{J}^{11}})^{-1} \hat{G}_{\mathcal{J}^{10}} \tilde{\theta}_{\mathcal{J}^0} \right. \\
&\quad \left. - r_T^{-1}(\hat{G}_{\mathcal{J}^{11}})^{-1} V(\hat{\theta}_{\mathcal{J}^1}^{(q)}) - \mathfrak{G}\{r_T^{-1}(\tilde{\theta} - \theta^*)\} \right\} \\
&\quad + 1_{A_T^c} \left\{ r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)_{\mathcal{J}^1} - \mathfrak{G}\{r_T^{-1}(\tilde{\theta} - \theta^*)\} \right\} \\
&= 1_{A_T} \left\{ (\mathfrak{G} - \mathfrak{G})\{r_T^{-1}(\tilde{\theta} - \theta^*)\} - r_T^{-1}(\hat{G}_{\mathcal{J}^{11}})^{-1} V(\hat{\theta}_{\mathcal{J}^1}^{(q)}) \right\} \\
&\quad + 1_{A_T^c} \left\{ r_T^{-1}(\hat{\theta}^{(q)} - \theta^*)_{\mathcal{J}^1} - \mathfrak{G}\{r_T^{-1}(\tilde{\theta} - \theta^*)\} \right\} \\
&\rightarrow^p 0.
\end{aligned}$$

□

2.4.4 Proof of Theorem 2.5

By (2.9) and Assumption 4, $\{\hat{u}\}$ is $L^{\infty-}$ -bounded. By (2.10) in the Proof of theorem 2.2, we obtain

$$\begin{aligned}
P[\hat{\theta}^{(q)} \in \partial\Theta] &\leq P[r_T^{-1}|\hat{\theta}^{(q)} - \theta^*| \geq r_T^{-1}\epsilon_0] \\
&\leq \frac{1}{(r_T^{-1}\epsilon_0)^L} E[|\hat{u}|^L],
\end{aligned}$$

for all $L \geq 1$. For $j > \mathfrak{p}^0$, by the equation (2.11) and the Markov's inequality, we have

$$\begin{aligned}
& P\left[\hat{\theta}_j^{(q)} \neq 0, \hat{\theta}^{(q)} \notin \partial\Theta\right] \\
&\leq P\left[2|\hat{G}^{(j)}| \cdot |\hat{u} - \tilde{u}| \geq r_T^{-1}\kappa_T^j q |r_T \hat{u}|^{-(1-q)}\right] \\
&\leq \frac{1}{r_T^{-(1-q+\epsilon)M}} 2^M q^{-M} E\left[|\hat{G}^{(j)}|^M |\hat{u} - \tilde{u}|^M |\hat{u}|^{M(1-q)} \left(\frac{1}{r_T^{\epsilon-1}\kappa_T^j}\right)^M\right],
\end{aligned}$$

where $M = M(L) = 2L(1 - q + \epsilon)^{-1}$ and $L > 0$ is an arbitrary constant. Here, by Hölder inequality, we have

$$\begin{aligned}
& E\left[|\hat{G}^{(j)}|^M |\hat{u} - \tilde{u}|^M |\hat{u}|^{M(1-q)} \left(\frac{1}{r_T^{\epsilon-1}\kappa_T^j}\right)^M\right] \\
&\leq E\left[|\hat{G}^{(j)}|^{4M}\right]^{\frac{1}{4}} E\left[|\hat{u} - \tilde{u}|^{4M}\right]^{\frac{1}{4}} E\left[|\hat{u}|^{4M(1-q)}\right]^{\frac{1}{4}} E\left[\left(\frac{1}{r_T^{\epsilon-1}\kappa_T^j}\right)^{4M}\right]^{\frac{1}{4}}. \quad (2.13)
\end{aligned}$$

Since

$$\begin{aligned} E \left[\left(\frac{1}{r_T^{\epsilon-1} \kappa_T^j} \right)^{4M} \right] &= E \left[\left(\frac{|\tilde{\theta}_j|^\gamma}{r_T^{\epsilon-1} \alpha_T} \right)^{4M} \right] \\ &\leq \left(\frac{1}{r_T^{-(1+\gamma-\epsilon)} \alpha_T} \right)^{4M} E \left[|\tilde{u}|^{4\gamma M} \right] \end{aligned}$$

and $\{\hat{u}\}_T, \{\tilde{u}\}_T$ and $\{\hat{G}\}$ are $L^{\infty-}$ -bounded, the right-hand side of (2.13) is bounded uniformly in T . Finally, the inequality (2.4) is obtained in a similar way as the Proof of Theorem 2.2. \square

2.4.5 Proof of Theorem 2.6

(a) Since $\hat{\mathcal{J}}_{I_p}^0 = \{\mathbf{p}^0 + 1, \dots, \mathbf{p}\}$ implies $\check{\theta}_{\mathcal{J}^1} \in \arg\min_{\phi} \bar{\mathbb{L}}_T(\phi)$, we obtain (2.6) by Theorem 2.2.

(b) Next, let $B_T = \{\hat{\mathcal{J}}_{I_p}^0 = \{\mathbf{p}^0 + 1, \dots, \mathbf{p}\}\}$ and $\text{diam}(\Theta) = \sup\{|\theta_1 - \theta_2|; \theta_1, \theta_2 \in \Theta\}$. By Theorem 2.5, $P[B_T^c]$ is evaluated by any power of r_T . Since

$$\begin{aligned} &\sup_T E[|r_T^{-1}(\check{\theta} - \theta^*)|^p] \\ &\leq \sup_T E[|r_T^{-1}(\bar{\phi} - \phi^*)|^p] + \sup_T \left\{ P[B_T^c] \cdot (r_T^{-1} \text{diam}(\Theta))^p \right\} < \infty \end{aligned}$$

for $\bar{\phi} = \check{\theta}_{\mathcal{J}^1}$ and all $p > 0$, we have $L^{\infty-}$ -boundedness of $\{r_T^{-1}(\check{\theta} - \theta^*)\}_T$. By the definition of $\check{\theta}$, $\hat{\mathcal{J}}_{I_p}^0 = \{\mathbf{p}^0 + 1, \dots, \mathbf{p}\}$ and $|\check{\theta} - \theta^*| < \min_{1 \leq j \leq \mathbf{p}^0} |\theta_j^*|$ imply $\check{\mathcal{J}}^1 = \{1, \dots, \mathbf{p}^0\}$. Thus

$$P[\check{\mathcal{J}}^1 \neq \{1, \dots, \mathbf{p}^0\}] \leq P[B_T^c] + P\left[|r_T^{-1}(\check{\theta} - \theta^*)| \geq r_T^{-1} \min_{1 \leq j \leq \mathbf{p}^0} |\theta_j^*|\right]$$

Therefore, by $L^{\infty-}$ -boundedness of $\{r_T^{-1}(\check{\theta} - \theta^*)\}_T$, we obtain the inequality (2.7). \square

Chapter 3

Applications

In the previous chapter, we discussed the theory of pLSA estimator. In this chapter, we consider the its application to analysis of stochastic processes. In particular, we are interested in the point processes and the diffusion processes. For the point process, first, we consider the general theory of ergodic intensity model. Then, as an example, we treat the Cox process and Hawkes process using the quasi-likelihood analysis (QLA) method. For the diffusion process, we consider the ergodic and non-ergodic diffusion processes. We also use QLA method in this case.

3.1 Point process

In this section, we will apply the results in Section 3 to a point process with parameters containing zero components. We consider a multivariate point process $N = (N_t^\alpha)_{\alpha \in \mathbf{I}, t \in \mathbb{R}_+}$ with intensity process $\lambda(t, \theta) = (\lambda^\alpha(t, \theta))_{\alpha \in \mathbf{I}}$, $t \in \mathbb{R}_+$, where $\mathbf{I} = \{1, 2, \dots, \mathbf{d}\}$ is an index set. More precisely, given a stochastic basis $\mathcal{B} = (\Omega, \mathcal{F}, \mathbf{F}, P)$ with a filtration $\mathbf{F} = (\mathcal{F}_t)_{t \in \mathbb{R}_+}$, we suppose that N and $\lambda(\cdot, \theta)$ are defined on \mathcal{B} , the simple counting process N is \mathbf{F} -adapted right-continuous, $\lambda(\cdot, \theta)$ is predictable locally integrable for every $\theta \in \Theta$, and that $N - \int_0^\cdot \lambda(s, \theta^*) ds$ is a \mathbf{d} -dimensional local martingale with respect to \mathbf{F} . Assume that the components of N have no common jumps. The parameter space Θ is a bounded open set in $\mathbb{R}^{\mathbf{p}}$ that admits Sobolev's inequality

$$\|f\|_\infty \leq C_\Theta \sum_{i=0,1} \|\partial_\theta^i f\|_{L^r(\Theta)}$$

for elements f of the Sobolev space $f \in W^{1,r}(\Theta)$, with a constant C_Θ independent of f , for $r > \mathbf{p}$. We suppose that $0 \in \mathbb{R}^{\mathbf{p}}$ is in Θ and that the mapping $\theta \mapsto \lambda(t, \theta)$ is continuously extended to $\bar{\Theta}$.

We will use the quasi-likelihood method ([12]) with the quasi-log likelihood function

$$\ell_T(\theta) = \sum_{\alpha \in \mathbf{I}} \int_0^T \log(\lambda^\alpha(t, \theta)) dN_t^\alpha - \sum_{\alpha \in \mathbf{I}} \int_0^T \lambda^\alpha(t, \theta) dt. \quad (3.1)$$

Then $\mathcal{L}_T(\theta) = -\ell_T(\theta)$ becomes a loss function. The conditions stated later ensure the existence of the function (3.1). For the initial estimator $\tilde{\theta}$, we can use, for example, the quasi-maximum likelihood estimator $\tilde{\theta}^M$ and the quasi-Bayesian estimator $\tilde{\theta}^B$ given by

$$\tilde{\theta}^M \in \operatorname{argmax}_{\theta \in \tilde{\Theta}} \ell_T(\theta)$$

and

$$\tilde{\theta}^B = \left[\int_{\Theta} \exp(\ell_T(\theta)) \pi(\theta) d\theta \right]^{-1} \int_{\Theta} \theta \exp(\ell_T(\theta)) \pi(\theta) d\theta,$$

respectively, where π is a prior density satisfying $0 < \inf_{\theta} \pi(\theta) \leq \sup_{\theta} \pi(\theta) < \infty$.

For ergodic point processes, asymptotic normality and convergence of moments of $\tilde{\theta}^M$ and $\tilde{\theta}^B$ were proved in [12]. We recall their results briefly. Hereafter $\theta^* \in \Theta$ denotes the true value of θ and the distribution of the data is expressed by a multivariate point process N with intensity process $\lambda(t, \theta^*)$. For a random variable X , we denote $\|X\|_p = E[|X|^p]^{\frac{1}{p}}$. We write $C_{\uparrow}(\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p)$ the set of functions f satisfying the following conditions : (i) f is continuous on $(\mathbb{R}_+ - \{0\}) \times (\mathbb{R}_+ - \{0\}) \times \mathbb{R}^p$, (ii) for any $(u, v, w) \in \mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p$, $f(0, v, w) = f(u, 0, w) = 0$ and (iii) f is of polynomial growth in $(u, v, w, \frac{1_{\{u>0\}}}{u}, \frac{1_{\{v>0\}}}{v})$, i.e. there exists a constant $C_0 > 0$ and $m_1, m_2, m_3, m_4, m_5 \in \mathbb{N}$ such that $|f(u, v, w)| \leq C_0(1 + u^{m_1} + v^{m_2} + |w|^{m_3} + (\frac{1_{\{u>0\}}}{u})^{m_4} + (\frac{1_{\{v>0\}}}{v})^{m_5})$ holds for all $(u, v, w) \in \mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p$.

Assumption 6. The mapping $\lambda : \Omega \times \mathbb{R}_+ \times \Theta \rightarrow \mathbb{R}_+^d$ is $\mathcal{F} \times \mathbb{B}(\mathbb{R}_+) \times \mathbb{B}(\Theta)$ -measurable and almost surely satisfies

- (i) for every $\theta \in \Theta$, the mapping $s \mapsto \lambda(s, \theta)$ is left continuous,
- (ii) for every $s \in \mathbb{R}_+$, the mapping $\theta \mapsto \lambda(s, \theta)$ is in $C^4(\Theta)$ and admits a continuous extension to $\tilde{\Theta}$.

Assumption 7. (i) $\sup_{t \in \mathbb{R}_+} \sum_{i=0}^4 \left\| \sup_{\theta \in \Theta} \partial_{\theta}^i \lambda(t, \theta) \right\|_p < \infty$ for every $p > 1$.

(ii) $\sup_{t \in \mathbb{R}_+} \left\| \sup_{\theta \in \Theta} |\lambda^\alpha(t, \theta)^{-1} 1_{\{\lambda^\alpha(t, \theta) \neq 0\}}| \right\|_p < \infty$ for $p > 1$ and $\alpha \in \mathbf{I}$.

(iii) For any $\theta \in \Theta$ and $\alpha \in \mathbf{I}$, $\lambda^\alpha(t, \theta) = 0$ if and only if $\lambda^\alpha(t, \theta^*) = 0$.

Assumption 8. For every $(\alpha, \theta) \in \mathbf{I} \times \Theta$, there exists a probability measure $\nu^\alpha(\cdot, \theta)$ on $\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p$ and $0 < \delta < \frac{1}{2}$ such that

$$\sup_{\theta \in \Theta} T^\delta \left\| \frac{1}{T} \int_0^T f(\lambda^\alpha(t, \theta^*), \lambda^\alpha(t, \theta), \partial_\theta \lambda^\alpha(t, \theta)) dt - \int f(x, y, z) \nu^\alpha(dx, dy, dz, \theta) \right\|_p \rightarrow 0$$

as $T \rightarrow \infty$ for $p > 1$ and $f \in C_\uparrow(\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p)$.

Let $\nu^\alpha(dx, dy, \theta) = \int_{\mathbb{R}^p} \nu^\alpha(dx, dy, dz, \theta)$. Define $\mathbb{Y}_T(\theta)$ by

$$\mathbb{Y}_T(\theta) = \frac{1}{T} (\ell_T(\theta) - \ell_T(\theta^*)),$$

and $\mathbb{Y}(\theta)$ by the limit in probability of $\mathbb{Y}_T(\theta)$, where

$$\mathbb{Y}(\theta) = \sum_{\alpha \in \mathbf{I}} \int_{\mathbb{R}_+ \times \mathbb{R}_+} 1_{\{x, y > 0\}} \{x \log(y/x) - (y - x)\} \nu^\alpha(dx, dy, \theta).$$

Remark. From the above expression of $\mathbb{Y}(\theta)$, we easily obtain $\mathbb{Y}(\theta^*) = 0$ and for all $\theta \in \Theta$,

$$\mathbb{Y}(\theta) \leq 0. \tag{3.2}$$

Then Lemma 3.10 of [12] gives

$$\sup_{\theta \in \Theta} |\mathbb{Y}_T(\theta) - \mathbb{Y}(\theta)| \xrightarrow{p} 0$$

as $T \rightarrow \infty$.

The index χ_0 is defined by

$$\chi_0 = \inf_{\theta \in \Theta \setminus \{\theta^*\}} \frac{-\mathbb{Y}(\theta)}{|\theta - \theta^*|^2}.$$

Then identifiability is ensured by the condition

Assumption 9. $\chi_0 > 0$.

The Fisher information matrix is well defined by

$$\Gamma = \sum_{\alpha \in \mathbf{I}} \int_{\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p} z^{\otimes 2} x^{-1} 1_{\{x>0\}} \nu^\alpha(dx, dy, dz, \theta^*).$$

The matrix Γ is non-degenerate by Assumption 9.

By Theorem 3.14 of [12], we have

Theorem 3.1. Suppose that Assumptions 6-9 are satisfied. Then for $\tilde{\theta} = \tilde{\theta}^M$ and $\tilde{\theta}^B$, the convergence

$$\lim_{T \rightarrow \infty} E[f(\sqrt{T}(\tilde{\theta} - \theta^*))] = E[f(\Gamma^{-1/2}\zeta)]$$

holds for all $f \in C(\mathbb{R}^p)$ of polynomial growth, where ζ is a p -dimensional standard normal random variable.

Now we are on the point of applying it to the penalized methods. Take $\tilde{\theta} = \tilde{\theta}^M$ or $\tilde{\theta}^B$. The penalized estimator will be denoted by $\hat{\theta}$. Let $r_T = T^{-\frac{1}{2}}$ and let

$$\hat{G} = -T^{-1} \partial_{\theta}^2 \ell_T(\tilde{\theta}) 1_{\{-\partial_{\theta}^2 \ell_T(\tilde{\theta}) \in \mathcal{S}_+\}} + T^{-1} I_p$$

where \mathcal{S}_+ is the set of $p \times p$ positive definite symmetric matrices.

It is easy to show

$$\lim_{T \rightarrow \infty} \|T^\delta (\hat{G} - \Gamma)\|_p = 0$$

for every $p > 1$ and $0 < \delta < \frac{1}{2}$. Therefore the conditions in Theorems 2.1-2.5 are fulfilled in this situation.

3.1.1 Cox type of process with ergodic covariates

Regularization methods for Cox proportional hazards model are proposed by, for example, [9] and [23]. Here, we consider the multivariate point process N in Section 3.1 with intensity processes

$$\lambda^\alpha(t, \theta) = \exp\left(\sum_{j \in \mathbf{J}} \theta_j^\alpha X_t^j\right), \quad (\alpha \in \mathbf{I}) \quad (3.3)$$

where $\mathbf{J} = \{1, \dots, J\}$ is an index set and $X^j = (X_t^j)_{t \in \mathbb{R}_+}$ ($j \in \mathbf{J}$) are left-continuous adapted stochastic covariate processes satisfying the following conditions.

Assumption 10. The J -dimensional process $(X^j)_{j \in \mathbf{J}}$ is stationary and $E[\exp(uX_0^j)] < \infty$ for all $u \in \mathbb{R}$ and $j \in \mathbf{J}$.

Denote by \mathcal{B}_I the σ -field generated by $\{X_t^j; t \in I, j \in \mathbf{J}\}$ for $I \subset \mathbb{R}_+$.
Let

$$\alpha(h) = \sup_{A \in \mathcal{B}_{[0,t]}, B \in \mathcal{B}_{[t+h,\infty)}} |P[A \cap B] - P[A]P[B]|$$

for $h > 0$.

Assumption 11. There exists $a > 0$ such that $\alpha(h) \leq a^{-1}e^{-ah}$ for all $h > 0$.

Let $X_t = (X_t^j)_{j \in \mathbf{J}}$. For the model (3.3), $\theta = (\theta_j^\alpha)_{\alpha \in \mathbf{I}, j \in \mathbf{J}}$, $\mathbf{p} = \mathbf{dJ}$ and

$$\hat{G} = \text{diag} [\hat{G}_1, \dots, \hat{G}_d]$$

where

$$\hat{G}_\alpha = \frac{1}{T} \int_0^T X_t^{\otimes 2} \exp\left(\sum_{j \in \mathbf{J}} \tilde{\theta}_j^\alpha X_t^j\right) dt + \frac{1}{T} I_J. \quad (3.4)$$

It should be remarked that the first term on the right hand side of (3.4) may degenerate in general. Under Assumptions 10 and 11, we obtain $\hat{G} \xrightarrow{p} \Gamma$ for $\Gamma = \text{diag} [\Gamma_1(\theta^*), \dots, \Gamma_d(\theta^*)]$, where

$$\Gamma_\alpha(\theta) = E \left[X_0^{\otimes 2} \exp\left(\sum_{j \in \mathbf{J}} \theta_j^\alpha X_0^j\right) \right].$$

Write $\Gamma(\theta) = \text{diag} [\Gamma_1(\theta), \dots, \Gamma_d(\theta)]$.

Assumption 12. $\inf_{\theta \in \Theta} \det \Gamma(\theta) > 0$.

We assume that Θ is an open bounded convex subset in $\mathbb{R}^{\mathbf{p}}$, that admits the Sobolev inequality in Section 3.1.

Lemma 3.2. Assumption 8 holds under Assumptions 10 and 11.

Proof. We remark that $\exp(|x|) < \exp(x) + \exp(-x)$ for all $x \in \mathbb{R}$. Thus, for all $(\theta, \alpha, j) \in \Theta \times \mathbf{I} \times \mathbf{J}$ and $p, q > 1$ and $t > 0$,

$$\begin{aligned} E \left[|X_t^j|^p \left\{ \exp(\theta_j^\alpha X_t^j) \right\}^q \right] &\leq E \left[\exp(p|X_t^j|) \exp(q|\theta_j^\alpha| |X_t^j|) \right] \\ &= E \left[\exp \left\{ (p + q|\theta_j^\alpha|) |X_0^j| \right\} \right] < C_{p,q}, \end{aligned} \quad (3.5)$$

where $C_{p,q}$ is a constant depend on p, q but not depending of θ, i, j . By the definition of $\lambda^\alpha(t, \theta)$, for all $\alpha \in \mathbf{I}$,

$$\lambda^\alpha(t, \theta) = \exp\left(\sum_{j \in \mathbf{J}} \theta_j^\alpha X_t^j\right),$$

$$\partial_{\theta^{\alpha'}} \lambda^\alpha(t, \theta) = \begin{cases} X_t \exp\left(\sum_{j \in \mathbf{J}} \theta_j^\alpha X_t^j\right) & \text{if } \alpha' = \alpha \\ 0 & \text{if } \alpha' \neq \alpha \end{cases},$$

where $\theta^\alpha = [\theta_j^\alpha]_j$. For $f \in C_\uparrow(\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p)$, $\alpha \in \mathbf{I}$ and $\theta \in \Theta$, define \tilde{f}_θ^α by

$$\tilde{f}_\theta^\alpha(x) = f(e^{\sum_j \theta_j^{*\alpha} x_j}, e^{\sum_j \theta_j^\alpha x_j}, [1_{\{\alpha'=\alpha\}} x e^{\sum_j \theta_j^{\alpha'} x_j}]_{\alpha' \in \mathbf{I}})$$

for $x \in \mathbb{R}^J$. Then we can write for all $\alpha \in \mathbf{I}$,

$$\tilde{f}_\theta^\alpha(X_t) = f(\lambda^\alpha(t, \theta^*), \lambda^\alpha(t, \theta), \partial_\theta \lambda^\alpha(t, \theta)).$$

By (3.5), we obtain for all $\alpha \in \mathbf{I}$ and $p > 1$,

$$\sup_{\theta \in \Theta} E[|\tilde{f}_\theta^\alpha(X_t)|^p] = \sup_{\theta \in \Theta} E[|\tilde{f}_\theta^\alpha(X_0)|^p] < \infty.$$

Here, for $(\alpha, \theta) \in \mathbf{I} \times \Theta$ we define a probability measure $\nu^\alpha(\cdot, \theta)$ by

$$\nu^\alpha(A, \theta) = P\left[\left(e^{\sum_j \theta_j^{*\alpha} X_0^j}, e^{\sum_j \theta_j^\alpha X_0^j}, [1_{\{\alpha'=\alpha\}} X_0 e^{\sum_j \theta_j^{\alpha'} X_0^j}]_{\alpha' \in \mathbf{I}}\right) \in A\right]$$

for $A \in \mathcal{B}(\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p)$. Since

$$\int f(x, y, z) \nu^\alpha(dx, dy, dz, \theta) = E[\tilde{f}_\theta^\alpha(X_0)]$$

for all $\alpha \in \mathbf{I}$, we may show that for all $p > 1$, $\alpha \in \mathbf{I}$ and \tilde{f}_θ^α ,

$$\sup_{\theta \in \Theta} T^{\frac{1}{2}} \left\| \frac{1}{T} \int_0^T \left\{ \tilde{f}_\theta^\alpha(X_t) - E[\tilde{f}_\theta^\alpha(X_0)] \right\} dt \right\|_p = O(1).$$

By Assumption 10, there exists a constant C_0 such that

$$\left\| \int_{s_1}^{s_2} \left(\tilde{f}_\theta^\alpha(X_t) - E[\tilde{f}_\theta^\alpha(X_0)] \right) dt \right\|_p \leq C_0 (s_2 - s_1)^p$$

for $s_1 < s_2$. Then, Lemma 4 in [46] implies under Assumption 11 that

$$\begin{aligned} & E \left[\left| \int_0^T \left(\tilde{f}_\theta^\alpha(X_t) - E[\tilde{f}_\theta^\alpha(X_0)] \right) dt \right|^p \right] \\ &= E \left[\left| \sum_{l=1}^{\lfloor T \rfloor} \int_{\frac{(l-1)T}{\lfloor T \rfloor}}^{\frac{lT}{\lfloor T \rfloor}} \left(\tilde{f}_\theta^\alpha(X_t) - E[\tilde{f}_\theta^\alpha(X_0)] \right) dt \right|^p \right] \\ &\leq C_1 [T]^{\frac{p}{2}} + C_2 [T] \\ &= O(T^{\frac{p}{2}}), \end{aligned}$$

for $T \geq 1$ and $p \geq 2$, where C_1 and C_2 are constants depending on a and p . Therefore, we have

$$\sup_{\theta \in \Theta} T^{\frac{1}{2}} \left\| \frac{1}{T} \int_0^T \left(\tilde{f}_\theta^\alpha(X_t) - E[\tilde{f}_\theta^\alpha(X_0)] \right) dt \right\|_p = O(1).$$

□

Next, we will give a sufficient condition for Assumption 9.

Lemma 3.3. We assume that Θ is convex. Then Assumption 9 follows from Assumptions 10 and 12.

Proof. By the definition of $\mathbb{Y}(\theta)$,

$$\mathbb{Y}(\theta) = \sum_{\alpha \in \mathbf{I}} \mathbb{Y}^\alpha(\theta),$$

where $\mathbb{Y}^\alpha(\theta)$ is given by

$$\begin{aligned} \mathbb{Y}^\alpha(\theta) = E \left[\exp \left(\sum_{j \in \mathbf{J}} \theta_j^{*\alpha} X_0^j \right) \left(\sum_{j \in \mathbf{J}} (\theta_j^\alpha - \theta_j^{*\alpha}) X_0^j \right) \right. \\ \left. - \left\{ \exp \left(\sum_{j \in \mathbf{J}} \theta_j^\alpha X_0^j \right) - \exp \left(\sum_{j \in \mathbf{J}} \theta_j^{*\alpha} X_0^j \right) \right\} \right] \end{aligned}$$

for $\alpha \in \mathbf{I}$. Thus we have

$$\partial_\theta \mathbb{Y}(\theta) = \begin{bmatrix} \partial_{\theta^1} \mathbb{Y}^1(\theta) \\ \partial_{\theta^2} \mathbb{Y}^2(\theta) \\ \vdots \\ \partial_{\theta^d} \mathbb{Y}^d(\theta) \end{bmatrix},$$

where

$$\partial_{\theta^\alpha} \mathbb{Y}^\alpha(\theta) = E \left[\left\{ \exp \left(\sum_{j \in \mathbf{J}} \theta_j^{*\alpha} X_0^j \right) - \exp \left(\sum_{j \in \mathbf{J}} \theta_j^\alpha X_0^j \right) \right\} X_0 \right].$$

Similarly,

$$\partial_\theta^2 \mathbb{Y}(\theta) = \text{diag} \left[\partial_{\theta^1}^2 \mathbb{Y}^1(\theta), \partial_{\theta^2}^2 \mathbb{Y}^2(\theta), \dots, \partial_{\theta^d}^2 \mathbb{Y}^d(\theta) \right],$$

where

$$\partial_{\theta^\alpha}^2 \mathbb{Y}^\alpha(\theta) = -E \left[X_0^{\otimes 2} \exp \left(\sum_{j \in \mathbf{J}} \theta_j^\alpha X_0^j \right) \right] = -\Gamma_\alpha(\theta).$$

Therefore, we have

$$\partial_\theta^2 \mathbb{Y}(\theta) = -\Gamma(\theta).$$

By Assumption 12, for all $\theta \in \bar{\Theta}$, $-\partial_\theta^2 \mathbb{Y}(\theta)$ is positive definite. Therefore, for all $\theta \in \Theta$, there exists a vector θ^\dagger satisfying

$$-\mathbb{Y}(\theta) = -\frac{1}{2} \partial_\theta^2 \mathbb{Y}(\theta^\dagger) [(\theta - \theta^*)^{\otimes 2}].$$

By positive definiteness of $-\partial_\theta^2 \mathbb{Y}(\theta)$ and Assumption 10,

$$\begin{aligned} \inf_{\theta \in V(\theta^*) \setminus \{\theta^*\}} \frac{-\mathbb{Y}(\theta)}{|\theta - \theta^*|^2} &= \inf_{\theta \in V(\theta^*) \setminus \{\theta^*\}} \frac{-\partial_\theta^2 \mathbb{Y}(\theta^\dagger) [(\theta - \theta^*)^{\otimes 2}]}{2|\theta - \theta^*|^2} \\ &\geq \inf_{\theta \in V(\theta^*)} \frac{1}{2} \tau_{\min}(\Gamma(\theta)) > 0. \end{aligned} \quad (3.6)$$

Therefore, by (3.6), we obtain $\chi_0 > 0$. \square

Assumption 6 follows from (3.3), and Assumption 7 follows from Assumption 10. Thus the conditions in Theorems 2.1-2.5 are fulfilled under Assumptions 9-11.

3.1.2 Hawkes process

Hawkes (1971a[21], 1971b[22]) introduced a multivariate model for point processes with mutually exciting components now referred to as the Hawkes model. Ordinary, it was motivated by modeling after shocks and seismological phenomena (cf. Vere-Jones 1970 [40], Vere-Jones and Ozaki 1982 [41], Ogata 1999 [32]). However, the usage of the Hawkes model has been more and more spread out to various research areas : high-frequency finance (c.f. [3]), crime activity (c.f. [30]) and analysis of social networks (c.f. [13], [7], [51], [45]).

Example. For example, consider a finite network with d nodes (each node corresponding to a user in a social network). For each node α in $\{1, \dots, d\}$, we observe the timestamps $\{t_{\alpha,1}, t_{\alpha,2}, \dots\}$ of actions of node α on the network (a message, a click etc.). To each node α is associated a counting process $N^\alpha(t) = \sum_{i \geq 1} 1_{\{t_{\alpha,i} \leq t\}}$. If $N = (N^\alpha)_\alpha$ is a Hawkes process with intensity (3.7), then we can quantify the influence of β on α by the function $h_{\alpha\beta}(t)$.

Now we turn to the definition of a Hawkes process. Let $N_t = (N_t^\alpha)_{\alpha \in \mathbf{I}}$, $\mathbf{I} = \{1, \dots, d\}$, $N_0 = 0$, be a multidimensional point process and write $\mathbf{F} = (\mathcal{F}_t^N)_{t \in \mathbb{R}_+}$, where we call that $\mathcal{F}_t^N = \sigma\{N_s | 0 \leq s \leq t\}$ is the canonical filtration of N . We say that N is a linear Hawkes process or Hawkes' self-exciting process starting from 0 if there exist functions $h_{\alpha\beta} : \mathbb{R}_+ \rightarrow \mathbb{R}_+^{\mathbf{d} \times \mathbf{d}}$ and $\nu \in \mathbb{R}^{\mathbf{d}}$ such that the \mathcal{F}^N -intensity $\lambda(t)$ of N writes

$$\lambda^\alpha(t) = \nu_\alpha + \sum_{\beta \in \mathbf{I}} \int_0^{t-} h_{\alpha\beta}(t-s) dN_s^\beta, \quad \alpha \in \mathbf{I}. \quad (3.7)$$

The baseline intensities ν_α 's represent the rate of spontaneous occurrences of events, while the kernels $h_{\alpha\beta}$'s model self-interaction in the system. Indeed, if a shock occurs at time t_0 on the covariate N^β , an aftershock will happen on the covariate N^α around time t_1 with high probability if $h_{\alpha\beta}(t_1 - t_0)$ is large. When $h_{\alpha\beta} = 0$, covariate N^β has no influence on the chain of events related to N^α .

Let us define the matrix $\Phi = [\phi_{\alpha\beta}]_{\alpha\beta}$ where

$$\phi = \int_0^\infty h_{\alpha\beta}(s) ds,$$

and write $\rho(\Phi)$ its spectral radius. Finally, given $A = [a_{\alpha\beta}]_{\alpha\beta}$ and $C = [c_{\alpha\beta}]_{\alpha\beta} \in \mathbb{R}_+^{\mathbf{d} \times \mathbf{d}}$, we say that N is an exponential Hawkes process if the kernel functions $h_{\alpha\beta}$ are of the form

$$h_{\alpha\beta}(s) = c_{\alpha\beta} e^{-a_{\alpha\beta} s}.$$

Note that the matrix Φ has the representation $\Phi = [\frac{c_{\alpha\beta}}{a_{\alpha\beta}}]_{\alpha\beta}$ in this case. Hereafter, we assume $\rho(\Phi) < 1$, where $\rho(\Phi)$ is a spectral radius of Φ .

Now, we consider the exponential Hawkes process. Then the intensity has the following representation.

$$\lambda^\alpha(t) = \nu_\alpha + \sum_{\beta \in \mathbf{I}} \int_0^{t-} c_{\alpha\beta} e^{-a_{\alpha\beta}(t-s)} dN_s^\beta, \quad \alpha \in \mathbf{I}. \quad (3.8)$$

Here let $\theta = (\nu, C, A) \in \mathbb{R}_+^{\mathbf{d}} \times \mathbb{R}_+^{\mathbf{d} \times \mathbf{d}} \times \mathbb{R}_+^{\mathbf{d} \times \mathbf{d}}$. First, we consider the usual case, i.e., the true value of the parameter $\theta^* = (\nu^*, C^*, A^*)$ is interior point of Θ .

Assumption 13. For every $(\alpha, \theta) \in \mathbf{I} \times \Theta$, there exists a probability measure $\nu^\alpha(\cdot, \theta)$ on $\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p$ and $0 < \delta < \frac{1}{2}$ such that

$$\frac{1}{T} \int_0^T f(\lambda^\alpha(t, \theta^*), \lambda^\alpha(t, \theta), \partial_\theta \lambda^\alpha(t, \theta)) dt \xrightarrow{p} \int f(x, y, z) \nu^\alpha(dx, dy, dz, \theta)$$

as $T \rightarrow \infty$ for $f \in C_b(\mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}^p)$.

Assumption 14. For any $\theta \in \bar{\Theta} - \{\theta^*\}$, $\mathbb{Y}(\theta) \neq 0$.

Proposition 3.4 (Theorem 3.9. of Clinet and Yoshida 2017 ([12])). Under the Assumptions 6, 7, 13 and 14, the QMLE $\tilde{\theta}$ is consistent.

$$\tilde{\theta} \xrightarrow{p} \theta^*.$$

Next, we are interested in the case that some component of Φ^* , the true value of Φ , is precisely 0. Then the graph structure has a form like the left side of Figure 3.1 in contrast with the right side which represents the full model.

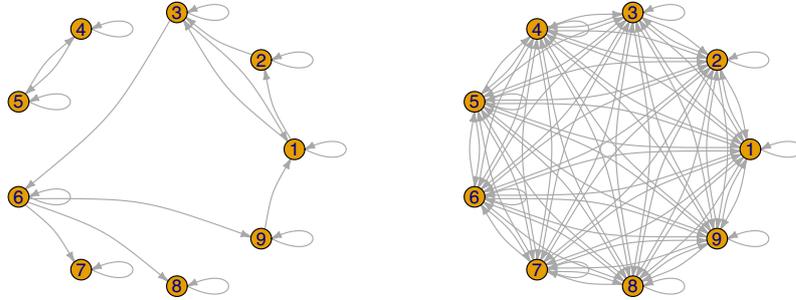


Figure 3.1: Graph structure; sparse model (left) and full model (right)

Now we consider the case where $\phi_{\alpha\beta}^* = 0$ for some (α, β) . Let

$$\mathcal{J}^* = \{(\alpha, \beta) \in \mathbf{I} \times \mathbf{I}; c_{\alpha\beta}^* = 0\}.$$

For a constant $\xi \in \mathbb{R}_+$ we define a function $p^\xi : \Theta \rightarrow \Theta$ by $p^\xi(\theta) = p^\xi(\nu, C, A) = (\nu, C, p_0^\xi(A))$, where $p_0^\xi(A) = [a_{\alpha\beta}^{(\xi)}]_{\alpha\beta}$ and

$$a_{\alpha\beta}^{(\xi)} = \begin{cases} a_{\alpha\beta} & (\alpha, \beta) \in \mathcal{J}^* \\ \xi & (\alpha, \beta) \notin \mathcal{J}^*. \end{cases}$$

If $c_{\alpha\beta}^* = 0$ for some (α, β) , then the value of $a_{\alpha\beta}^*$ does not affect the model and Proposition 3.4 does not hold. Thus we modify the Assumption 14.

Assumption 15. For any $\theta \in \bar{\Theta} - \{p^\xi(\theta^*) \in \Theta; \xi \in \mathbb{R}_+\}$, $\mathbb{Y}(\theta) \neq 0$.

Proposition 3.5. Under the assumptions 6, 7, 13 and 15, (the QMLE) $p^\xi(\tilde{\theta})$ is consistent.

$$p^\xi(\tilde{\theta}) \xrightarrow{p} p^\xi(\theta^*)$$

In particular,

$$\tilde{C} \xrightarrow{p} C^*.$$

Since the structure of C (or Φ) determine the model, define the pLSA estimator \hat{C} by

$$\hat{C} \in \operatorname{argmin}_C \hat{G}[(C - \tilde{C})^{\otimes 2}] + \sum_{\alpha, \beta} \kappa_T^{\alpha\beta} |c_{\alpha\beta}|^q,$$

where $\kappa_T^{\alpha\beta} = \alpha_T(|c_{\alpha\beta}| + \eta_T)^{-\gamma}$, $\eta_T \rightarrow 0$ ($T \rightarrow \infty$) and $\hat{G} = -\frac{1}{T} \partial_C^2 \ell_T(\theta)$. Then by the Theorem 2.2, we obtain the selection consistency of the pLSA estimator \hat{C} .

Remark. By the Theorem 4.6. of [12], the exponential Hawkes model verifies the Assumptions 6-9 in the usual situation. By using the discussion in Appendix, same results are possibly derived in the situation where the true value of parameter $\theta^* \in \partial\Theta$. However, since it is difficult to treat that case, we omit the detailed description here.

3.2 Diffusion process

3.2.1 Ergodic case

Given a stochastic basis $(\Omega, \mathcal{F}, \mathbf{F}, P)$, $\mathbf{F} = (\mathcal{F}_t)_{t \in \mathbb{R}_+}$, we consider a d -dimensional process $X = (X_t)_{t \in \mathbb{R}_+}$ adapted to the filtration $\mathbf{F} = (\mathcal{F}_t)_{t \in \mathbb{R}_+}$ and satisfying the following stochastic integral equation

$$X_t = X_0 + \int_0^t a(X_s, \theta_2) ds + \int_0^t b(X_s, \theta_1) dW_s, \quad t \in \mathbb{R}_+$$

where W is an r -dimensional standard \mathbf{F} -Wiener process, $\theta = (\theta_1, \theta_2) \in \Theta_1 \times \Theta_2 = \Theta$ with Θ_1 and Θ_2 being bounded domains of \mathbb{R}^{p_1} and \mathbb{R}^{p_2} , respectively, moreover $b : \mathbb{R}^d \times \Theta_1 \rightarrow \mathbb{R}^d \otimes \mathbb{R}^r$ and $a : \mathbb{R}^d \times \Theta_2 \rightarrow \mathbb{R}^d$. We define the function B by $B(x, \theta_1) = b(x, \theta_1)b(x, \theta_1)'$ and assume that $B(x, \theta_1)$ is invertible. We denote the true value of $\theta = (\theta_1, \theta_2)$ by $\theta^* = (\theta_1^*, \theta_2^*)$ and the number of active parameters of θ_k^* by p_k^0 for $k = 1, 2$. We assume that each parameter space have a locally Lipschitz boundary.

In this subsection, we assume that the process X is ergodic. That is, there exists a unique invariant probability measure $\mu = \mu_{\theta^*}$ such that for any bounded measurable function $g : \mathbb{R}^d \rightarrow \mathbb{R}$, the convergence

$$\frac{1}{T} \int_0^T g(X_t) dt \rightarrow^p \int_{\mathbb{R}^d} g(x) \mu(dx)$$

holds.

We suppose that $0 \in \mathbb{R}^{p_1+p_2}$ is in Θ . Here we have the discrete-time observations $(X_{t_i}, Y_{t_i})_{i=0}^n$ where $t_i = ih$ with $h = h_n$ depending on n . We will consider the situation when $h_n \rightarrow 0$ and $nh_n^p \rightarrow 0$ as $n \rightarrow \infty$, and there exists $\epsilon_0 \in (0, \frac{p-1}{p})$ such that $n^{\epsilon_0} \leq nh_n$ for large n .

Here, we assume the following properties of an initial estimator $\tilde{\theta} = (\tilde{\theta}_{1,n}, \tilde{\theta}_{2,n})$:

$$\begin{aligned} & (\sqrt{n}(\tilde{\theta}_{1,n} - \theta_1^*), \sqrt{nh}(\tilde{\theta}_{2,n} - \theta_2^*)) \\ & \rightarrow^d (\Gamma_1^{-\frac{1}{2}} \zeta_1, \Gamma_2^{-\frac{1}{2}} \zeta_2) \sim N_{p_1+p_2}(0, \text{diag}(\Gamma_1^{-1}, \Gamma_2^{-1})) \end{aligned}$$

and

$$\sup_n \left(\|\sqrt{n}(\tilde{\theta}_{1,n} - \theta_1^*)\|_p + \|\sqrt{nh}(\tilde{\theta}_{2,n} - \theta_2^*)\|_p \right) < \infty$$

for every $p > 1$, where ζ_1 and ζ_2 are p_1 and p_2 -dimensional standard normal variables respectively, and $\Gamma_1 = (\Gamma_1^{ij})_{i,j=1,\dots,p_1}$ and $\Gamma_2 = (\Gamma_2^{ij})_{i,j=1,\dots,p_2}$ with

$$\begin{aligned} \Gamma_1^{ij} &= \frac{1}{2} \int_{\mathbb{R}^d} \text{Tr} \left\{ (\partial_{\theta_1^i} B(x, \theta_1^*)) B^{-1}(x, \theta_1^*) (\partial_{\theta_1^j} B(x, \theta_1^*)) B^{-1}(x, \theta_1^*) \right\} \mu(dx), \\ \Gamma_2^{ij} &= \int_{\mathbb{R}^d} (\partial_{\theta_2^i} a(x, \theta_2^*))' B(x, \theta_1^*)^{-1} \partial_{\theta_2^j} a(x, \theta_2^*) \mu(dx). \end{aligned}$$

We assume integrability and non-degeneracy of Γ_1 and Γ_2 . It is known that the quasi-maximum likelihood estimator, the quasi-Bayesian estimator and the hybrid type estimators possess these properties under certain mild conditions ([46], [37], [39], [25]). For instance, if we use the hybrid multistep estimator $\tilde{\theta}^H = (\tilde{\theta}_{1,n}^H, \tilde{\theta}_{2,n}^H)$ by Uchida and Kamatani ([25]) as an initial estimator $\tilde{\theta}$, then above conditions are satisfied by Theorem 1 of [25].

For $q_1, q_2 \in (0, 1]$, we define the objective functions $Q_{1,n}^{(q_1)}$ and $Q_{2,n}^{(q_2)}$ by

$$Q_{1,n}^{(q_1)} = \hat{G}_{1,n}[(\theta_1 - \tilde{\theta}_{1,n})^{\otimes 2}] + \sum_{i_1=1}^{p_1} \kappa_{1,n}^{i_1} |\theta_1^{i_1}|^{q_1}$$

and

$$Q_{2,n}^{(q_2)} = \hat{G}_{2,n}[(\theta_2 - \tilde{\theta}_{2,n})^{\otimes 2}] + \sum_{i_2=1}^{p_2} \kappa_{2,n}^{i_2} |\theta_2^{i_2}|^{q_2},$$

respectively, where $\hat{G}_{k,n}$ ($k = 1, 2$) are some $\mathbf{p}_k \times \mathbf{p}_k$ random matrices such that $\hat{G}_{k,n} \xrightarrow{p} \Gamma_k$ and that the family $\{|\hat{G}_{k,n}| + (\det \hat{G}_{k,n})^{-1}\}_{k,n}$ is L^∞ -bounded, and $\kappa_{k,n}^{i_k} = \alpha_{k,n} |\tilde{\theta}_{k,n}^{i_k}|^{-\gamma_k}$ for some numbers $\gamma_k > -(1 - q_k)$ and some sequences $(\alpha_{1,n})_n$ and $(\alpha_{2,n})_n$ satisfying

$$(\sqrt{n})^{2-q_1+\gamma_1} \alpha_{1,n} \rightarrow \infty, \sqrt{n} \alpha_{1,n} \rightarrow 0$$

and

$$(\sqrt{nh})^{2-q_2+\gamma_2} \alpha_{2,n} \rightarrow \infty, \sqrt{nh} \alpha_{2,n} \rightarrow 0$$

respectively. Then we have the penalized LSA estimators $\hat{\theta}_{1,n}^{(q_1)}$ and $\hat{\theta}_{2,n}^{(q_2)}$ satisfying

$$\hat{\theta}_{1,n}^{(q_1)} \in \operatorname{argmin}_{\theta_1 \in \Theta_1} Q_{1,n}^{(q_1)}(\theta_1)$$

and

$$\hat{\theta}_{2,n}^{(q_2)} \in \operatorname{argmin}_{\theta_2 \in \Theta_2} Q_{2,n}^{(q_2)}(\theta_2).$$

For these estimators $\hat{\theta}_{1,n}^{(q_1)}$ and $\hat{\theta}_{2,n}^{(q_2)}$, Theorems 2.1-2.5 hold respectively. Additionally, we consider the limit distribution of the joint variable $((\hat{\theta}_{1,n}^{(q_1)})_{\mathcal{J}_1^1}, (\hat{\theta}_{2,n}^{(q_2)})_{\mathcal{J}_2^1})$. Here, $\mathcal{J}_k^1, \mathcal{J}_k^{11}$ and \mathcal{J}_k^{10} are defined similarly to $\mathcal{J}^1, \mathcal{J}^{11}$ and \mathcal{J}^{10} , respectively, for each $k = 1, 2$.

Now we can rephrase Theorems 1-5 in the present situation. In particular,

Proposition 3.6. The convergence

$$\begin{aligned} \left(\sqrt{n}(\hat{\theta}_{1,n}^{(q_1)} - \theta_1^*)_{\mathcal{J}_1^1}, \sqrt{nh}(\hat{\theta}_{2,n}^{(q_2)} - \theta_2^*)_{\mathcal{J}_2^1} \right) &\rightarrow^d \left(\mathfrak{G}_1 \Gamma_1^{-\frac{1}{2}} \zeta_1, \mathfrak{G}_2 \Gamma_2^{-\frac{1}{2}} \zeta_2 \right) \\ &\sim N_{\mathbf{p}_1^0 + \mathbf{p}_2^0} \left(0, \operatorname{diag} \left(((\Gamma_1)_{\mathcal{J}_1^{11}})^{-1}, ((\Gamma_2)_{\mathcal{J}_2^{11}})^{-1} \right) \right) \end{aligned}$$

holds, where $\mathfrak{G}_k = [I_{\mathbf{p}_k^0} \quad ((\Gamma_k)_{\mathcal{J}_k^{11}})^{-1} (\Gamma_k)_{\mathcal{J}_k^{10}}], k = 1, 2$.

Proof. By Theorem 2.3, we have

$$\sqrt{n}(\hat{\theta}_{1,n}^{(q_1)} - \theta_1^*)_{\mathcal{J}_1^1} - \mathfrak{G}_1\{\sqrt{n}(\tilde{\theta}_{1,n} - \theta_1^*)_{\mathcal{J}_1^1}\} \rightarrow^p 0$$

and

$$\sqrt{nh}(\hat{\theta}_{2,n}^{(q_2)} - \theta_2^*)_{\mathcal{J}_2^1} - \mathfrak{G}_2\{\sqrt{nh}(\tilde{\theta}_{2,n} - \theta_2^*)_{\mathcal{J}_2^1}\} \rightarrow^p 0.$$

Therefore,

$$\begin{aligned} & \begin{bmatrix} \sqrt{n}(\hat{\theta}_{1,n}^{(q_1)} - \theta_1^*)_{\mathcal{J}_1^1} \\ \sqrt{nh}(\hat{\theta}_{2,n}^{(q_2)} - \theta_2^*)_{\mathcal{J}_2^1} \end{bmatrix} \rightarrow^d \begin{bmatrix} \mathfrak{G}_1 \Gamma_1^{-\frac{1}{2}} \zeta_1 \\ \mathfrak{G}_2 \Gamma_2^{-\frac{1}{2}} \zeta_2 \end{bmatrix} \\ & \sim N_{\mathbf{p}_1^0 + \mathbf{p}_2^0} \left(0, \text{diag} \left(((\Gamma_1)_{\mathcal{J}_1^{11}})^{-1}, ((\Gamma_2)_{\mathcal{J}_2^{11}})^{-1} \right) \right) \end{aligned}$$

□

3.2.2 Non-ergodic case

In this subsection, we will deal with the case where the Fisher information matrix is not deterministic. We consider the following stochastic regression model

$$Y_t = Y_0 + \int_0^t b_s ds + \int_0^t \sigma(X_s, \theta) dW_s, \quad t \in [0, T], \quad (3.9)$$

where W is an r -dimensional standard Wiener process independent of the initial value of Y_0 , X and b are progressively measurable processes with values in \mathbb{R}^d and \mathbb{R}^m , respectively. σ is an $\mathbb{R}^m \otimes \mathbb{R}^r$ -valued measurable function defined on $\mathbb{R}^d \times \Theta$, and Θ is a bounded domain in \mathbb{R}^p with a locally Lipschitz boundary. Additionally, we define $S = \sigma^{\otimes 2} = \sigma\sigma'$. The data set consists of discrete observations $(X_{t_j}, Y_{t_j})_{j=0}^n$ with $t_j = jT/n$ and T is fixed.

Here, we assume that there exists an estimator $\tilde{\theta}_n$ of θ^* such that

$$\sqrt{n}(\tilde{\theta}_n - \theta^*) \rightarrow^{d_s} \Gamma^{-\frac{1}{2}} \zeta$$

as $n \rightarrow \infty$, and for any continuous functions $f : \mathbb{R}^p \rightarrow \mathbb{R}$ of at most polynomial growth,

$$E[f(\sqrt{n}(\tilde{\theta}_n - \theta^*))] \rightarrow E[f(\Gamma^{-\frac{1}{2}} \zeta)],$$

where Γ is the Fisher information matrix given by

$$\Gamma = \frac{1}{2T} \int_0^T \text{Tr} \left((\partial_\theta S) S^{-1} (\partial_\theta S) S^{-1} (X_t, \theta^*) \right) dt,$$

ζ is a p -dimensional standard normal random variable independent of Γ and \rightarrow^{ds} means the $\sigma(\Gamma)$ -stable convergence in distribution. Here we remark that the Fisher information matrix Γ is not necessarily deterministic. In fact, Uchida and Yoshida [38] proved that the quasi-maximum likelihood estimator and the quasi-Bayesian estimator have these properties under mild regularity conditions. An essential condition in their argument is the non-degeneracy of a key index χ_0 :

Assumption 16. For every $L > 0$, there exists $c_L > 0$ such that

$$P[\chi_0 \leq r^{-1}] \leq \frac{c_L}{r^L} \quad (r > 0)$$

where

$$\chi_0 = \inf_{\theta \neq \theta^*} \frac{-\mathbb{Y}(\theta)}{|\theta - \theta^*|^2}$$

with

$$\begin{aligned} \mathbb{Y}(\theta) = & -\frac{1}{2T} \int_0^T \left\{ \log \left(\frac{\det S(X_t, \theta)}{\det S(X_t, \theta^*)} \right) \right. \\ & \left. + \text{Tr} \left(S^{-1}(X_t, \theta) S(X_t, \theta^*) - I_d \right) \right\} dt. \end{aligned}$$

For the initial estimator, for example, we can take the maximum likelihood type estimator $\tilde{\theta}_n^M$ that satisfies

$$\mathbb{H}_n(\tilde{\theta}_n^M) = \sup_{\theta \in \Theta} \mathbb{H}_n(\theta),$$

or the Bayes type estimator $\tilde{\theta}_n^B$ for a prior density $\pi : \Theta \rightarrow \mathbb{R}_+$ with respect to the quadratic loss defined by

$$\tilde{\theta}_n^B = \left(\int_{\Theta} \exp(\mathbb{H}_n(\theta)) \pi(\theta) d\theta \right)^{-1} \int_{\Theta} \theta \exp(\mathbb{H}_n(\theta)) \pi(\theta) d\theta,$$

where $\mathbb{H}_n(\theta)$ is a quasi-log likelihood function defined by

$$\mathbb{H}_n(\theta) = -\frac{1}{2} \sum_{i=1}^n \left\{ \log \det S(X_{t_{i-1}}, \theta) + \frac{1}{h} S(X_{t_{i-1}}, \theta)^{-1} [(\Delta_i Y)^{\otimes 2}] \right\}. \quad (3.10)$$

Then we can use the QLA method in [38] to show the stable convergence and the L^p -boundedness of the estimators. In order to verify Assumption 16 in practice, we may apply one of criteria given in [38].

Here we define the objective function

$$Q_n^{(q)}(\theta) = \hat{G}_n[(\theta - \tilde{\theta}_n)^{\otimes 2}] + \sum_{i=1}^p \kappa_n^i |\theta_i|^q$$

and penalized LSA estimator $\hat{\theta}_n^{(q)} \in \operatorname{argmin}_{\theta \in \Theta} Q_n^{(q)}(\theta)$. We can take

$$\hat{G}_n = -\frac{1}{n} \partial_{\theta}^2 \mathbb{H}_n(\tilde{\theta}^M) 1_{\{-\partial_{\theta}^2 \mathbb{H}_n(\tilde{\theta}^M) \in \mathcal{S}_+\}} + \frac{1}{n} I_p \quad (3.11)$$

when we use the QMLE as an initial estimator. Let $\kappa_n^i = \alpha_n |\tilde{\theta}_{n,i}|^{-\gamma}$ for the number $\gamma > -(1 - q)$ and the sequence α_n satisfying the conditions

$$(\sqrt{n})^{2-q+\gamma} \alpha_n \rightarrow \infty, \quad \sqrt{n} \alpha_n \rightarrow 0.$$

Similarly to the previous sections, we can show that Theorems 2.1-4 hold for this penalized LSA estimator. On the other hand, we should choose $\hat{G}_n = I_p$, in place of (3.11), for the P-O estimator.

Chapter 4

Simulations

In this section we report three simulations. The first one is the Cox process in Section 3.1.1, the second one is the Hawkes process in Section 3.1.2 and the third one is the non-ergodic diffusion type process in Section 3.2.2. For each simulation we perform 1000 Monte Carlo replications. $\%$ (*method*) in this chapter denotes the number of times, in percentage over 1000 iterations, that the estimator obtained by the *method* chooses the true model. In the first and third cases, for comparison we use the unified LASSO type estimator and the Bridge type estimator. The unified LASSO type estimator is the special case of penalized LSA estimator where $q = 1$. The Bridge estimator is not the special case of penalized LSA estimator, but we use this phraseology when the penalty has the form $r_T^{-1} \sum_i |\theta_i|^q$, i.e., $r = 1, \gamma = 0$ and $q < 1$.

For the convenience of calculation, in the first and third cases, we use identity matrix as a coefficient matrix \hat{G} . Thus, the penalized LSA estimator is not efficient, and we use the P-O estimator in order to obtain the efficient estimator.

It has been shown through the simulation studies that the penalized LSA estimator can select the correct model if we choose appropriate tuning parameters and the P-O estimator has good performance for the active parameters. It is an important thing to give a “good” tuning parameter, however, we do not refer to how to select the tuning parameter.

4.1 Simulation for the Cox model

We consider the Cox model (3.3) in Section 3.1.1 with $\alpha = 1$. Let $\mathbf{p} = 20$, then the parameter space Θ is $[-10, 10]^{20}$. The covariate process $X = (X_t)_{t \in [0, T]}$ is a 20-dimensional OU process satisfying the following stochastic differential

equations

$$dX_t^i = -\mathbf{a}_i X_t^i dt + 0.4dW_t^i, \quad X_0 = 0, \quad t \in [0, T]$$

where \mathbf{a}_i ($i = 1, \dots, 20$) are constants given by

$$\begin{aligned} \mathbf{a}_1 &= \mathbf{a}_6 = \mathbf{a}_{11} = \mathbf{a}_{16} = 0.15, \\ \mathbf{a}_2 &= \mathbf{a}_7 = \mathbf{a}_{12} = \mathbf{a}_{17} = 0.2, \\ \mathbf{a}_3 &= \mathbf{a}_8 = \mathbf{a}_{13} = \mathbf{a}_{18} = 0.25, \\ \mathbf{a}_4 &= \mathbf{a}_9 = \mathbf{a}_{14} = \mathbf{a}_{19} = 0.3, \\ \mathbf{a}_5 &= \mathbf{a}_{10} = \mathbf{a}_{15} = \mathbf{a}_{20} = 0.35. \end{aligned}$$

and $W = (W^i)_{i=1, \dots, 20}$ is a 20-dimensional standard Wiener process. Figure 4.1 shows a sample path of the covariate process. Data $N = (N_t)_{t \in [0, T]}$ is a sample path of the point process with intensity $\lambda(t, \theta^*)$ in (3.3), where the true values θ^* of the parameter is

$$\theta^* = [2, -1, 1, -0.5, -1.5, 1.5, 0.5, 0.75, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]'$$

Figure 4.2 and 4.3 show sample paths of intensity process $\lambda(t, \theta^*)$ and counting process N_t on $[0, 200]$ and $[0, 3]$ respectively.

Let $\mathcal{L}_T(\theta) = \mathbb{L}_T(\theta) = -\ell_T(\theta)$ in (3.1) and we use QMLE for the initial estimator θ . Then the objective function is denoted by

$$Q_T^{(q)}(\theta) = (\theta - \tilde{\theta})'(\theta - \tilde{\theta}) + \sum_{j=1}^{20} \kappa_T^j |\theta_j|^q.$$

where $\kappa_T^j = \alpha_T |\tilde{\theta}_j|^{-\gamma}$, $\alpha_T = (\frac{1}{\sqrt{T}})^r$, $1 < r < 2 - q + \gamma$. Let the triplet of tuning parameters $(\gamma, r, q) = (1, 1.2, 0.3)$. We will consider the cases $T = 50, 100, 200$ and 400. Table 1 compares the results of the variable selection of the penalized LSA estimator, the unified LASSO type estimator and the Bridge type estimator. Here, we define the unified LASSO type estimator and the Bridge type estimator as the penalized LSA estimator with tuning parameter $(\gamma, r, q) = (1, 1.2, 1)$ and $(\gamma, r, q) = (0, 1, 0.3)$, respectively.

Table 2 compares the means and standard deviations (parentheses) for the three estimators (initial estimator, penalized LSA estimator and P-O estimator) and shows the results of the variable selection for penalized LSA estimator in the case $T = 200$.

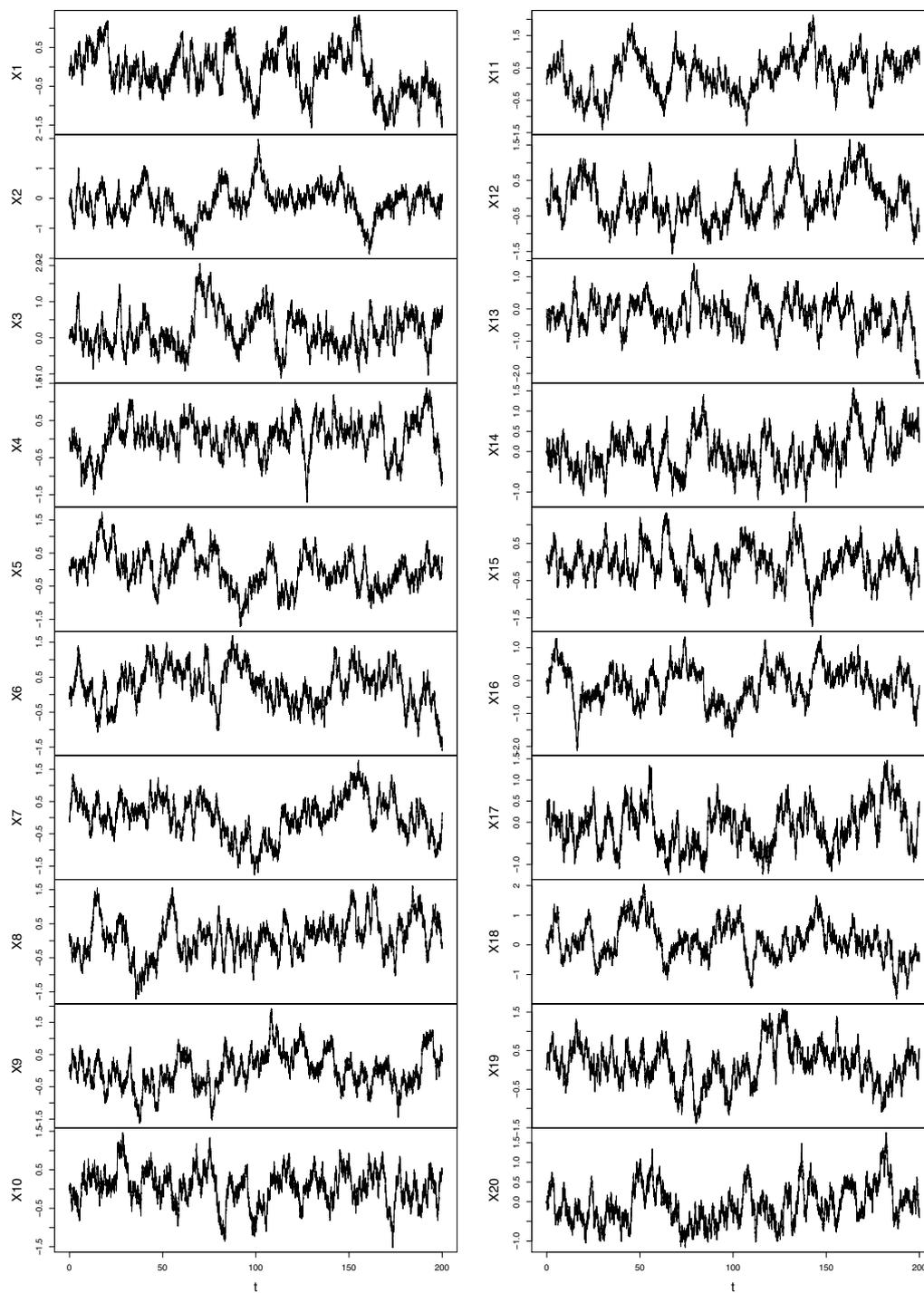
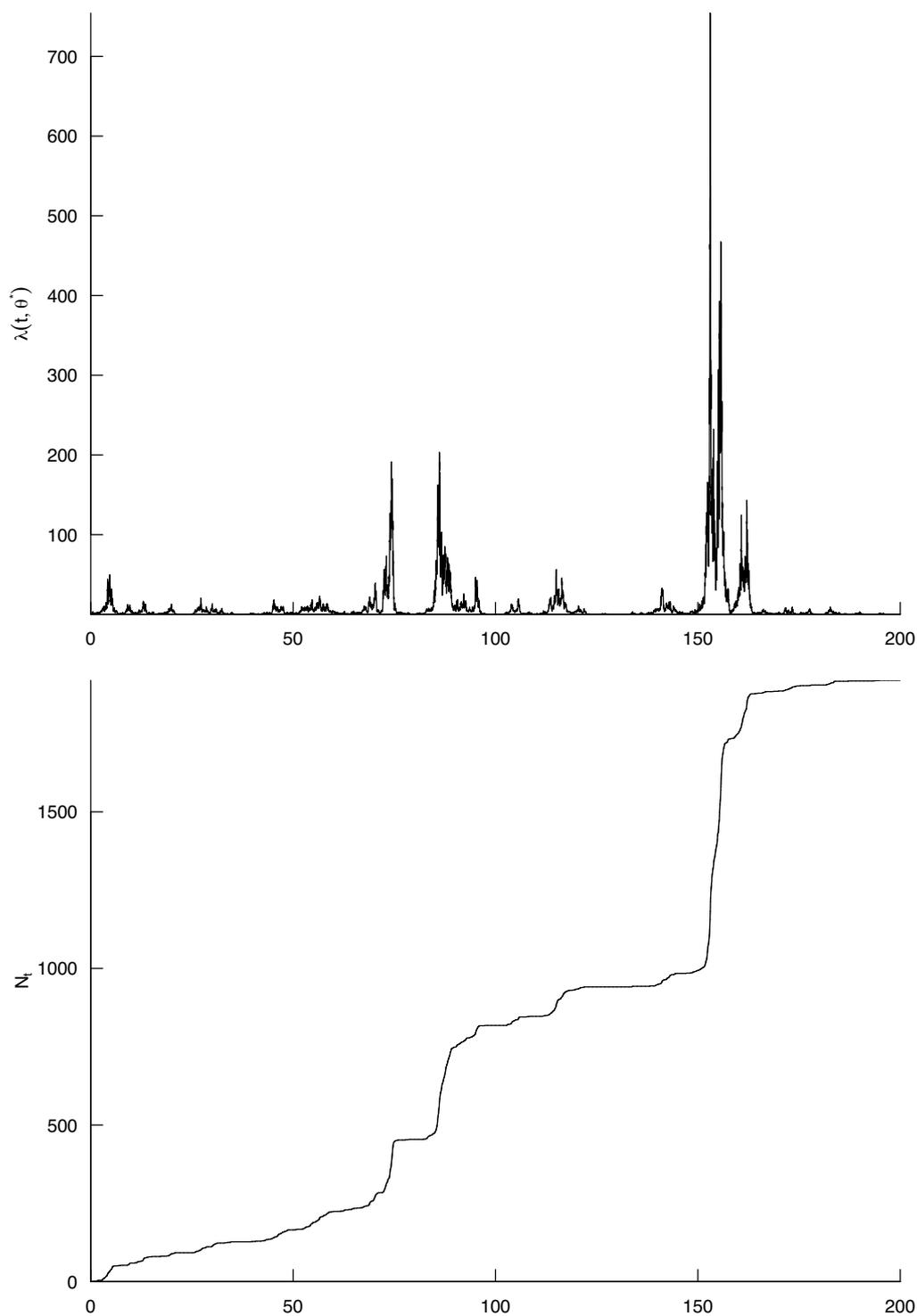


Figure 4.1: Sample path of covariate processes X

Figure 4.2: Sample paths of intensity $\lambda(t)$ and counting processes N (global)

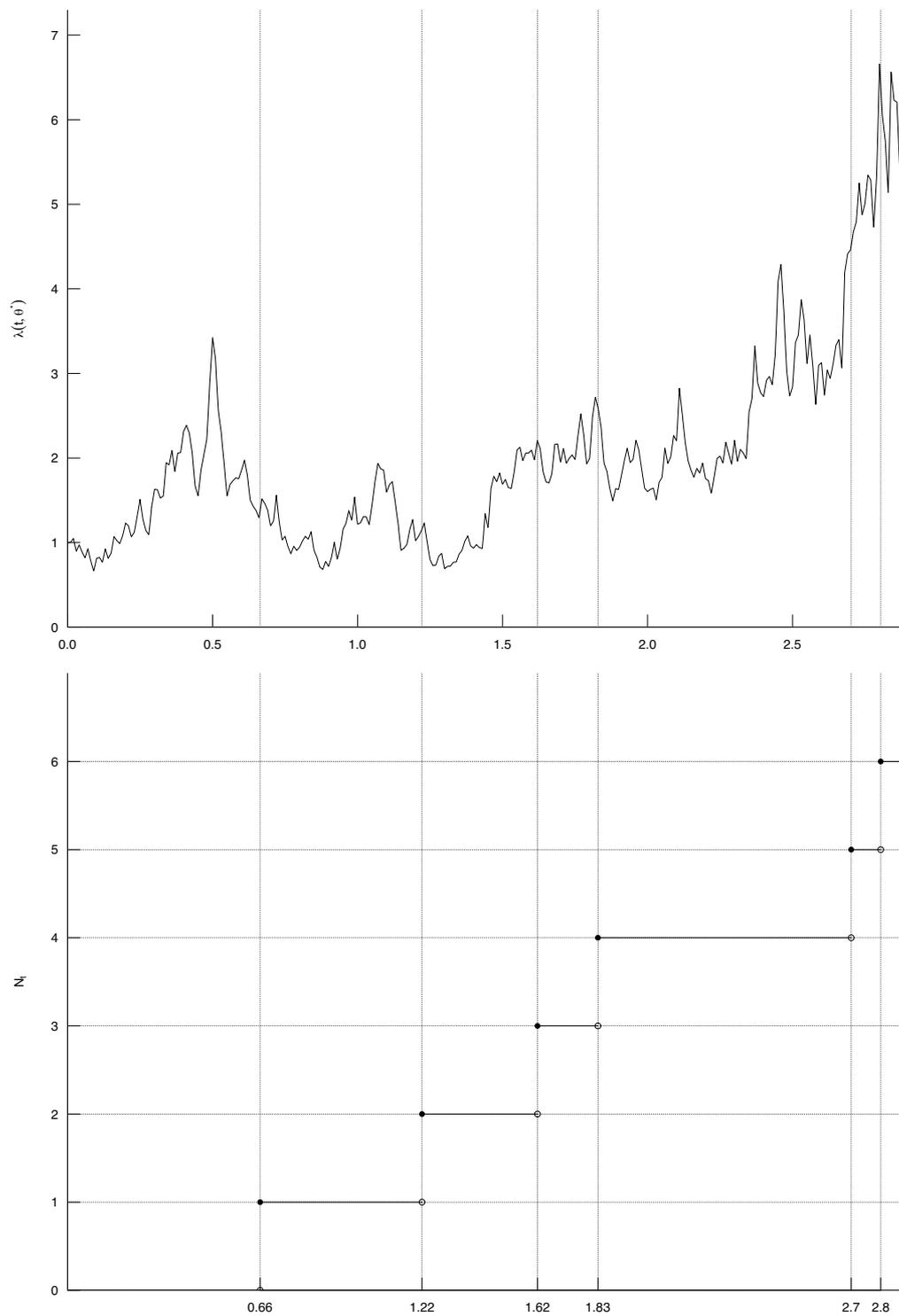


Figure 4.3: Sample paths of intensity $\lambda(t)$ and counting processes N (local)

Table 4.1. Results of the variable selection under $T=50, 100, 200, 400$.

	(γ, r, q)	$T = 50$	$T = 100$	$T = 200$	$T = 400$
%(p-LSA)	(1, 1.2, 0.3)	32.1	70.4	96.8	99.9
%(unified LASSO)	(1, 1.2, 1)	5.9	17.4	47.1	79.6
%(Bridge type)	(0, 1, 0.3)	8.9	21.9	52.3	80.3

Table 4.2. The summary of results for the simulation under $T = 200$.

	θ_1	θ_2	θ_3	θ_4	θ_5	θ_6	θ_7	θ_8	θ_9	θ_{10}
true	2	-1	1	-0.5	-1.5	1.5	0.5	0.75	0	0
initial	1.9938 (0.0722)	-0.9936 (0.0830)	0.9941 (0.0858)	-0.5003 (0.0889)	-1.4909 (0.0962)	1.4978 (0.0745)	0.5009 (0.0851)	0.7490 (0.0880)	0.0001 (0.0908)	-0.0002 (0.0974)
p-LSA	1.9918 (0.0728)	-0.9872 (0.0840)	0.9877 (0.0868)	-0.4758 (0.1035)	-1.4877 (0.0965)	1.4946 (0.0748)	0.4750 (0.1049)	0.7383 (0.0904)	-0.0001 (0.0265)	-0.0008 (0.0340)
P-O	1.9959 (0.0565)	-0.9971 (0.0682)	0.9964 (0.0713)	-0.4979 (0.0892)	-1.4931 (0.0823)	1.4998 (0.0602)	0.4946 (0.0915)	0.7523 (0.0730)	-0.0003 (0.0228)	-0.0007 (0.0307)
%(p-LSA)	100.0	100.0	100.0	98.8	100.0	100.0	98.4	100.0	99.4	99.2

	θ_{11}	θ_{12}	θ_{13}	θ_{14}	θ_{15}	θ_{16}	θ_{17}	θ_{18}	θ_{19}	θ_{20}
true	0	0	0	0	0	0	0	0	0	0
initial	0.0014 (0.0760)	-0.0038 (0.0794)	0.0007 (0.0872)	0.0026 (0.0947)	0.0063 (0.0971)	-0.0015 (0.0701)	0.0042 (0.0802)	0.0013 (0.0896)	-0.0048 (0.0914)	-0.0014 (0.0941)
p-LSA	0.0000 (0.0000)	0.0003 (0.0185)	0.0007 (0.0252)	0.0002 (0.0240)	-0.0004 (0.0248)	0.0000 (0.0000)	0.0011 (0.0194)	-0.0003 (0.0239)	0.0013 (0.0247)	0.0004 (0.0208)
P-O	0.0000 (0.0000)	-0.0003 (0.0152)	0.0007 (0.0221)	0.0000 (0.0188)	-0.0003 (0.0259)	0.0000 (0.0000)	0.0007 (0.0134)	-0.0004 (0.0214)	0.0010 (0.0199)	0.0004 (0.0213)
%(p-LSA)	100.0	99.7	99.4	99.4	99.5	100.0	99.7	99.5	99.7	99.5

4.2 Simulations for the Hawkes process

We consider the Hawkes model in Section 3.1.2. Let $d = 3$ and the parameter space

$$\Theta = \{(\nu, C, A) \in \mathbb{R}_+ \times \mathbb{R}_+^{3 \times 3} \times \mathbb{R}_+^{3 \times 3}; \rho(\Phi) < 1, \nu_\alpha > 0, a_{\alpha\beta} > 0, \forall(\alpha, \beta)\}.$$

For $\alpha \in \mathbf{I}$ we generated the data $(t_i^\alpha)_{i=1,2,\dots}$ which is jump point of the counting process $(N_t^\alpha)_t$ with the intensity

$$\lambda^\alpha(t) = \nu_\alpha + \sum_{\beta \in \mathbf{I}} \int_0^{t^-} c_{\alpha\beta} e^{-a_{\alpha\beta}(t-s)} dN_s^\beta,$$

where $\theta^* = (\nu^*, C^*, A^*)$ is the true value of the parameter;

$$\nu^* = \begin{bmatrix} 1 \\ 0.5 \\ 1 \end{bmatrix}, \quad C^* = \begin{bmatrix} 1 & 0 & 0 \\ 0.5 & 1 & 0 \\ 0 & 0.5 & 1 \end{bmatrix}, \quad A^* = \begin{bmatrix} 2 & * & * \\ 2 & 3 & * \\ * & 3 & 4 \end{bmatrix}.$$

In this simulation, we take $\xi = 1$ (in subsection 3.1.2), i.e.,

$$A^* = \begin{bmatrix} 2 & 1 & 1 \\ 2 & 3 & 1 \\ 1 & 3 & 4 \end{bmatrix}.$$

Then the matrix Φ^* is

$$\Phi^* = \begin{bmatrix} 1/2 & 0 & 0 \\ 1/4 & 1/3 & 0 \\ 0 & 1/6 & 1/4 \end{bmatrix}.$$

Figure 4.4 shows the graph structure of this model.

Let $\mathcal{L}_T(\theta) = -\ell_T(\theta)$ in (3.1) and we use QMLE for the initial estimator $\tilde{\theta} = (\tilde{\nu}, \tilde{C}, \tilde{A})$. Then the objective function is denoted by

$$Q_T^{(q)}(C) = \hat{G}[(C - \tilde{C})^{\otimes 2}] + \sum_{\alpha, \beta \in \mathbf{I}} \kappa_T^{\alpha\beta} |c_{\alpha\beta}|^q,$$

where $\kappa_T^{\alpha\beta} = \alpha_T(|c_{\alpha\beta}| + 1/T)^{-\gamma}$, $\alpha_T = T^{-r/2}$, $1 < r < 2 - q + \gamma$ and $\hat{G} = \partial_C^2 \ell_T(\theta)$. In this simulation, we take a tuning parameter $(\gamma, r, q) = (3, 1.2, 1)$.

Tables 4.3 and 4.4 show the means and standard deviations (parentheses) for the QMLE (with $T = 100$, $T = 200$, $T = 400$ and $T = 600$) and Table 4.5 shows the results of the variable selection for the QMLE and the penalized LSA estimator.

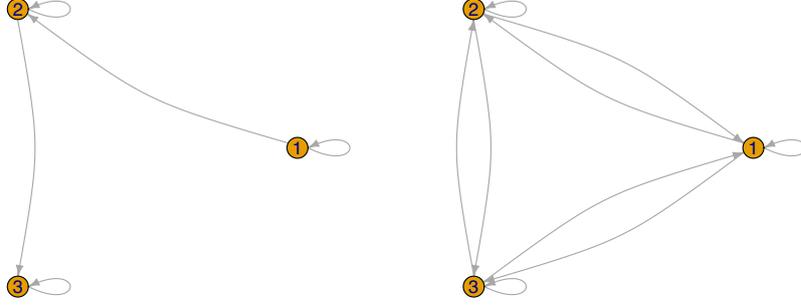


Figure 4.4: Graph structure; true model (left) and full model (right)

Here, we refer to the algorithm for generating the sample path of Hawkes process. Ogata [31] gave a proposition that states the simulation of multivariate point processes, by distributing accepted points to each dimension with probabilities proportional to their intensities.

Lemma 4.1 (Ogata [31]). Let $N_t = (N_t^1, \dots, N_t^d)$ be a d -variate point process on an interval $[0, T]$ with stochastic intensities $\lambda^\alpha(t) = \lambda^\alpha(t|\mathcal{F}_t^N)$ for $\alpha \in \mathbf{I}$. Suppose there is a one-dimensional \mathcal{F}^N -predictable process $\bar{\lambda}(t)$ which is defined path-wisely satisfying

$$\sum_{\alpha \in \mathbf{I}} \lambda^\alpha(t) \leq \bar{\lambda}(t), \quad 0 < t \leq T$$

and set

$$\lambda^0(t) = \bar{\lambda}(t) - \sum_{\alpha \in \mathbf{I}} \lambda^\alpha(t).$$

Let $\bar{t}_1, \dots, \bar{t}_{\bar{N}_T} \in (0, T]$ be the points of the process \bar{N}_t with stochastic intensity $\bar{\lambda}(t)$. For each of the points $\bar{t}_k, k = 1, \dots, \bar{N}_T$, attach a mark $\alpha = 0, 1, \dots, d$ with probability $\lambda^\alpha(\bar{t}_k)/\bar{\lambda}(\bar{t}_k)$, respectively. Then the points with marks $\alpha \in \mathbf{I}$, provide a d -variate point process with stochastic intensities $\lambda^\alpha(t)$.

The simulation of a multivariate Hawkes process with exponential decays on a fixed interval is similar to the univariate case, with only one extra step

that decides which dimension an accepted point belongs to. By Lemma 4.1, given that a point is accepted at time s , it should be distributed to dimension α with probability $\lambda^\alpha(s)/\sum_{\alpha \in \mathbf{I}} \lambda^\alpha(s)$, for $\alpha \in \mathbf{I}$. The procedure for simulating a d -variate Hawkes process on the interval $[0, T]$ is summarized below, where \mathcal{T}^α represents the ordered set of accepted points in dimension α and n^α counts the number of points \mathcal{T}^α , for $\alpha \in \mathbf{I}$. As before, s is always the newest candidate point generated.

- (1) Set $\mathcal{T}^1 = \dots = \mathcal{T}^d = \emptyset$, $s = 0$ and $n^1 = \dots = n^d = 0$.
- (2) Repeat the following until $s > T$.
 - (i) Set $\bar{\lambda} = \sum_{\alpha \in \mathbf{I}} \lambda^\alpha(s) = \sum_{\alpha \in \mathbf{I}} \left(\nu_\alpha + \sum_{\beta \in \mathbf{I}} \sum_{\tau \in \mathcal{T}^\beta} c_{\alpha\beta} e^{-a_{\alpha\beta}(s-\tau)} \right)$.
 - (ii) Generate u from a uniform distribution on $[0, 1]$.
 - (iii) Generate $w = -\log u/\bar{\lambda}$ as the interarrival to the next candidate point.
 - (iv) Set the new candidate point $s = s + w$.
 - (v) Generate D from a uniform distribution on $[0, 1]$.
 - (a) If $D \leq \sum_{\alpha \in \mathbf{I}} \lambda^\alpha(s)/\bar{\lambda}$, then do the following:
 - (a1) Find $\beta \in \mathbf{I}$ such that $\sum_{\alpha=1}^{\beta-1} \lambda^\alpha(s) < D\bar{\lambda} \leq \sum_{\alpha=1}^\beta \lambda^\alpha(s)$.
 - (a2) Assign candidate point s to dimension β by setting $n^\beta = n^\beta + 1$, $t_{n^\beta}^\beta = s$ and $\mathcal{T}^\beta = \mathcal{T}^\beta \cup \{t_{n^\beta}^\beta\}$.
 - (b) else, do nothing.
- (3) If $t_{n^\beta}^\beta > T$, then $\mathcal{T}^1, \dots, \mathcal{T}^\beta - \{t_{n^\beta}^\beta\}, \dots, \mathcal{T}^d$ contain the simulated points for each dimension;
- (4) else $\mathcal{T}^\alpha, \alpha \in \mathbf{I}$ contain the simulated points.

Figures 4.5-4.7 show an example sample path of intensity $\lambda^1(t)$, $\lambda^2(t)$ and $\lambda^3(t)$ and associated counting process $N = (N^\alpha)_{\alpha \in \mathbf{I}}$.

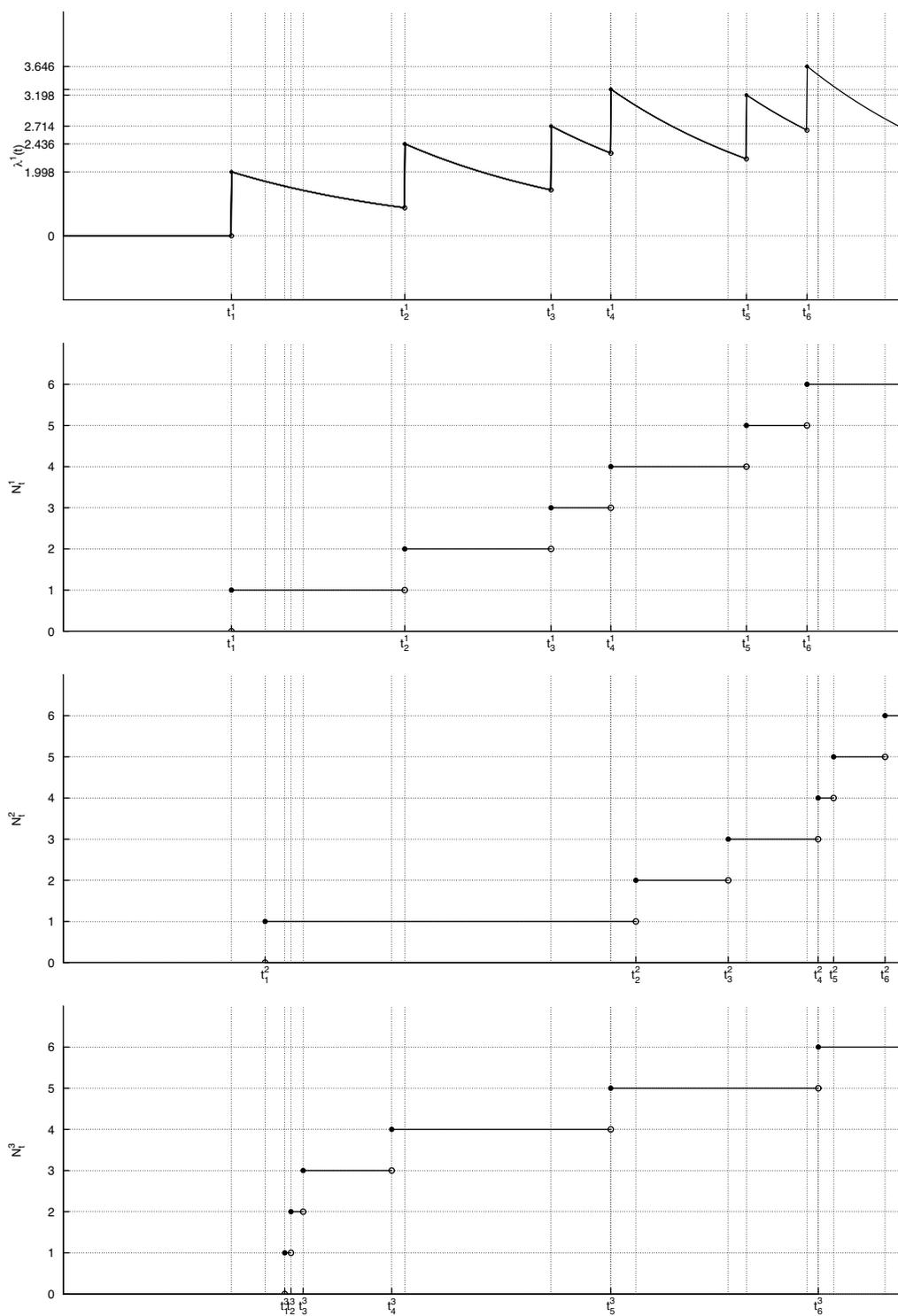


Figure 4.5: Sample paths of intensity $\lambda^1(t)$ and counting processes N

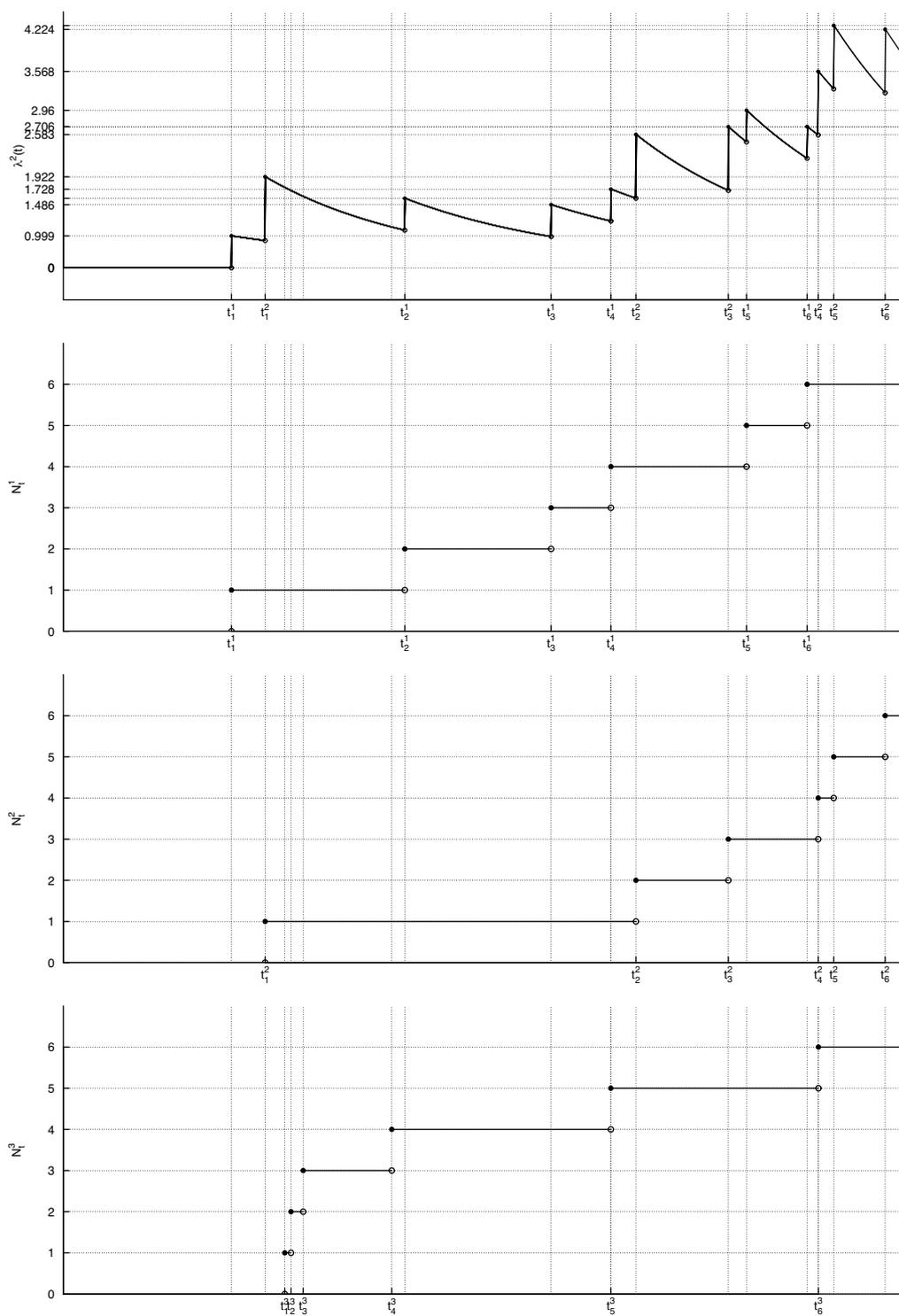


Figure 4.6: Sample paths of intensity $\lambda^2(t)$ and counting processes N

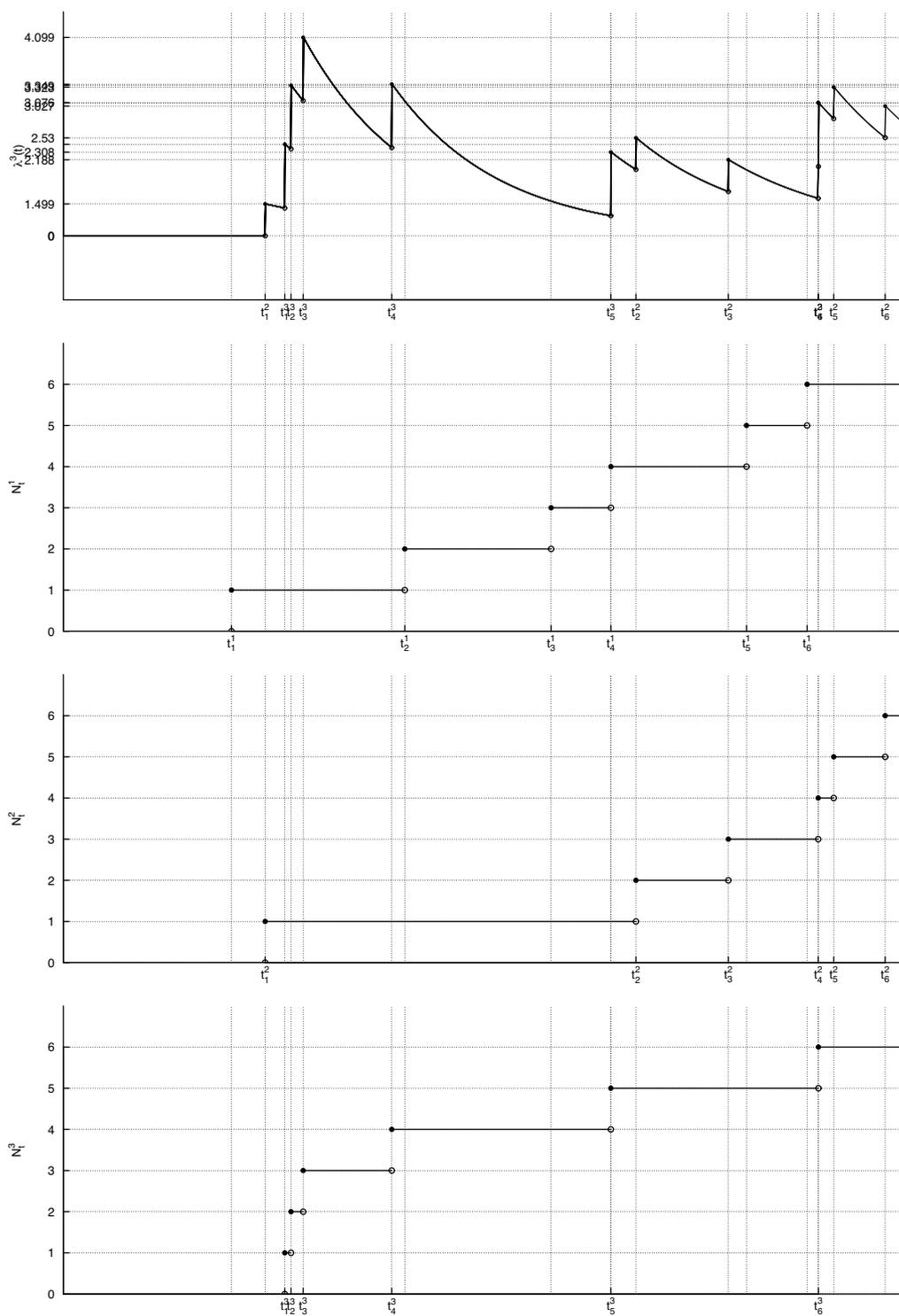


Figure 4.7: Sample paths of intensity $\lambda^3(t)$ and counting processes N

Table 4.3. normal methods (ν , QMLE) ; $T = 100, 200, 400, 600$.

	ν_1	ν_2	ν_3
true	1	0.5	1
$T = 100$	0.9404 (0.2638)	0.4581 (0.2098)	0.9375 (0.2438)
$T = 200$	0.9622 (0.1748)	0.4806 (0.1547)	0.9828 (0.1710)
$T = 400$	0.9785 (0.1252)	0.5008 (0.1013)	1.0037 (0.1132)
$T = 600$	0.9837 (0.1040)	0.5074 (0.0808)	1.0122 (0.0901)

Table 4.4. normal methods (C and A , QMLE) ; $T = 100, 200, 400, 600$.

	C_{11}	C_{12}	C_{13}	C_{21}	C_{22}	C_{23}	C_{31}	C_{32}	C_{33}
true	1	0	0	0.5	1	0	0	0.5	1
$T = 100$	0.9869 (0.3206)	0.1437 (0.3375)	0.1174 (0.2509)	0.5395 (0.3024)	1.0300 (0.4968)	0.1003 (0.2224)	0.0995 (0.2369)	0.5598 (0.4454)	1.0266 (0.5151)
$T = 200$	0.9737 (0.2131)	0.0871 (0.1781)	0.0625 (0.1480)	0.4984 (0.1783)	0.9841 (0.3004)	0.0522 (0.1126)	0.0558 (0.1406)	0.5248 (0.3173)	0.9848 (0.3592)
$T = 400$	0.9736 (0.1493)	0.0467 (0.0962)	0.0414 (0.0862)	0.4760 (0.1136)	0.9609 (0.2003)	0.0294 (0.0723)	0.0279 (0.0711)	0.4976 (0.2069)	0.9515 (0.2300)
$T = 600$	0.9721 (0.1170)	0.0305 (0.0670)	0.0300 (0.0644)	0.4725 (0.0954)	0.9557 (0.1628)	0.0194 (0.0526)	0.0202 (0.0507)	0.4934 (0.1739)	0.9479 (0.1802)
	A_{11}	A_{12}	A_{13}	A_{21}	A_{22}	A_{23}	A_{31}	A_{32}	A_{33}
true	2	***	***	2	3	***	***	3	4
$T = 100$	2.3240 (1.2544)	3.4120 (3.7475)	3.2599 (3.8204)	2.3013 (2.5430)	3.9122 (3.0128)	3.3814 (3.4455)	3.4714 (4.0253)	4.0366 (4.3408)	4.9993 (3.659)
$T = 200$	2.0683 (0.5400)	3.1594 (2.9022)	2.8740 (2.3451)	2.0451 (1.1171)	3.2604 (1.3562)	2.9903 (1.9802)	3.1011 (2.5508)	3.7092 (3.0898)	4.3864 (2.5230)
$T = 400$	2.0024 (0.3594)	2.8727 (1.8543)	2.8210 (1.7061)	1.9237 (0.5999)	3.0091 (0.8115)	2.8593 (1.7341)	2.8526 (1.7710)	3.3250 (1.9000)	3.9466 (1.3386)
$T = 600$	1.9713 (0.2776)	2.7551 (1.6391)	2.6996 (1.4824)	1.8992 (0.4821)	2.9361 (0.6134)	2.7922 (1.5474)	2.8203 (1.4925)	3.2191 (1.5721)	3.8592 (0.9821)

Table 4.5. The number of times, in percentage over 1000 iterations, that the estimator chooses the true model (coordinatewise); $T = 100, 200, 400, 600$.

true	T	Φ_{11} active	Φ_{12} 0	Φ_{13} 0	Φ_{21} active	Φ_{22} active	Φ_{23} 0	Φ_{31} 0	Φ_{32} active	Φ_{33} active	all
initial	100	100.0	54.4	54.6	100.0	99.9	55.9	55.9	97.8	99.8	9.4
	200	100.0	56.1	59.9	100.0	100.0	60.1	60.4	99.6	100.0	11.9
	400	100.0	58.7	59.7	100.0	100.0	64.2	64.2	100.0	100.0	12.9
	600	100.0	63.8	61.1	100.0	100.0	66.8	68.3	100.0	100.0	17.3
pLSA	100	99.7	78.6	81.3	93.0	99.1	82.6	81.7	80.1	97.0	31.0
	200	100.0	84.8	88.5	97.8	99.9	90.1	91.2	86.1	99.4	49.6
	400	100.0	90.8	91.3	99.7	100.0	94.1	95.3	93.4	100.0	69.5
	600	100.0	94.0	93.6	99.8	100.0	96.1	95.4	96.7	100.0	77.8

4.3 Simulation for the diffusion type process

We consider the model (3.9) in Section 3.2.2. Let $\mathbf{p} = 10$, i.e., the parameter space Θ is $[-10, 10]^{10}$. The process Y is defined by

$$Y_t = \int_0^t \sigma(X_s, \theta) dW_s, \quad t \in [0, 1],$$

where X is a 10-dimensional OU process satisfying the following stochastic differential equation

$$dX_t = -0.2X_t dt + 0.5I_{10}dw_t, \quad X_0 = 0, \quad t \in [0, 1]$$

and $\sigma(x, \theta) = \exp(\sum_{j=1}^{10} \theta_j x_j) \wedge M_0$, $M_0 = 10^5$. Here w is a 10-dimensional standard Wiener process independent of W . We generated the data

$(Y_{t_i}, X_{t_i})_{i=0,1,\dots,n}$, $t_i = \frac{i}{n}$ with

$$\theta^* = [1, 1, -1, -1, 0.5, 0, 0, 0, 0, 0]'$$

Let $\mathcal{L}_n(\theta) = \mathbb{L}_n(\theta) = -\mathbb{H}_n(\theta)$ in (3.10). We used the QMLE for the initial estimator $\tilde{\theta}$.

In order to apply our methods to this model, we use Theorem 6 in ([38]). By the definition of $\sigma(x, \theta)$, [A1] holds. [B2] is satisfied if we choose the stopping time $\tau \equiv 0$. Now we need to check [A3']. Since $\text{supp}\mathcal{L}\{X_0\} = \{0\}$, we can take $0 \in U \subset \{x \in \mathbb{R}^{10}; \sigma(x, \theta) < M, \forall \theta\}$. If we define $f(x, \theta) = (m_0 \sum_j (\theta_j - \theta_j^*) x_j) / |\theta - \theta^*|$ for sufficiently small m_0 when $\theta \neq 0$ and $f(x, 0) = \epsilon_0$ for some positive number ϵ_0 , then (i) is satisfied for $\rho = 2$. Next we take a covering $\{\Theta_k\}_{k=1,\dots,11}$ such that $\Theta_k = \{\theta \in \Theta - \{\theta^*\}; |\theta_k - \theta_k^*| \geq |\theta_j - \theta_j^*|, \forall j\}$ for $k = 1, \dots, 10$ and $\Theta_{11} = \{\theta^*\} \subset \Theta$. For $\Theta_k (k = 1, \dots, 10)$ if we take $\xi_0 = e_k$ and $\Psi(P_{\xi_0}^\perp x, \theta) = (\sum_{j \neq k} (\theta_j - \theta_j^*) x_j) / (\theta_k - \theta_k^*)$, then $|f(x, \theta)| \geq \frac{1}{\sqrt{10}} (\xi_0 \cdot x + \Psi(P_{\xi_0}^\perp x, \theta))$, and (ii) holds.

Then the objective function is denoted by

$$Q_n^{(q)}(\theta) = (\theta - \tilde{\theta})'(\theta - \tilde{\theta}) + \sum_{j=1}^{10} \kappa_n^j |\theta_j|^q$$

where $\kappa_n^j = \alpha_n |\tilde{\theta}_j|^{-\gamma}$, $\alpha_n = (\frac{1}{\sqrt{n}})^r$, $1 < r < 2 - q - \gamma$. We considered the cases where $n = 2500, 50000, 10000, 20000$ and the triplet of tuning parameters $(\gamma, r, q) = (3.2, 1.2, 0.3)$. In the same way as Table 1, Table 3 compares the results of the variable selection of the penalized LSA estimator, the unified LASSO type estimator and the Bridge type estimator.

Table 4 compares the means and the standard deviations (parentheses) for the three estimators (initial estimator, penalized LSA estimator and P-O estimator) in the case $n = 10000$.

Table 4.6. Results for the variable selection under $n=2500, 5000, 10000, 20000$.

	(γ, r, q)	$n = 2500$	$n = 5000$	$n = 10000$	$n = 20000$
% (p-LSA)	(3.2, 1.2, 0.3)	64.8	86.3	97.8	99.9
% (unified LASSO)	(3.2, 1.2, 1)	54.0	77.2	93.8	98.6
% (Bridge type)	(0, 1, 0.3)	0.7	0.9	1.0	1.7

Table 4.7. The summary of results for the simulation under $n = 10000$.

	θ_1	θ_2	θ_3	θ_4	θ_5	θ_6	θ_7	θ_8	θ_9	θ_{10}
true	1	1	-1	-1	0.5	0	0	0	0	0
initial	1.0020 (0.0879)	0.9952 (0.0896)	-0.9981 (0.0835)	-1.0017 (0.0847)	0.5021 (0.0842)	-0.0040 (0.0885)	-0.0006 (0.0896)	-0.0027 (0.0849)	0.0046 (0.0839)	0.0019 (0.0896)
p-LSA	1.0013 (0.0881)	0.9945 (0.0899)	-0.9974 (0.0838)	-1.0010 (0.0849)	0.4856 (0.1097)	-0.0003 (0.0107)	0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)	0.0003 (0.0000)
P-O	1.0017 (0.0663)	0.9980 (0.0694)	-0.9960 (0.0731)	-0.9978 (0.0676)	0.4936 (0.0916)	-0.0002 (0.0058)	0.0000 (0.0000)	0.0000 (0.0000)	0.0000 (0.0000)	0.0002 (0.0049)
%(p-LSA)	100.0	100.0	100.0	100.0	98.0	99.9	100.0	100.0	100.0	99.9

Bibliography

- [1] Akritas, M.G., Johnson, R.A.: Asymptotic inference in lévy processes of the discontinuous type. *The Annals of Statistics* pp. 604–614 (1981)
- [2] Athreya, K.B., Keiding, N.: Estimation theory for continuous-time branching processes. *Sankhyā: The Indian Journal of Statistics, Series A* pp. 101–123 (1977)
- [3] Bacry, E., Delattre, S., Hoffmann, M., Muzy, J.F.: Modelling microstructure noise with mutually exciting point processes. *Quantitative Finance* **13**(1), 65–77 (2013)
- [4] Basawa, I.V., Scott, D.J.: Asymptotic optimal inference for non-ergodic models, vol. 17. Springer Science & Business Media (2012)
- [5] Bickel, P.J., Ritov, Y., Tsybakov, A.B., et al.: Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics* **37**(4), 1705–1732 (2009)
- [6] Billingsley, P.: Statistical inference for Markov processes, vol. 2. University of Chicago Press (1961)
- [7] Blundell, C., Beck, J., Heller, K.A.: Modelling reciprocating relationships with hawkes processes. In: *Advances in Neural Information Processing Systems*, pp. 2600–2608 (2012)
- [8] Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J., et al.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning* **3**(1), 1–122 (2011)
- [9] Bradic, J., Fan, J., Jiang, J.: Regularization for Cox ’ s proportional hazards model with np-dimensionality. *Annals of Statistics* **39**(6), 3092 (2011)

- [10] Brillinger, D.R.: Statistical inference for stationary point processes. In: Selected Works of David Brillinger, pp. 499–543. Springer (2012)
- [11] Brown, M.: Statistical analysis of non-homogeneous poisson processes. Stochastic point processes pp. 67–89 (1972)
- [12] Clinet, S., Yoshida, N.: Statistical inference for ergodic point processes and application to limit order book. Stochastic Processes and their Applications **127**(6), 1800–1839 (2017)
- [13] Crane, R., Sornette, D.: Robust dynamic classes revealed by measuring the response function of a social system. Proceedings of the National Academy of Sciences **105**(41), 15,649–15,653 (2008)
- [14] De Gregorio, A., Iacus, S.M.: Adaptive lasso-type estimation for multivariate diffusion processes. Econometric Theory **28**(4), 838–860 (2012)
- [15] Efron, B., Hastie, T., Johnstone, I., Tibshirani, R., et al.: Least angle regression. The Annals of Statistics **32**(2), 407–499 (2004)
- [16] Fan, J., Li, R.: Variable selection via nonconcave penalized likelihood and its oracle properties. Journal of the American Statistical Association **96**(456), 1348–1360 (2001)
- [17] Feigin, P.D.: A note on maximum likelihood estimation for simple branching processes. Australian Journal of Statistics **19**(2), 152–154 (1977)
- [18] Frank, L.E., Friedman, J.H.: A statistical view of some chemometrics regression tools. Technometrics **35**(2), 109–135 (1993)
- [19] Friedman, J., Hastie, T., Tibshirani, R.: Sparse inverse covariance estimation with the graphical lasso. Biostatistics **9**(3), 432–441 (2008)
- [20] Friedman, J., Hastie, T., Tibshirani, R.: Regularization paths for generalized linear models via coordinate descent. Journal of Statistical Software **33**(1), 1 (2010)
- [21] Hawkes, A.G.: Point spectra of some mutually exciting point processes. Journal of the Royal Statistical Society: Series B (Methodological) **33**(3), 438–443 (1971)
- [22] Hawkes, A.G.: Spectra of some self-exciting and mutually exciting point processes. Biometrika **58**(1), 83–90 (1971)

- [23] Huang, J., Sun, T., Ying, Z., Yu, Y., Zhang, C.H.: Oracle inequalities for the lasso in the cox model. *Annals of Statistics* **41**(3), 1142 (2013)
- [24] Jeganathan, P.: On the asymptotic theory of estimation when the limit of the log-likelihood ratios is mixed normal. *Sankhyā Ser. A* **44**(2), 173–212 (1982)
- [25] Kamatani, K., Uchida, M.: Hybrid multi-step estimators for stochastic differential equations based on sampled data. *Statistical Inference for Stochastic Processes* **18**(2), 177–204 (2015)
- [26] Knight, K., Fu, W.: Asymptotics for lasso-type estimators. *Annals of Statistics* pp. 1356–1378 (2000)
- [27] Luschgy, H.: Asymptotic inference for semimartingale models with singular parameter points. *Journal of Statistical Planning and Inference* **39**(2), 155–186 (1994)
- [28] Masuda, H., Shimizu, Y.: Moment convergence in regularized estimation under multiple and mixed-rates asymptotics. *Mathematical Methods of Statistics* **26**(2), 81–110 (2017)
- [29] Meinshausen, N., Bühlmann, P., et al.: High-dimensional graphs and variable selection with the lasso. *The Annals of Statistics* **34**(3), 1436–1462 (2006)
- [30] Mohler, G.O., Short, M.B., Brantingham, P.J., Schoenberg, F.P., Tita, G.E.: Self-exciting point process modeling of crime. *Journal of the American Statistical Association* **106**(493), 100–108 (2011)
- [31] Ogata, Y.: On lewis’ simulation method for point processes. *IEEE Transactions on Information Theory* **27**(1), 23–31 (1981)
- [32] Ogata, Y.: Seismicity analysis through point-process modeling: A review. In: *Seismicity Patterns, their Statistical Significance and Physical Meaning*, pp. 471–507. Springer (1999)
- [33] Park, M.Y., Hastie, T.: L1-regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **69**(4), 659–677 (2007)
- [34] Suzuki, T., Yoshida, N.: Penalized least squares approximation methods and their applications to stochastic processes. *Japanese Journal of Statistics and Data Science* pp. 1–29 (2020)

- [35] Tibshirani, R.: Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 267–288 (1996)
- [36] Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., Knight, K.: Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **67**(1), 91–108 (2005)
- [37] Uchida, M., Yoshida, N.: Adaptive estimation of an ergodic diffusion process based on sampled data. *Stochastic Processes and their Applications* **122**(8), 2885–2924 (2012)
- [38] Uchida, M., Yoshida, N.: Quasi likelihood analysis of volatility and nondegeneracy of statistical random field. *Stochastic Processes and their Applications* **123**(7), 2851–2876 (2013)
- [39] Uchida, M., Yoshida, N.: Adaptive bayes type estimators of ergodic diffusion processes from discrete observations. *Statistical Inference for Stochastic Processes* **17**(2), 181–219 (2014)
- [40] Vere-Jones, D.: Stochastic models for earthquake occurrence. *Journal of the Royal Statistical Society: Series B (Methodological)* **32**(1), 1–45 (1970)
- [41] Vere-Jones, D., Ozaki, T.: Some examples of statistical estimation applied to earthquake data. *Annals of the Institute of Statistical Mathematics* **34**(1), 189–207 (1982)
- [42] Wainwright, M.J.: Sharp thresholds for high-dimensional and noisy sparsity recovery using-constrained quadratic programming (lasso). *IEEE Transactions on Information Theory* **55**(5), 2183–2202 (2009)
- [43] Wang, H., Leng, C.: Unified lasso estimation by least squares approximation. *Journal of the American Statistical Association* **102**(479), 1039–1048 (2007)
- [44] Witten, D.M., Tibshirani, R., Hastie, T.: A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* **10**(3), 515–534 (2009)
- [45] Yang, S.H., Zha, H.: Mixture of mutually exciting processes for viral diffusion. In: *International Conference on Machine Learning*, pp. 1–9 (2013)

- [46] Yoshida, N.: Polynomial type large deviation inequalities and quasi-likelihood analysis for stochastic differential equations. *Annals of the Institute of Statistical Mathematics* **63**(3), 431–479 (2011)
- [47] Yuan, M., Lin, Y.: Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**(1), 49–67 (2006)
- [48] Yuan, M., Lin, Y.: Model selection and estimation in the gaussian graphical model. *Biometrika* **94**(1), 19–35 (2007)
- [49] Zhang, C.H., et al.: Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics* **38**(2), 894–942 (2010)
- [50] Zhao, P., Yu, B.: On model selection consistency of lasso. *Journal of Machine Learning Research* **7**(Nov), 2541–2563 (2006)
- [51] Zhou, K., Zha, H., Song, L.: Learning social infectivity in sparse low-rank networks using multi-dimensional hawkes processes. In: *Artificial Intelligence and Statistics*, pp. 641–649 (2013)
- [52] Zou, H.: The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* **101**(476), 1418–1429 (2006)
- [53] Zou, H., Hastie, T., Tibshirani, R.: Sparse principal component analysis. *Journal of Computational and Graphical Statistics* **15**(2), 265–286 (2006)

appendix

A1. Polynomial type large deviation for sparse estimation

Polynomial type large deviation inequality

Let Θ be a bounded open set in \mathbb{R}^d . The closure of Θ in \mathbb{R}^d is denoted by $\bar{\Theta}$. Our interest will be on inference for the parameter θ in Θ . Given a probability space (Ω, \mathcal{F}, P) , we consider a sequence of random fields $\mathbb{H}_T : \Omega \times \bar{\Theta} \rightarrow \mathbb{R}$, where \mathbb{T} is a subset in $(0, \infty)$ such that $\sup \mathbb{T} = \infty$; e.g. \mathbb{T} is $\mathbb{R}_+ = [0, \infty)$, $\mathbb{Z}_+ = \{0, 1, 2, \dots\}$, $\mathbb{N} = \{1, 2, \dots\}$ etc. We suppose that the mapping $\bar{\Theta} \ni \theta \mapsto \mathbb{H}_T(\omega, \theta) \in \mathbb{R}$ is continuous for all $\omega \in \Omega$.

Let $\theta^* \in \bar{\Theta}$. Let $a_T \in GL(\mathfrak{p}, \mathbb{R})$ ($T \in \mathbb{T}$) such that $\lim_{T \rightarrow \infty} |a_T| = 0$, where $|a_T| = \{\text{Tr}(a_T a_T')\}^{1/2}$. Let $\mathbb{U}_T = \{u \in \mathbb{R}^d; \theta^* + a_T u \in \bar{\Theta}\}$. Define a statistical random field $\mathbb{Z}_T : \Omega \times \mathbb{U}_T \rightarrow (0, \infty)$ by

$$\mathbb{Z}_T(u) = \exp\left\{\mathbb{H}_T(\theta_T^\dagger(u)) - \mathbb{H}_T(\theta^*)\right\} \quad (u \in \mathbb{U}_T),$$

where $\theta_T^\dagger(u) = \theta^* + a_T u$.

Some notation will be necessary. We write

$$T[u_1, \dots, u_k] = T[u_1 \otimes \dots \otimes u_k] = \sum_{i_1, \dots, i_k} T_{i_1, \dots, i_k} u_1^{i_1} \dots u_k^{i_k}$$

for a tensor-valued tensor $T = (T_{i_1, \dots, i_k})_{i_1, \dots, i_k}$ and $u_1 = (u_1^{i_k})_{i_k}$. The brackets $[\]$ will often be used for multilinear mappings. We simply denote the r times product of u by $u^{\otimes r} = u \otimes \dots \otimes u$.

It is useful to consider the situation where \mathbb{Z} is locally asymptotically quadratic (LAQ). Let Δ_T be a d -dimensional random variable and let Γ be a $d \times d$ nonnegative-definite symmetric random matrix. We defined the random

field $r_T : \Omega \times \mathbb{U}_T \rightarrow \mathbb{R}$ by

$$\mathbb{Z}_T(u) = \exp\left(\Delta_T[u] - \frac{1}{2}\Gamma[u^{\otimes 2}] + r_T(u)\right) \quad (u \in \mathbb{U}_T).$$

It should be remarked that $0 \in \mathbb{R}^d$ is in \mathbb{U}_T but any d -dimensional ball centered at 0 may not be included in \mathbb{U}_T even for large T . The restriction on \mathbb{U}_T is common if u is banned from moving in all directions. Even in such a case, the LAN property makes sense if the random field \mathbb{H}_T is smooth in θ in restricted direction, as we will later see in application.

Though θ^* may be in Θ , we are interested in the case where $\theta^* \in \partial\Theta$, since the former case was treated in Yoshida [46]. The aim of this note is to remark that the same polynomial type large deviation as [46] holds, and to illustrate an application in the context of sparse estimation. In what follows, for simplicity, we will treat a single probability measure P , while a uniform estimate in ξ is considered in [46]. Besides, we will not treat the random field with a nuisance parameter τ in [46].

Let $b_T = \lambda_{\min}(a'_T a_T)^{-1}$, where λ_{\min} denotes the minimum eigenvalue of a real symmetric matrix. Since $|a_T| \rightarrow 0, 0 < b_T \rightarrow \infty$ as $T \rightarrow \infty$. We assume that

$$\sup_T \{b_T \lambda_{\max}(a'_T a_T)\} < \infty,$$

where λ_{\max} is the maximum eigenvalue of a real symmetric matrix.

Let

$$\mathbb{Y}_T(\theta) = b_T^{-1} \{\mathbb{H}_T(\theta) - \mathbb{H}_T(\theta^*)\}.$$

Fix a constant $\alpha \in (0, 1)$. Let

$$\mathbb{U}_T(r) = \{u \in \mathbb{U}_T; r \leq |u| \leq b_T^{(1-\alpha)/2}\}.$$

Let ρ, ρ_1, ρ_2 and β_2 be real constants. Let $L > 0$. The following conditions were assumed in [46].

[A1] There exists a constant C_L such that

$$\sup_{T \in \mathbb{T}} P \left[\sup_{u \in \mathbb{U}_T(r)} (1 + |u|^2)^{-1} |r_T(u)| \geq r^{-\rho_1} \right] \leq C_L r^{-L} \quad (r > 0)$$

The supremum on the empty set reads $-\infty$ by convention.

[A2] There exists a constant C_L such that

$$P[\lambda_{\min}(\Gamma) < 4r^{-\rho_1}] \leq C_L r^{-L} \quad (r > 0)$$

Suppose that a random field $\mathbb{Y} : \Omega \times \bar{\Theta} \rightarrow \mathbb{R}$ is given.

[A3] There exists a positive random variable χ_0 such that

$$\mathbb{Y}(\theta) = \mathbb{Y}(\theta) - \mathbb{Y}(\theta^*) \leq -\chi_0 |\theta - \theta^*|^\rho \quad (\theta \in \bar{\Theta}).$$

[A4] The following inequalities hold:

$$0 < \rho_1 < 1, \quad \alpha\rho < \rho_2, \quad \beta_2 \geq 0, \quad 1 - 2\beta_2 - \rho_2 > 0.$$

[A5] There exists a constant C_L such that

$$P[\chi_0 \leq r^{-(\rho_2 - \alpha\rho)}] \leq C_L r^{-L} (r > 0)$$

[A6] $\sup_{T \in \mathbb{T}} E[|\Delta_T|^{M_1}] < \infty$ for $M_1 = L(1 - \rho_1)^{-1}$. Moreover,

$$\sup_{T \in \mathbb{T}} E \left[\left\{ \sup_{\theta \in \bar{\Theta}, |\theta - \theta^*| \geq b_T^{-\alpha/2}} b_T^{\frac{1}{2} - \beta_2} |\mathbb{Y}_T(\theta) - \mathbb{Y}(\theta)| \right\}^{M_2} \right] < \infty$$

for $M_2 = L(1 - 2\beta_2 - \rho_2)^{-1}$.

Even in the case where $\theta^* \in \partial\Theta$, the same proof as Theorem 1 of [46] can be applied to the following theorem.

Theorem 4.2. Suppose that Conditions [A1]-[A6] are satisfied for a given positive number L . Then there exists a constant C_L such that

$$P \left[\sup_{u \in \mathbb{V}_T(r)} \mathbb{Z}_T(u) \geq \exp\left(-\frac{1}{2} r^{2 - (\rho_1 \vee \rho_2)}\right) \right] \leq C_L r^{-L} \quad (r > 0, T \in \mathbb{T}), \quad (4.1)$$

where $\mathbb{V}_T(r) = \{u \in \mathbb{U}_T; |u| \geq r\}$.

We will consider simplified version of Theorem 2.1. Among many possibilities, we only treat \mathbb{H}_T of class C^3 . In what follows, we only consider a bounded convex open set Θ . An extension of the results in this appendix to a finite sum of bounded convex open sets is straightforward. Let $C^k(\bar{\Theta})$ be the space of functions $f : \bar{\Theta} \rightarrow \mathbb{R}$ such that $f|_\Theta \in C^k(\Theta)$

and the derivatives $\partial_\theta^i f \in C(\Theta; (\mathbb{R}^d)^{\otimes i})$ are continuously extended to $\partial\Theta$ for $i \in \{0, 1, \dots, k\}$, $(\mathbb{R}^d)^{\otimes 0} = \mathbb{R}$. For $f \in C^k(\bar{\Theta})$, Taylor's formula is valid:

$$f(\theta) = \sum_{i=0}^{k-1} \frac{1}{k!} \partial_\theta^i f(\theta_0) [(\theta - \theta_0)^{\otimes i}] \\ + \frac{1}{(k-1)!} \int_0^1 (1-s)^{k-1} \partial_\theta^k f(\theta_0 + s(\theta - \theta_0)) ds [(\theta - \theta_0)^{\otimes k}]$$

for $\theta, \theta_0 \in \bar{\Theta}$ since $\bar{\Theta}$ is convex and the derivatives $\partial_\theta^i f$ are continuously extended to $\partial\Theta$.

We will assume that the mapping $\mathbb{H}_T(\omega, \cdot) \in C^3(\bar{\Theta})$ for all T and $\omega \in \Omega$. Set $\beta = \alpha/(1-\alpha)$. We replace Condition [A4] by

[A4'] Parameters β_1, ρ_1, ρ_2 and β_2 satisfy the following inequalities: $0 < \beta_1 < 1/2, 0 < \rho_1 < \min\{1, \beta, 2\beta_1/(1-\alpha)\}, \alpha\rho < \rho_2, \beta_2 \geq 0$ and $1 - 2\beta_2 - \rho_2 > 0$.

Furthermore, we assume the following additional conditions:

[A1'] For $M_3 = L(\beta - \rho_1)^{-1}$,

$$\sup_{T \in \mathbb{T}} E \left[\left(b_T^{-1} \sup_{\theta \in \bar{\Theta}} |\partial_\theta^3 \mathbb{H}_T(\theta)| \right)^{M_3} \right] < \infty.$$

Moreover, for $M_4 = L(\frac{2\beta_1}{1-\alpha} - \rho_1)^{-1}$,

$$\sup_{T \in \mathbb{T}} E \left[b_T^{\beta_1} |\Gamma_T(\theta^*) - \Gamma| \right] < \infty,$$

where a random matrix $\Gamma_T(\theta)$ is defined by

$$\Gamma_T(\theta)[u, u] = -\partial_\theta^2 \mathbb{H}_T(\theta)[a_T u, a_T u]$$

for $u \in \mathbb{R}^d$.

Now it is easy to obtain a counterpart of Theorems 2 of [46].

Theorem 4.3. Let $L > 0$. Suppose that Conditions [A1'], [A2], [A3], [A4'], [A5] and [A6] are satisfied. Then there exists a constant C_L such that Inequality (4.1) holds for all $T > 0$ and $r > 0$.

Restricted tangents

For limit theorems, for simplicity, we will only consider the situation where \mathbb{U}_T is asymptotically locally similar in the following sense. Let \mathbb{S} be a subset of \mathbb{R}^d .

[S1] For every $R > 0$, there exists $T_R \in \mathbb{T}$ such that

$$\{u \in \mathbb{U}_T; |u| \leq R\} = \mathbb{S} \cap \{u \in \mathbb{R}^d; |u| \leq R\}$$

for all $T \in \mathbb{T}$ with $T \geq T_R$.

Write $\overline{B(0, R)} = \{u \in \mathbb{R}^d; |u| \leq R\}$. Assume [S1]. Define $\mathbb{S}(R)$ by

$$\mathbb{S}(R) = \{u \in \mathbb{U}_{T_R}; |u| \leq R\} = \mathbb{S} \cup \overline{B(0, R)}.$$

Then

$$\mathbb{S} = \bigcup_{R>0} \mathbb{S}(R) = \bigcup_{R>0} \{u \in \mathbb{U}_{T_R}; |u| \leq R\} = \bigcup_{R>0} \bigcup_{T \geq T_R} \{u \in \mathbb{U}_T; |u| \leq R\}.$$

Since $\mathbb{S}(R) = \mathbb{U}_{T_R} \cap \overline{B(0, R)}$ is convex and increasing in R , the set \mathbb{S} is convex. BY definition, $0 \in \mathbb{S}$. If $\theta^* \in \Theta$, then $\mathbb{S} = \mathbb{R}^d$. Conditions [S1] imposes restrictions on the shape of Θ and a_T when $\theta^* \in \partial\Theta$.

Typically, the random field \mathbb{Z}_T is locally asymptotically mixed normal. To simplify description, we set $r_T(u) = 0$ for $u \in \mathbb{R}^d - \mathbb{U}_T$. Let \mathcal{G} be a sum σ -field of \mathcal{F} such that Γ is \mathcal{G} -measurable. Let Δ be a d -dimensional random variable defined on an extension $(\bar{\Omega}, \bar{\mathcal{F}}, \bar{P})$ of (Ω, \mathcal{F}, P) . Random variables on Ω are naturally extended on $\bar{\Omega}$. The \mathcal{G} -stable convergence is denoted by $d_s(\mathcal{G})$. Let

$$\mathbb{Z}(u) = \exp\left(\Delta[u] - \frac{1}{2}\Gamma[u^{\otimes 2}]\right) \quad (u \in \mathbb{R}^d).$$

We denote $\mathbb{C}(R) = C(\mathbb{S}(R))$ for $R > 0$ and equip $\mathbb{C}(R)$ with the supremum norm. Then $(\mathbb{C}(R), \mathbb{B}[\mathbb{C}(R)])$ is a measurable space with the Borel σ -field $\mathbb{B}[\mathbb{C}(R)]$.

Proposition 4.4. Suppose that [S1] and the following conditions are satisfied.

(i) For every $R > 0$, $\sup_{u \in \mathbb{S}(R)} |r_T(u)| \xrightarrow{p} 0$ as $T \rightarrow \infty$.

(ii) $\Delta_T \xrightarrow{d_s(\mathcal{G})} \Delta$ as $T \rightarrow \infty$.

Then

$$\mathbb{Z}_T|_{\mathbb{C}(R)} \xrightarrow{d_s(\mathcal{G})} \mathbb{Z}|_{\mathbb{C}(R)} \quad (T \rightarrow \infty)$$

for every $R > 0$.

quasi-maximum likelihood estimator

A measurable mapping $\tilde{\theta}_T^M : \Omega \rightarrow \bar{\Theta}$ is called a quasi-maximum likelihood estimator (QMLE) with respect to \mathbb{H}_T if

$$\mathbb{H}_T(\tilde{\theta}_T^M) = \max_{\theta \in \bar{\Theta}} \mathbb{H}_T(\theta).$$

Such a mapping exists according to the measurable selection theorem.

If Γ is positive definite a.s., then by convexity, \mathbb{Z}_T takes its maximum in \mathbb{S} at a single point $\tilde{u}_\mathbb{S}^M \in \mathbb{S}$ a.s.:

$$\tilde{u}_\mathbb{S}^M = \operatorname{argmax}_{u \in \mathbb{S}} \mathbb{Z}(u).$$

Let $\tilde{u}_T^M = a_T^{-1}(\tilde{\theta}_T^M - \theta^*)$.

Proposition 4.5. Suppose that [S1] and the following conditions are fulfilled.

(i) $\mathbb{Z}_T|_{\mathbb{C}(R)} \rightarrow^{d_s(\cdot)} \mathbb{Z}|_{\mathbb{C}(R)}$ as $T \rightarrow \infty$ for every $R > 0$.

(ii) $\{\tilde{u}_T^M\}_T$ is tight.

Then $\tilde{u}_T^M \rightarrow^{d_s(\mathcal{F})} \tilde{u}_\mathbb{S}^M$ as $T \rightarrow \infty$.

Combining Theorem 3.2 with the PLD inequality, we obtain convergence of moments of \tilde{u}_T^M as well as its tightness.

Theorem 4.6. Let $L > p > 0$. Suppose that Condition [S1] and the following conditions are fulfilled.

(i) $\Delta_T \rightarrow^{d_s(\mathcal{G})} \Delta$ as $T \rightarrow \infty$.

(ii) For every $R > 0$, $\sup_{u \in \mathbb{S}} |r_T(u)| \rightarrow^p 0$ as $T \rightarrow \infty$.

(iii) There exists a constant C_L such that

$$P \left[\sup_{u \in \mathbb{V}_T(r)} \mathbb{Z}_T(u) \geq 1 \right] \geq C_L r^{-L} \quad (r > 0, T \in \mathbb{T})$$

Then

$$E[f(\tilde{u}_T^M)Y] \rightarrow E[f(\tilde{u}_\mathbb{S}^M)] \quad (T \rightarrow \infty).$$

for every continuous function $f : \mathbb{S} \rightarrow \mathbb{R}$ satisfying $\sup_{u \in \mathbb{S}} (1 + |u|)^{-p} |f(u)| < \infty$ and every bounded \mathcal{G} -measurable random variable Y .

Proof. We apply Propositions 3.1 and 3.2. Condition (iii) gives L^p boundedness of $\{\tilde{u}_T^M\}_{T \in \mathbb{T}}$ for $p < L$. \square