# Weierstraß-Institut
## für Angewandte Analysis und Stochastik

# Parameter estimation in time series analysis

Vladimir Spokoiny

Weierstrass Institute
for Applied Analysis and Stochastics
Mohrenstr. 39
10117 Berlin Germany
E-Mail: spokoiny@wias-berlin.de

**Abstract**

The paper offers a novel unified approach to studying the accuracy of parameter estimation for a time series. Important features of the approach are: (1) The underlying model is not assumed to be parametric. (2) The imposed conditions on the model are very mild and can be easily checked in specific applications. (3) The considered time series need not to be ergodic or stationary. The approach is equally applicable to ergodic, unit root and explosive cases. (4) The parameter set can be unbounded and non-compact. (5) No conditions on parameter identifiability are required. (6) The established risk bounds are nonasymptotic and valid for large, moderate and small samples. (7) The results describe confidence and concentration sets rather than the accuracy of point estimation. The whole approach can be viewed as complementary to the classical one based on the asymptotic expansion of the log-likelihood. In particular, it claims a consistency of the considered estimate in a rather general sense, which usually is assumed to be fulfilled in the asymptotic analysis. In standard situations under ergodicity conditions, the usual rate results can be easily obtained as corollaries from the established risk bounds. The approach and the results are illustrated on a number of popular time series models including autoregressive, Generalized Linear time series, ARCH and GARCH models and meadian/quantile regression.

# 1 Introduction

Estimation of parameters of a time series is one of the most popular statistical problems which is included as an important building block in almost any econometric analysis. It is well known that statistical inference for time series is much more involved than the similar problem in the i.i.d. or regression set-up. The established results require quite strong conditions which are rather difficult to check in particular applications; see e.g. Brockwell and Davis (1991), Fan and Yao (2003). The aim of this paper is to offer a rather general and unified approach to measuring the quality in statistical estimation problem for time series which delivers meaningful and informative results under mild assumptions. We focus on the parametric modeling. It is however worth noting that any parametric assumption is only an approximation of reality and it is not precisely fulfilled in many particular situations. One can say that, in long run, any

fixed parametric specification is not flexible enough to describe the real structure of the data. The presented approach continues to apply in the cases when the underlying model *does not follow the parametric specification.* In some sense this approach refocuses the statistical paradigm: in many situations it might be useful and reasonable to apply a misspecified model with nice geometric properties rather than trying to precisely mimic the underlying model specifications. Typical examples of this pragmatic procedure are given by least squares, least absolute deviation or quintile regression: all of them can be viewed as quasi maximum likelihood estimates with a specific parametric structure.

One more nice feature of the presented approach is that the model assumptions are very general and non-restrictive and can be easily checked in specific applications. In particular, there is no any identifiability requirements, the results apply even if the parameter of the model is not identifiable. The parameter set can be unbounded and non-compact. Also no conditions like mixing, ergodicity, stationarity etc. are required: the observed time series can be *non-stationary*, *non-ergodic*, *non-mixing*, etc. The approach equally applies to ergodic, unit root and explosive time series. This enables, for instance, to analyze the quality of estimation and testing procedures for the unit root or cointegration analysis in a unified way; cf. Brockwell and Davis (1991), Johansen (1995), Johansen (2002). The required conditions are very mild and can be easily checked in particular applications.

The established risk bounds are *nonasymptotic* can be used for large, moderate and small samples. The results describe nonasymptotic *confidence and concentration sets* rather than the accuracy of point estimation. In the most of examples, the usual consistency and rate results can be easily obtained as corollaries from the established risk bounds.

The obtained exponential bound have been already used in various econometric studies. Spokoiny (2007) offered a local change point volatility estimation method, Čížek et al. (2007) discussed the estimation problem for varying coefficient ARCH and GARCH models, while Giacomini et al. (2007) focused on time varying copulae, Chen and Spokoiny (2007) considered the problem of robust risk management for non-normal and non-stationary market using stagewise aggregation procedure. All these and many other procedures are based on the multiple model check. The crucial issue in practical applications of such methods is the choice of the related parameters like thresholds or critical values in a data-driven way. This choice as well as the related theoretical analysis require to bound from above the probability of a wrong choice which can be done by the results presented below.

The paper is organized as follows. The next section describes the considered time series framework. Particularly, possible violations of the parametric assumption are discussed.

Section 3 presents the main results in form of penalized exponential bounds on the (quasi) maximum likelihood. Section 3.2 demonstrates some implications of the obtained results. We especially focus on the concentration properties of the estimates and on the likelihood based confidence sets. Section 4 illustrates how the general results can be specified for a number of popular time series models like Generalized Linear time series regression, linear autoregression, median and quantile regression. The main result given in Theorem 3.3 is obtained as a specification of general penalized exponential bound for the maximum of a random field from Section 5.

## 2  Modeling approach

This section describes the considered model and the modeling approach. Let the observed process $Y_t, t = 1, \ldots, n$ be progressively measurable w.r.t. a filtration $\mathcal{F} = (\mathcal{F}_t)$. Typically $\mathcal{F}_t$ stands for the information available at the moment $t$. One way of describing the joint distribution of the sample $\boldsymbol{Y}$ is by specifying the conditional distribution $Q_t = \mathcal{L}(Y_t | \mathcal{F}_{t-1})$ of every observation of $Y_t$ given the "past" $\mathcal{F}_{t-1}$. The parametric approach discussed below allows to reduce the whole description of the model to a few parameters which have to be estimated from the data.

### 2.1  A parametric model

The parametric time series modeling usually includes two important components, see e.g. Anderson (1994), Brockwell and Davis (1991), Kedem and Fokianos (2002), Fan and Yao (2003). One of them describes the type of conditional distribution $Q_t$ of $Y_t$ given the "past" $\mathcal{F}_{t-1}$ and the other one explains the dynamics of the corresponding parameter. The standard approach assumes that the conditional distribution $Q_t$ belongs to some given parametric family $\mathcal{P} = (P_v, v \in \mathcal{U})$, but the corresponding parameter $v$ may change in time and even be a random predictable process $f_t \sim \mathcal{F}_{t-1}$. We write this relation in the form

$$Q_t \stackrel{\text{def}}{=} \mathcal{L}(Y_t | \mathcal{F}_{t-1}) = P_{f_t} \in \mathcal{P}. \tag{2.1}$$

The second structural component of the parametric modeling concerns the driving (dynamic) process $f_t$. Namely, it is assumed that this process is uniquely described by a finite dimensional parameter $\boldsymbol{\theta}_0 \in \Theta \subseteq \mathbb{R}^p$, that is, $f_t = f_t(\boldsymbol{\theta}_0)$ for some $\boldsymbol{\theta}_0 \in \Theta$. These two assumptions lead to the parametric model in the form

$$Q_t = P_{f_t(\boldsymbol{\theta}_0)}. \tag{2.2}$$

Some typical examples of such parametric specifications are given in Section 2.3 and continued in Section 4.

For estimating $\boldsymbol{\theta}_0$, we apply the quasi maximum likelihood (quasi-MLE) approach. Let the family $\mathcal{P}$ be dominated by a measure $P_0$. Denote by $\ell(y,v)$ the corresponding log-density

$$\ell(y,v) = \log \frac{dP_v}{dP_0}(y).$$

The quasi log-likelihood $L(\boldsymbol{Y},\boldsymbol{\theta})$ for the model (2.1)–(2.2) can be represented in the form

$$L(\boldsymbol{Y},\boldsymbol{\theta}) = \sum_t \ell\big(Y_t, f_t(\boldsymbol{\theta})\big).$$

Here in in what follows, $\sum_t$ means summation over the whole time interval $t = 1, \ldots, n$. We define the quasi-MLE estimate $\widetilde{\boldsymbol{\theta}}$ of the parameter $\boldsymbol{\theta}$ by maximizing the log-likelihood $L(\boldsymbol{\theta})$:

$$\widetilde{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{Y},\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_t \ell\big(Y_t, f_t(\boldsymbol{\theta})\big), \tag{2.3}$$

and denote by $L(\boldsymbol{Y},\widetilde{\boldsymbol{\theta}})$ the corresponding maximum.

It is important to stress that the parametric assumption (2.2) is only an approximation of the underlying data distribution $\boldsymbol{P}$ which justifies the estimation procedure (2.3). In reality, the modeling assumption (2.2) can be violated in one or even both parts. One of the aims of our study is to address the questions of what is estimated and with which accuracy if the parametric assumption $Q_t = P_{f_t(\boldsymbol{\theta}_0)}$ is not precisely fulfilled.

## 2.2 Violation of the parametric assumption

The parametric model (2.2) can be violated by two different reasons. One is due to misspecified conditional distribution and the other one due to a wrong parametric dynamics $f_t = f_t(\boldsymbol{\theta})$.

### 2.2.1 Misspecified conditional distribution

The model (2.2) assumes that the conditional distribution $Q_t$ of $Y_t$ given $\mathcal{F}_{t-1}$ belongs to the given family $\mathcal{P}$. This assumption can be well justified for many examples from categorical data analysis, for instance, for binary or discrete observations; see Fokianos and Kedem (2003). However, this assumption could be too restrictive for many other applications. We present here a couple of examples of this sort. First consider a

stochastic dynamic system described by the equation $Y_t = f_t + \varepsilon_t$ whose drift $f_t$ is a predictable process and innovations $\varepsilon_t$ are martingale differences. In the case of conditionally standard normal innovations $\varepsilon_t$, the distribution $Q_t$ is normal with the mean $f_t$ leading to the log-density function $\ell(y, v) = -(y - v)^2/2$ (up to an unimportant constant term $-\frac{1}{2}\log(2\pi)$). Parametric dynamics $f_t = f_t(\boldsymbol{\theta})$ leads to the log-likelihood $L(\boldsymbol{\theta}) = -\sum_t |Y_t - f_t(\boldsymbol{\theta})|^2/2$. In the case of non-normal innovations, this expression is a quasi log-likelihood leading to the least squares solution.

Another typical example is given by the volatility modeling. The log-returns $r_t$ are described by the conditional heteroscedasticity model: $r_t = \sigma_t \varepsilon_t$. The case of standard Gaussian innovations $\varepsilon_t$ and the parametric dynamics $f_t = f_t(\boldsymbol{\theta})$ for the volatility $f_t = \sigma_t^2$ leads to the log-likelihood $L(\boldsymbol{\theta}) = -\frac{1}{2}\sum_t \{ r_t^2/f_t(\boldsymbol{\theta}) + \log(2\pi f_t(\boldsymbol{\theta})) \}$. In the case of, say, heavy tailed innovations, one can still try to maximize this expression which becomes a quasi log-likelihood.

### 2.2.2 Misspecified parametric dynamics and the best parametric fit

Suppose for a moment that the conditional distribution $Q_t$ belongs to the given family $\mathcal{P}$ almost surely for all $t$. Then the data $\boldsymbol{Y}$ can be described by the model $\mathcal{L}(Y_t | \mathcal{F}_{t-1}) = P_{f_t} \in \mathcal{P}$ for some predictable process $f_t$. The parametric assumption means that the process $f_t$ belongs to a parametric family of processes $(f_t(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p)$. This assumption is very useful for the analysis but it is usually only an idealization of reality. An interesting question in this respect is what is estimated in the situation when the process $f_t$ does not follow the parametric dynamics $f_t = f_t(\boldsymbol{\theta})$ whatever $\boldsymbol{\theta}$ is. Below we show that in such cases the quasi log-likelihood approach leads to an estimate of the projection of the given model on the parametric subspace of models. One also speaks about the best parametric fit $f_t(\boldsymbol{\theta})$ to the model $f_t$ which can be defined as a solution of the optimization problem

$$\boldsymbol{\theta}_0 \overset{\text{def}}{=} \operatorname*{argmax}_{\boldsymbol{\theta} \in \Theta} \boldsymbol{E}L(\boldsymbol{\theta}) = \operatorname*{argmin}_{\boldsymbol{\theta} \in \Theta} \mathcal{K}(\boldsymbol{P}, \boldsymbol{P_\theta}),$$

where $\boldsymbol{P} = \prod_t Q_t = \prod_t P_{f_t}$ is the true measure, $\boldsymbol{P_\theta} = \prod_t P_{f_t(\boldsymbol{\theta})}$ is its parametric counterpart, and $\mathcal{K}(\boldsymbol{P}, \boldsymbol{P}') \overset{\text{def}}{=} \boldsymbol{E}\log(d\boldsymbol{P}/d\boldsymbol{P}')$ is the Kullback-Leibler divergence between two measures $\boldsymbol{P}$ and $\boldsymbol{P}'$.

The interpretation of $\boldsymbol{\theta}_0$ as the "best parametric fit" continues to apply even in the case when also the assumption $Q_t \in \mathcal{P}$ a.s. is violated. However, $\boldsymbol{\theta}_0$ cannot be defined as a minimizer of the Kullback-Leibler divergence anymore.

## 2.3 Examples

This section presents some popular time series models. Later in Section 4 we illustrate the obtained results on these examples.

### 2.3.1 Linear autoregression

Let the family $\mathcal{P}$ be a Gaussian shift. This case corresponds to the model $Y_t = f_t + \varepsilon_t$ in which the innovations $\varepsilon_t$ are assumed to be i.i.d Gaussian: $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$. The corresponding (quasi) maximum likelihood approach leads back to the least square estimate $\widetilde{\boldsymbol{\theta}}$: with $L(\boldsymbol{\theta}) = -(2\sigma^2)^{-1} \sum_t \{Y_t - f_t(\boldsymbol{\theta})\}^2$

$$\widetilde{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname*{argmin}_{\boldsymbol{\theta}} \sum_t \{Y_t - f_t(\boldsymbol{\theta})\}^2.$$

Linear autoregression means the structural equation $f_t(\boldsymbol{\theta}) = \alpha_1 Y_{t-1} + \ldots + \alpha_p Y_{t-p}$ for $\boldsymbol{\theta} = (\alpha_1, \ldots, \alpha_p)^\top$ leading to maximization of the quadratic functional $\sum_t (Y_t - \Psi_t^\top \boldsymbol{\theta})^2$ with $\Psi_t = (Y_{t-1}, \ldots, Y_{t-p})^\top$ which admits a closed form solution:

$$\widetilde{\boldsymbol{\theta}} = \left( \sum_t \Psi_t \Psi_t^\top \right)^{-1} \sum_t Y_t \Psi_t = B^{-1} \sum_t Y_t \Psi_t$$

$$L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) = \frac{1}{2\sigma^2} (\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^\top B (\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0),$$

with $B \overset{\text{def}}{=} \sum_t \Psi_t \Psi_t^\top$. Note, however, that through the closed form solution for $\widetilde{\boldsymbol{\theta}}$ is available, the analysis in time series context remains a difficult task, especially in the non-ergodic case, cf. Dickey and Fuller (1981), Basawa and Brockwell (1984), Chan and Wei (1988), Fountis and Dickey (1989), Cox and Llatas (1991), Koul and Saleh (1993), Phillips and Xu (2006). All the mentioned papers studied the asymptotic properties of the estimates. There are very few papers concerning nonasymptotic results, see Chan and Wei (1988). Fixed accuracy sequential procedures are discussed in Lai and Siegmund (1983), Sriram (1987), Shiryaev and Spokoiny (1997), Konev and Pergamenshchikov (1997). More robust estimates like minimum distance, M- or quantile estimates have been considered in Wang (1986), Chan and Wei (1988).

Our approach based on the exponential bounds for the (quasi) likelihood continues to apply. Moreover, it does not assume that the innovations are independent or conditionally Gaussian. Neither we require that the underlying process $f_t$ follows the linear structural equation $f_t = \Psi_t^\top \boldsymbol{\theta}$.

### 2.3.2 GARCH(1,1) estimation

GARCH-modeling introduced in Bollerslev (1986) is very popular in analysis of financial time series. A number of GARCH extensions is proposed to make the model even more flexible; for example, EGARCH Nelson (1991), QGARCH Sentana (1995), among many others. We focus on the classical GARCH(1,1) model although most of conclusions can be extended to more general specifications. The underlying modeling assumption is that the observed squared log-returns $R_t$ follows the conditional heteroscedasticity equation:

$$R_t = X_t \varepsilon_t^2,$$

where $\varepsilon_t$ are standardized innovations satisfying $\boldsymbol{E}\big(\varepsilon_t \big| \mathcal{F}_{t-1}\big) = 0$, $\boldsymbol{E}\big(\varepsilon_t^2 \big| \mathcal{F}_{t-1}\big) = 1$, and $X_t$ is a predictable volatility process. The parametric GARCH(1,1) assumption means that the volatility process $X_t$ follows the equation

$$X_t = \omega + \alpha R_{t-1} + \beta X_{t-1}. \tag{2.4}$$

For simplicity we assume that the initial value $X_0$ is fixed, e.g. $X_0 = R_0$. Then for every $\boldsymbol{\theta} = (\omega, \alpha, \beta)^\top$ we can recursively apply the structural equation (2.4) yielding the process $X_t(\boldsymbol{\theta})$ with

$$X_t(\boldsymbol{\theta}) = \omega + \alpha R_{t-1} + \beta X_{t-1}(\boldsymbol{\theta}), \qquad t \geq 1, \qquad X_0(\boldsymbol{\theta}) = X_0.$$

With this process we associate a (quasi) log likelihood

$$L(\boldsymbol{\theta}) = -\frac{1}{2} \sum_{t=1}^{n} \big\{ \log(2\pi X_t(\boldsymbol{\theta})) + R_t / X_t(\boldsymbol{\theta}) \big\}.$$

This expression becomes the log-likelihood if the innovations $\varepsilon_t$ are conditionally on $\mathcal{F}_{t-1}$ standard normal and the structural equation (2.4) is fulfilled for some combination of parameters. Asymptotic properties of such estimates are well studied, see e.g. Lee and Hansen (1994), Fan and Yao (2003), Sun and Stengos (2006), Francq and Zakoian (2007), and references therein.

### 2.3.3 Median and quantile time series estimation

Median or more generally quantile estimation is known to be more robust and stable against outliers and it is frequently used in econometric studies; see Koenker (2005), Koenker and Xiao (2006). The corresponding approach explains the observations $Y_t$ by the regression equation $Y_t = f_t + \varepsilon_t$ where the individual errors $\varepsilon_t$ are not assumed to fulfill $\boldsymbol{E}(\varepsilon_t | \mathcal{F}_{t-1}) = 0$. Instead, one imposes the constraint $\boldsymbol{P}(\varepsilon_t > 0 | \mathcal{F}_{t-1}) = \alpha$ for a given $\alpha$. The median regression corresponds to $\alpha = 1/2$. Under the parametric

assumption $f_t = f_t(\boldsymbol{\theta}_0)$ for a given parametric class of predictable processes $f_t(\boldsymbol{\theta})$, the corresponding estimate can be defined by maximizing the quasi log-likelihood

$$L(\boldsymbol{\theta}) = \sum_t \ell_\alpha \big( Y_t - f_t(\boldsymbol{\theta}) \big)$$

with $\ell_\alpha(x) = (1-\alpha)x_- - \alpha x_+$.

In this case, $\mathcal{P}$ is the family with the log-density $\ell(y, v) = \ell_\alpha(y - v)$. In particular, the median regression for $\alpha = 1/2$ corresponds to the Laplacian shift family.

### 2.3.4 Categorical time series

Let $\mathcal{P}$ be an exponential family with the canonical parametrization (EFC) which means that the corresponding log-likelihood function can be written in the form

$$\ell(y, v) = yv - d(v) + \ell(y)$$

where $d(\cdot)$ is a given convex function; see e.g. McCullagh and Nelder (1989), Green and Silverman (1994). The term $\ell(y)$ is unimportant and it cancels in the log-likelihood ratio. Such families are often used in the categorical time series analysis for describing the conditional distribution $Q_t$ of the observed data; see Fokianos and Kedem (2003). The corresponding model can be written as

$$\mathcal{L}\big(Y_t | \mathcal{F}_{t-1}\big) = P_{f_t} \in \mathcal{P}. \tag{2.5}$$

Parametric modeling assumes a specific structure of the driving process $f_t$ leading to the parametric log-likelihood function $L(\boldsymbol{\theta}) = \sum_t \ell(Y_t, f_t(\boldsymbol{\theta}))$:

$$L(\boldsymbol{\theta}) = \sum_t \ell(Y_t, f_t(\boldsymbol{\theta})) = \sum_t \big\{ Y_t f_t(\boldsymbol{\theta}) - d\big(f_t(\boldsymbol{\theta})\big) \big\}. \tag{2.6}$$

Usually $f_t(\boldsymbol{\theta})$ can be represented in the form $f_t(\boldsymbol{\theta}) = m(X_t, \boldsymbol{\theta})$ for some *regression function* $m(\cdot, \cdot)$, an explanatory process $X_t$ and a parameter vector $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p$. One popular example of linear regression is discussed in the next section.

Our approach allows to account for the both kinds of model misspecification: $Q_t \notin \mathcal{P}$ and/or $f_t \neq f_t(\boldsymbol{\theta})$. However, to be more specific, we consider below the case when $Q_t \in \mathcal{P}$. Then the focus of analysis is the best parametric approximation of the true regression function $f_t$ by a parametric model $f_t(\boldsymbol{\theta})$.

### 2.3.5 Generalized Linear time series

Let $\mathcal{P}$ be again an EFC. A parametric *generalized linear* specification for the model (2.5) is given by the following set of structural equations:

$$\mathcal{L}\big(Y_t\big|\mathcal{F}_{t-1}\big) = P_{f_t} \in \mathcal{P}, \qquad f_t = g(\Psi_t), \qquad \Psi_t = A(\boldsymbol{\theta})\Psi_{t-1} \tag{2.7}$$

where $\Psi_t$ is a predictable $R^d$-dimensional explanatory process, $g(\cdot)$ is a given mapping from $I\!\!R^d$ to $I\!\!R$, $A(\boldsymbol{\theta})$ is a given $d \times d$-matrix linearly depending on the parameter vector $\boldsymbol{\theta} \in \Theta \subset I\!\!R^p$. Such models are widely used in statistical modeling. A popular example is given by the equations $f_t = g(X_t)$ and

$$X_t = \omega + \alpha_1 Y_{t-1} + \ldots \alpha_p Y_{t-p} + \beta_1 X_{t-1} + \ldots \beta_q X_{t-q}. \tag{2.8}$$

Here $\Psi_t = (X_t, \ldots, X_{t-q+1}, Y_t, \ldots, Y_{t-p+1}, 1)^\top$, $\boldsymbol{\theta} = (\beta_1, \ldots, \beta_q, \alpha_1, \ldots, \alpha_p, \omega)^\top$, and the first row of $A(\boldsymbol{\theta})$ is just $\boldsymbol{\theta}$. For $\beta_1 = \ldots = \beta_q = 0$ the value $X_t$ is a linear combination of the past observations and (2.7) becomes an autoregressive type model. If there is at least one coefficient $\beta_j \neq 0$, then $X_t$ is an unobservable (hidden/latent/exogeneous) component and (2.7) is of ARMA type; see e.g. Fokianos and Kedem (2003).

Let $\boldsymbol{\theta}$ be the parameter vector. Then, given $\boldsymbol{\theta}$, the observations $Y_1, \ldots, Y_t$, and the pre-history $\Psi^0, Y^0$, one can uniquely reconstruct the process $\Psi_t = \Psi_t(\boldsymbol{\theta})$ and then $f_t(\boldsymbol{\theta}) = g(\Psi_t(\boldsymbol{\theta}))$ for $t \geq 1$ by recurrently applying the relation $\Psi_t = A(\boldsymbol{\theta})\Psi_{t-1}$. The process $f_t(\boldsymbol{\theta})$ leads to the (quasi) log-likelihood $L(\boldsymbol{\theta}) = \sum_t \big\{Y_t f_t(\boldsymbol{\theta}) - d\big(f_t(\boldsymbol{\theta})\big)\big\}$. Inference for GLM's has been discussed in many papers and books. We only mention Green and Silverman (1994), Chen (1995), Chen et al. (1999), Sun et al. (2000), Fokianos and Kedem (2003), Kedem and Fokianos (2002), Fan and Yao (2003), among many others. Our analysis further in Section 4 essentially differs from all the mentioned studies. In particular, it does not assume that any of imposed parametric specifications from (2.7) is really fulfilled. The methods and results are non-asymptotic.

In all the examples, the true model is still given by (2.1). By $\boldsymbol{\theta}_0$ we denote the parameter corresponding to the best parametric fit: $\boldsymbol{\theta}_0 = \operatorname{argmax}_{\boldsymbol{\theta}} \boldsymbol{E}L(\boldsymbol{\theta})$. The parameter $\boldsymbol{\theta}_0$ is estimated by maximizing the objective function $L(\boldsymbol{\theta})$.

## 3 Exponential risk bounds

In this section we first introduce the basic notions and conditions on the model and then state the main results in form of general exponential bounds for the supremum of the quasi log-likelihood function $L(\boldsymbol{\theta})$. The quality of estimation of $\boldsymbol{\theta}$ is measured in terms

of the maximum $L(\widetilde{\boldsymbol{\theta}}) = \max_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta})$ rather than the point of maximum $\widetilde{\boldsymbol{\theta}}$, where $L(\boldsymbol{\theta})$ from (2.5). More precisely, we define the point

$$\boldsymbol{\theta}_0 \overset{\text{def}}{=} \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \boldsymbol{E} L(\boldsymbol{\theta})$$

which is the true value in the parametric situation and can be viewed as the parameter of the best parametric fit in the general case. Now the aim of our study is to establish some exponential bounds on the supremum in $\boldsymbol{\theta}$ of the random field

$$L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \overset{\text{def}}{=} L(\boldsymbol{\theta}) - L(\boldsymbol{\theta}_0) = \sum_t \big\{ \ell(Y_t, f_t(\boldsymbol{\theta})) - \ell(Y_t, f_t(\boldsymbol{\theta}_0)) \big\}.$$

Later in Section 3.2 we comment how the accuracy of estimation of $\boldsymbol{\theta}_0$ by $\widetilde{\boldsymbol{\theta}}$ relates to the value $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0)$. We will also see that the bound for $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0)$ yields the confidence sets for the parameter $\boldsymbol{\theta}_0$ and concentration sets for the estimate $\widetilde{\boldsymbol{\theta}}$.

Define for $\boldsymbol{\theta} \in \Theta$

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \overset{\text{def}}{=} \sum_t \mathfrak{m}_t\big(\mu, f_t(\boldsymbol{\theta}), f_t(\boldsymbol{\theta}_0)\big). \tag{3.1}$$

where for $\upsilon, \upsilon' \in \mathcal{U}$

$$\mathfrak{m}_t(\mu, \upsilon, \upsilon') \overset{\text{def}}{=} - \log \boldsymbol{E}\big[\exp\big\{\mu\ell(Y_t, \upsilon) - \mu\ell(Y_t, \upsilon')\big\}\big|\mathcal{F}_{t-1}\big].$$

This definition assumes the following condition:

**(E)**   *There exists some $\mu > 0$ such that for all $\boldsymbol{\theta} \in \Theta$ and all $t$ the value $\mathfrak{m}_t\big(\mu, f_t(\boldsymbol{\theta}), f_t(\boldsymbol{\theta}_0)\big)$ is finite.*

Note that this condition is automatically fulfilled with $\mu \leq 1$ if $\boldsymbol{P} = \boldsymbol{P}_{\boldsymbol{\theta}_0}$ and $L(\boldsymbol{\theta})$ is indeed a log-likelihood function.

The main observation behind the definition (3.1) is that

$$\boldsymbol{E} \exp\big\{\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)\big\} = 1.$$

Our main goal is to get an exponential bound for the maximum of the random field $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ over $\boldsymbol{\theta} \in \Theta$. Unfortunately, this maximum may explode and we consider the penalized expression $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \text{pen}(\boldsymbol{\theta})$ where the penalty function $\text{pen}(\boldsymbol{\theta})$ should provide some bounded exponential moments for

$$\sup_{\boldsymbol{\theta} \in \Theta}\big[\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \text{pen}(\boldsymbol{\theta})\big].$$

More precisely, we present a penalty function $\text{pen}(\boldsymbol{\theta})$ that ensures under rather general conditions for every $\varrho < 1$ that the value

$$\mathfrak{Q}(\varrho) \stackrel{\text{def}}{=} \boldsymbol{E} \sup_{\boldsymbol{\theta} \in \Theta} \exp\big\{\varrho\big[\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \text{pen}(\boldsymbol{\theta})\big]\big\}$$

is bounded by a fixed constant.

We consider the following decomposition of the log-likelihood process $L(\boldsymbol{\theta})$ into the martingale-difference and predictable parts:

$$L(\boldsymbol{\theta}) \stackrel{\text{def}}{=} M(\boldsymbol{\theta}) + \zeta(\boldsymbol{\theta})$$

where

$$M(\boldsymbol{\theta}) \quad \stackrel{\text{def}}{=} \quad \sum_t m_t\big(f_t(\boldsymbol{\theta})\big),$$

$$\zeta(\boldsymbol{\theta}) \quad \stackrel{\text{def}}{=} \quad \sum_t \zeta_t\big(Y_t, f_t(\boldsymbol{\theta})\big)$$

with $m_t(v) \stackrel{\text{def}}{=} \boldsymbol{E}\big[\ell(Y_t, v)\big|\mathcal{F}_{t-1}\big]$, $\zeta_t(Y_t, v) \stackrel{\text{def}}{=} \ell(Y_t, v) - m_t(v)$ for $v \in \mathcal{U}$. Below we assume that the (random) functions $m_t(v)$ and $\zeta_t(v)$ are differentiable w.r.t. $v$ and denote $\dot{m}_t(v) = dm_t(v)/dv$ and $\dot{\zeta}_t(y, v) = d\zeta_t(y, v)/dv$.

Suppose also that the random function $f_t(\boldsymbol{\theta})$ is differentiable in $\boldsymbol{\theta}$ and denote $\nabla f_t(\boldsymbol{\theta}) = \partial f_t(\boldsymbol{\theta})/\partial\boldsymbol{\theta} \in \mathbb{R}^p$. Define

$$\nabla M(\boldsymbol{\theta}) \quad \stackrel{\text{def}}{=} \quad \sum_t \dot{m}_t\big(f_t(\boldsymbol{\theta})\big)\nabla f_t(\boldsymbol{\theta})$$

$$\nabla \zeta(\boldsymbol{\theta}) \quad \stackrel{\text{def}}{=} \quad \sum_t \dot{\zeta}_t\big(Y_t, f_t(\boldsymbol{\theta})\big)\nabla f_t(\boldsymbol{\theta}).$$

Condition $(E)$ assumes that the quasi log-likelihood has bounded exponential moments. We also assume a similar property for its gradient.

**(ED)** *There exist some deterministic symmetric matrix $V(\boldsymbol{\theta})$ and a constant $\lambda^* > 0$ such that for all $\lambda \leq \lambda^*$*

$$\sup_{\boldsymbol{\gamma} \in \mathcal{S}^p} \sup_{\boldsymbol{\theta} \in \Theta} \log \boldsymbol{E} \exp\left\{2\lambda\frac{\boldsymbol{\gamma}^\top \nabla\zeta(\boldsymbol{\theta})}{\sqrt{\boldsymbol{\gamma}^\top V(\boldsymbol{\theta})\boldsymbol{\gamma}}}\right\} \leq 2\lambda^2, \tag{3.2}$$

*and*

$$\sup_{\boldsymbol{\gamma} \in \mathcal{S}^p} \sup_{\boldsymbol{\theta} \in \Theta} \boldsymbol{E} \exp\left\{2\lambda\frac{\boldsymbol{\gamma}^\top [\nabla M(\boldsymbol{\theta}) - \boldsymbol{E}\nabla M(\boldsymbol{\theta})]}{\sqrt{\boldsymbol{\gamma}^\top V(\boldsymbol{\theta})\boldsymbol{\gamma}}}\right\} \leq 2\lambda^2. \tag{3.3}$$

This condition is usually simple to check. Below we present some simple sufficient conditions for (3.2) and (3.3).

**Lemma 3.1.** *Suppose that there exist a constant $\lambda_1^* > 0$ and a random function $\mathfrak{n}_t(v) \sim \mathcal{F}_{t-1}$ such that for all $t$ and $\lambda \leq \lambda_1^*$*

$$\log \boldsymbol{E}\left[\exp\left\{2\lambda \frac{\dot{\zeta}_t(Y_t, v)}{\mathfrak{n}_t(v)}\right\} \middle| \mathcal{F}_{t-1}\right] \leq 2\lambda^2, \qquad v \in \mathcal{U}. \tag{3.4}$$

*Let also there exist a deterministic matrix function $V(\boldsymbol{\theta}) \geq I$ and the value $\lambda^* \geq \lambda_1^*$ such that it holds almost surely for any $\boldsymbol{\theta} \in \Theta$ and $\boldsymbol{\gamma} \in \mathcal{S}^p$*

$$B(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_t \mathfrak{n}_t^2(f_t(\boldsymbol{\theta})) \nabla f_t(\boldsymbol{\theta}) \nabla f_t(\boldsymbol{\theta})^\top \leq V(\boldsymbol{\theta}),$$

$$\mathfrak{n}_t(f_t(\boldsymbol{\theta})) \left|\boldsymbol{\gamma}^\top \nabla f_t(\boldsymbol{\theta})\right| \leq \frac{\lambda_1^*}{\lambda^*} \sqrt{\boldsymbol{\gamma}^\top V(\boldsymbol{\theta})\boldsymbol{\gamma}}. \tag{3.5}$$

*Then (3.2) is fulfilled with this $V(\boldsymbol{\theta})$ for $\lambda \leq \lambda^*$.*

*Proof.* By definition

$$2\lambda \frac{\boldsymbol{\gamma}^\top \nabla \zeta(\boldsymbol{\theta})}{\sqrt{\boldsymbol{\gamma}^\top V(\boldsymbol{\theta})\boldsymbol{\gamma}}} - 2\lambda^2 \frac{\boldsymbol{\gamma}^\top B(\boldsymbol{\theta})\boldsymbol{\gamma}}{\boldsymbol{\gamma}^\top V(\boldsymbol{\theta})\boldsymbol{\gamma}} = \sum_t \left\{2\lambda c_t \frac{\dot{\zeta}_t(Y_t, f_t(\boldsymbol{\theta}))}{\mathfrak{n}_t(f_t(\boldsymbol{\theta}))} - 2\lambda^2 c_t^2\right\}$$

where $c_t = \mathfrak{n}_t(f_t(\boldsymbol{\theta})) \boldsymbol{\gamma}^\top \nabla f_t(\boldsymbol{\theta}) \left[\boldsymbol{\gamma}^\top V(\boldsymbol{\theta})\boldsymbol{\gamma}\right]^{-1/2}$ so that $\lambda \leq \lambda^*$ implies $\lambda c_t \leq \lambda_1^*$ in view of (3.5). Now by (3.4)

$$\boldsymbol{E}\left[\exp\left\{2\lambda c_t \frac{\dot{\zeta}_t(Y_t, f_t(\boldsymbol{\theta}))}{\mathfrak{n}_t(f_t(\boldsymbol{\theta}))} - 2\lambda^2 c_t^2\right\} \middle| \mathcal{F}_{t-1}\right] \leq 1$$

and the result follows by induction arguments starting from $t = n$. $\qquad\square$

**Lemma 3.2.** *Let for some $\overline{\lambda} > 0$, the function*

$$\boldsymbol{E} \exp\left\{2\overline{\lambda}\boldsymbol{\gamma}^\top[\nabla M(\boldsymbol{\theta}) - \boldsymbol{E}\nabla M(\boldsymbol{\theta})]\right\}$$

*be uniformly continuous in $(\boldsymbol{\theta}, \boldsymbol{\gamma}) \in \Theta \times \mathcal{S}^p$. Let also a matrix $V_0(\boldsymbol{\theta})$ be uniformly continuous $\boldsymbol{\theta} \in \Theta$ and satisfy $V_0(\boldsymbol{\theta}) \geq I$ and*

$$\text{Var}\left[\nabla M(\boldsymbol{\theta})\right] \leq V_0(\boldsymbol{\theta}), \qquad \boldsymbol{\theta} \in \Theta.$$

*Then for every $\lambda^* < \overline{\lambda}$ there exists a constant $C_1 = C_1(\lambda^*, \overline{\lambda})$ such that (3.3) is fulfilled with $V(\boldsymbol{\theta}) = C_1 V_0(\boldsymbol{\theta})$.*

The result is an easy corollary of Lemma 5.8 from Golubev and Spokoiny (2009).

Define for every $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \Theta$, $u = \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|$ and $\boldsymbol{\gamma} = (\boldsymbol{\theta}' - \boldsymbol{\theta})/u$

$$\mathcal{S}^2(\boldsymbol{\theta}, \boldsymbol{\theta}') \stackrel{\text{def}}{=} u^2 \int_0^1 \boldsymbol{\gamma}^\top V(\boldsymbol{\theta} + tu\boldsymbol{\gamma})\boldsymbol{\gamma}\,dt.$$

Next, introduce for every $\boldsymbol{\theta}^\circ \in \Theta$ the local vicinity $\mathcal{B}(\epsilon, \boldsymbol{\theta}^\circ)$ such that $\mathcal{S}(\boldsymbol{\theta}, \boldsymbol{\theta}^\circ) \leq \epsilon$ for all $\boldsymbol{\theta} \in \mathcal{B}(\epsilon, \boldsymbol{\theta}^\circ)$.

Let also the matrix function $V(\cdot)$ satisfy the following regularity condition:

**(V)**   *There exist constants $\epsilon > 0$ and $\nu_1 < 1$ such that*

$$\sup_{\boldsymbol{\theta},\boldsymbol{\theta}^\circ \in \Theta:\ \mathfrak{S}(\boldsymbol{\theta},\boldsymbol{\theta}^\circ)\leq \epsilon}\ \sup_{\boldsymbol{\gamma}\in S^p} \frac{\boldsymbol{\gamma}^\top V(\boldsymbol{\theta})\boldsymbol{\gamma}}{\boldsymbol{\gamma}^\top V(\boldsymbol{\theta}^\circ)\boldsymbol{\gamma}} \leq \nu_1\,.$$

The next result presents the claimed exponential bound. It is a specification of a more general result from Theorem 5.5 in Section 5.

**Theorem 3.3.** *Assume $(E)$, $(ED)$ with some $\lambda^* > 0$ and $(V)$ with some $\nu_1$ and $\epsilon$. Let $\varrho < 1$ be such that $\varrho\epsilon/(1-\varrho) \leq \lambda^*$. If the function $\mathrm{pen}(\boldsymbol{\theta})$ fulfills*

$$\mathfrak{H}_\epsilon(\varrho) \stackrel{\mathrm{def}}{=} \log\left\{\omega_p^{-1}\epsilon^{-p}\int_\Theta \sqrt{\det(V(\boldsymbol{\theta}))}\exp\{-\varrho\,\mathrm{pen}_\epsilon(\boldsymbol{\theta})\}d\boldsymbol{\theta}\right\} < \infty \qquad (3.6)$$

*with $\mathrm{pen}_\epsilon(\boldsymbol{\theta}^\circ) = \inf_{\boldsymbol{\theta}\in\mathcal{B}(\epsilon,\boldsymbol{\theta}^\circ)}\mathrm{pen}(\boldsymbol{\theta})$ and $\omega_p$ being the volume of the unit ball in $I\!\!R^p$, then*

$$\boldsymbol{E}\exp\left\{\sup_{\boldsymbol{\theta}\in\Theta}\varrho\big[\mu L(\boldsymbol{\theta},\boldsymbol{\theta}_0) + \mathfrak{M}(\mu,\boldsymbol{\theta},\boldsymbol{\theta}_0) - \mathrm{pen}(\boldsymbol{\theta})\big]\right\} \leq \mathfrak{Q}(\varrho), \qquad (3.7)$$

*with*

$$\log\mathfrak{Q}(\varrho) = \frac{2\epsilon^2\varrho^2}{1-\varrho} + (1-\varrho)\mathbb{Q}_p + \mathfrak{H}_\epsilon(\varrho) + p\log(\nu_1)$$

*where $\mathbb{Q}_p$ is the usual entropy number for the Euclidean ball in $I\!\!R^p$.*

## 3.1   Penalty via the norm $\|\sqrt{V^*}(\boldsymbol{\theta}-\boldsymbol{\theta}_0)\|$

The choice of the penalty function $\mathrm{pen}(\boldsymbol{\theta})$ can be made more precise if the condition $(ED)$ can be checked with a constant matrix $V(\boldsymbol{\theta}) \equiv V^*$ for a fixed matrix $V^*$ and all $\boldsymbol{\theta}$. This section describes how the penalty function can be defined in terms of the norm $\|\sqrt{V^*}(\boldsymbol{\theta}-\boldsymbol{\theta}_0)\|$.

**Theorem 3.4.** *Let the conditions $(E)$ and $(ED)$ be fulfilled with $V(\boldsymbol{\theta}) \equiv V^*$ for some matrix $V^*$ for all $\boldsymbol{\theta} \in \Theta$. Let $\varrho \in (0,1)$ and $\epsilon > 0$ be fixed to ensure $\varrho\epsilon/(1-\varrho) \leq \lambda^*$. Suppose that $\varkappa(\mathfrak{r})$ is a monotonously decreasing positive function on $[0,+\infty)$ satisfying*

$$\mathfrak{P}^* \stackrel{\mathrm{def}}{=} \omega_p^{-1}\int_{I\!\!R^p}\varkappa(\|\boldsymbol{\theta}\|)d\boldsymbol{\theta} = p\int_0^\infty \varkappa(t)t^{p-1}dt < \infty. \qquad (3.8)$$

*Define*

$$\mathrm{pen}(\boldsymbol{\theta}) = -\varrho^{-1}\log\varkappa\big(\epsilon^{-1}\|\sqrt{V^*}(\boldsymbol{\theta}-\boldsymbol{\theta}_0)\| + 1\big). \qquad (3.9)$$

*Then the assertion (3.7) holds with*

$$\log\mathfrak{Q}(\varrho) = \frac{2\epsilon^2\varrho^2}{1-\varrho} + (1-\varrho)\mathbb{Q}_p + \log(\mathfrak{P}^*).$$

*Proof.* This result is a straightforward corollary of Theorem 3.3 applied with $V(\boldsymbol{\theta}) \equiv V^*$ and thus, condition $(V)$ is fulfilled with $\nu_1 = 1$. $\hfill\square$

Here two natural ways of defining the penalty function $\text{pen}(\boldsymbol{\theta})$: quadratic or logarithmic in $\left\|\sqrt{V^*}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\right\|$. The functions $\varkappa(\cdot)$ and the corresponding $\mathfrak{P}^*$-values are:

$$
\begin{aligned}
\varkappa_1(u) &= e^{-\delta_1(t-1)_+^2}, & \mathfrak{P}_1^* &= 1 + \omega_p^{-1}(\pi/\delta_1)^{p/2}, \\
\varkappa_2(t) &= (t+1)^{-p-\delta_2}, & \mathfrak{P}_2^* &= p/\delta_2,
\end{aligned}
\tag{3.10}
$$

where $\delta_1, \delta_2 > 0$ are some constant and $[a]_+$ means $\max\{a, 0\}$. The corresponding penalties read as:

$$
\begin{aligned}
\text{pen}_1(\boldsymbol{\theta}) &= \varrho^{-1}\delta_1\,\epsilon^{-2}\left\|\sqrt{V^*}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\right\|^2. \\
\text{pen}_2(\boldsymbol{\theta}) &= -\varrho^{-1}(p + \delta_2)\log\left(\epsilon^{-1}\left\|\sqrt{V^*}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\right\| + 2\right).
\end{aligned}
$$

## 3.2   Some corollaries

The result of Theorem 3.3 means that the value $\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \text{pen}(\boldsymbol{\theta})$ is uniformly in $\boldsymbol{\theta} \in \Theta$ stochastically bounded. In particular, one can plug the estimate $\widetilde{\boldsymbol{\theta}}$ in place of $\boldsymbol{\theta}$: with some $\varrho < 1$

$$
\boldsymbol{E}\exp\left\{\varrho\left[\mu L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \text{pen}(\widetilde{\boldsymbol{\theta}})\right]\right\} \le \mathfrak{Q}(\varrho).
\tag{3.11}
$$

Below we present some corollaries of this result.

To simplify the presentation, we consider the case when there is a deterministic function $\overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ such that the following bound holds almost sure:

$$
\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \ge \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0), \qquad \boldsymbol{\theta} \in \Theta
\tag{3.12}
$$

### 3.2.1   Concentration properties of the estimator $\widetilde{\boldsymbol{\theta}}$

Define for every subset $A$ of the parameter set $\Theta$ the value

$$
\mathfrak{z}(A) \overset{\text{def}}{=} \inf_{\boldsymbol{\theta} \notin A}\{\overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \text{pen}(\boldsymbol{\theta})\}.
\tag{3.13}
$$

The next result shows that the estimator $\widetilde{\boldsymbol{\theta}}$ deviates out of the set $A$ with an exponentially small probability of order $\exp\{-\varrho\mathfrak{z}(A)\}$.

**Corollary 3.5.** *Suppose (3.11). Then for any set $A \subset \Theta$*

$$
\boldsymbol{P}\left(\widetilde{\boldsymbol{\theta}} \notin A\right) \le \mathfrak{Q}(\varrho)\mathrm{e}^{-\varrho\mathfrak{z}(A)}.
$$

14

*Proof.* If $\widetilde{\boldsymbol{\theta}} \notin A$, then $\mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \mathrm{pen}(\widetilde{\boldsymbol{\theta}}) \geq \mathfrak{z}(A)$. As $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) \geq 0$, it follows

$$
\begin{aligned}
\mathfrak{Q}(\varrho) &\geq \boldsymbol{E} \exp\Big\{ \varrho\big[\mu L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \mathrm{pen}(\widetilde{\boldsymbol{\theta}})\big]\Big\} \\
&\geq \boldsymbol{E} \exp\Big\{ \varrho\big[\mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \mathrm{pen}(\widetilde{\boldsymbol{\theta}})\big]\Big\} \geq \mathrm{e}^{\mathfrak{z}(A)} \boldsymbol{P}\big(\widetilde{\boldsymbol{\theta}} \notin A\big)
\end{aligned}
$$

as required. $\qquad\square$

Two particular choices of the set $A$ can be mentioned:

$$
\begin{aligned}
A &= \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0) = \{\boldsymbol{\theta} : \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \leq \mathfrak{r}\}, \\
A &= \mathcal{A}'(\mathfrak{r}, \boldsymbol{\theta}_0) = \{\boldsymbol{\theta} : \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \mathrm{pen}(\boldsymbol{\theta}) \leq \mathfrak{r}\},
\end{aligned}
$$

For the set $\mathcal{A}'(\mathfrak{r}, \boldsymbol{\theta}_0)$, Corollary 3.5 yields

$$
\boldsymbol{P}\big(\widetilde{\boldsymbol{\theta}} \notin \mathcal{A}'(\mathfrak{r}, \boldsymbol{\theta}_0)\big) = \boldsymbol{P}\big(\mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \mathrm{pen}(\widetilde{\boldsymbol{\theta}}) \geq \mathfrak{r}\big) \leq \mathfrak{Q}(\varrho)\mathrm{e}^{-\varrho\mathfrak{r}}.
$$

For the set $\mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)$, define additionally the value $\mathfrak{b}(\mathfrak{r})$ by the relation

$$
\overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \mathrm{pen}(\boldsymbol{\theta}) \geq \mathfrak{r} - \mathfrak{b}(\mathfrak{r}), \qquad \boldsymbol{\theta} \in \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0),
$$

or, equivalently,

$$
\mathfrak{b}(\mathfrak{r}) = \sup_{\boldsymbol{\theta} \in \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)} \big\{ \mathfrak{r} + \mathrm{pen}(\boldsymbol{\theta}) - \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \big\}.
$$

**Corollary 3.6.** *Suppose (3.11). Then for any* $\mathfrak{r} > 0$

$$
\boldsymbol{P}\big(\widetilde{\boldsymbol{\theta}} \notin \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)\big) = \boldsymbol{P}\big(\overline{\mathfrak{M}}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) \geq \mathfrak{r}\big) \leq \mathfrak{Q}(\varrho)\mathrm{e}^{-\varrho[\mathfrak{r} - \mathfrak{b}(\mathfrak{r})]}.
$$

In typical situations the value $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ and thus, $\overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ is nearly proportional to the sample size $n$ and is nearly quadratic in $\boldsymbol{\theta} - \boldsymbol{\theta}_0$ so that and each set $\mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)$ corresponds to a root-$n$ neighborhood of the point $\boldsymbol{\theta}_0$, and the concentration property becomes a non-asymptotic analog of root-n consistency. See below Section 3.4 for a precise formulation.

It is important to stress that for applying the result of Corollary 3.6, it is not required to compute the rate function $\mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0)$ and the penalty function $\mathrm{pen}(\boldsymbol{\theta})$. It only suffices to obtain some rough upper bound for the penalty function and deterministic lower bound for the rate function. The result claims that the estimate well localizes on a vicinity $\mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)$ of the point $\boldsymbol{\theta}_0$.

Another remark concerns the identifiability issue. It was already mentioned in the introduction that the results do not require any identifiability condition. However, if the model parameter is not well identifiable this leads to the situation that the rate function

$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ is very flat and its level sets are quite big. Therefore, a poor parametrization leads to a less informative concentration property. In particular, the set $\mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)$ can be unbounded or disconnected.

The concentration property is very useful in combination with a more fine analysis based on the Taylor expansion of the (quasi) log-likelihood. Indeed, it ensures that the estimate belongs with a high probability to a small vicinity of $\boldsymbol{\theta}_0$ and in this vicinity the classical asymptotic technique based on the second order approximation of the process $L(\boldsymbol{\theta})$ can be used to address the issues of asymptotic distribution and asymptotic efficiency.

### 3.2.2 Confidence sets based on $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta})$

Next we discuss how the exponential bound can be used for establishing some risk bounds and for constructing the confidence sets for the target $\boldsymbol{\theta}_0$ based on the maximized value $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta})$. The inequality (3.11) claims that $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0)$ is stochastically bounded with finite exponential moments.

Define

$$\mathfrak{b} \stackrel{\text{def}}{=} \mathfrak{b}(0) = \sup_{\boldsymbol{\theta}}[\text{pen}(\boldsymbol{\theta}) - \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)]_+ \,. \tag{3.14}$$

**Corollary 3.7.** *Suppose (3.11) and let $\mathfrak{b}$ from (3.14) be finite. Then*

$$\boldsymbol{E} \exp\{\varrho\mu L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0)\} \leq e^{\varrho\mathfrak{b}} \mathfrak{Q}(\varrho).$$

*Proof.* Observe that

$$\boldsymbol{E} \exp\{\varrho\mu L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0)\} \leq e^{\varrho\mathfrak{b}} \boldsymbol{E} \exp\{\varrho[\mu L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \text{pen}(\widetilde{\boldsymbol{\theta}})]\} \leq e^{\varrho\mathfrak{b}} \mathfrak{Q}(\varrho).$$

This obviously yields the assertion. $\qquad\qquad\square$

By the same reasons, one can construct confidence sets based on the (quasi) likelihood process. Define

$$\mathcal{E}(\mathfrak{z}) = \{\boldsymbol{\theta} \in \Theta : L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \leq \mathfrak{z}\}.$$

The bound for $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0)$ ensures that $\boldsymbol{\theta}_0$ belongs to this set with a high probability provided that $\mathfrak{z}$ is large enough. The next result claims that $\mathcal{E}(\mathfrak{z})$ does not cover the true value $\boldsymbol{\theta}_0$ with a probability which decreases exponentially with $\mathfrak{z}$.

**Corollary 3.8.** *Suppose (3.11). For any $\mathfrak{z} > 0$*

$$\boldsymbol{P}(\boldsymbol{\theta}_0 \notin \mathcal{E}(\mathfrak{z})) \leq \mathfrak{Q}(\varrho) \exp\{-\varrho\mu\mathfrak{z} + \varrho\mathfrak{b}\}.$$

*Proof.* The bound (3.11) implies for the event $\{\boldsymbol{\theta}_0 \notin \mathcal{E}(\mathfrak{z})\} = \{L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) > \mathfrak{z}\}$

$$
\begin{aligned}
\boldsymbol{P}\big\{\boldsymbol{\theta}_0 \notin \mathcal{E}(\mathfrak{z})\big\} &\leq \boldsymbol{P}\big\{\varrho\big[\mu L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \mathrm{pen}(\widetilde{\boldsymbol{\theta}})\big] > \varrho\mu\mathfrak{z} - \varrho\mathfrak{b}\big\} \\
&\leq \exp\big\{-\varrho\mu\mathfrak{z} + \varrho\mathfrak{b}\big\}\boldsymbol{E}\exp\big\{\varrho\mu L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) - \mathrm{pen}(\widetilde{\boldsymbol{\theta}})\big\} \\
&\leq \mathfrak{Q}(\varrho)\exp\big\{-\varrho\mu\mathfrak{z} + \varrho\mathfrak{b}\big\}
\end{aligned}
$$

as required. $\qquad\square$

The result of Corollary 3.8 only presents an upper bound for the coverage probability of the value $\boldsymbol{\theta}_0$ by the set $\mathcal{E}(\mathfrak{z})$. The given exponential bound contains some implicit constant and is rather rough, and therefore, it can hardly be used for computing the coverage probability and for fixing the constant $\mathfrak{z}_\alpha$ which ensures the coverage level $1 - \alpha$. It would be unrealistic to obtain a universal non-asymptotic sharp bound for the coverage level which applies in such a general situation. However, the result is meaningful because it suggests the form of the confidence set and guarantees that the choice a sufficiently big but fixed threshold $\mathfrak{z}$ ensures the prescribed coverage probability. A precise value can be found by the Monte-Carlo simulations, see e.g. Spokoiny (2007) for some examples.

## 3.3 Identifiability condition

Until this point no any identifiability condition on the model has been used, that is, the presented results apply even for a very poor parametrization. Actually, a particular parametrization of the parameter set plays no role as long as the value of maximum is considered. If we want to derive any quantitative result on the point of maximum $\widetilde{\boldsymbol{\theta}}$, then the parametrization matters and an identifiability condition is really necessary. Here we follow the usual path by applying the quadratic lower bound for the rate function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ in a vicinity of the point $\boldsymbol{\theta}_0$.

Finally we discuss a risk bound in terms of the classical loss $\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0$. The idea is to apply the quadratic lower bound for the rate function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ in a vicinity of the point $\boldsymbol{\theta}_0$ and to use the concentration property of the estimator $\widetilde{\boldsymbol{\theta}}$. To explain the conditions imposed below suppose that the log-likelihood function is two times continuously differentiable in $\boldsymbol{\theta}$. This implies the differentiability in $\boldsymbol{\theta}$ of the moment generating function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = -\sum_t \mathfrak{m}_t\{\mu, f_t(\boldsymbol{\theta}), f_t(\boldsymbol{\theta}_0)\}$. Obviously $\mathfrak{M}(\mu, \boldsymbol{\theta}_0, \boldsymbol{\theta}_0) = 0$ and a simple algebra yields for the gradient $\nabla\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = d\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)/d\boldsymbol{\theta}$:

$$
\boldsymbol{E}\nabla\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} = -\mu\boldsymbol{E}\nabla L(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} = -\mu\nabla\boldsymbol{E}L(\boldsymbol{\theta}_0) = 0
$$

because $\boldsymbol{\theta}_0$ is the point of maximum of $\boldsymbol{E}L(\boldsymbol{\theta})$. The same holds automatically for the lower bound $\overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$. The Taylor expansion of the second order in a vicinity of $\boldsymbol{\theta}_0$ yields for all $\boldsymbol{\theta}$ close to $\boldsymbol{\theta}_0$ the following approximation:

$$\overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \approx \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top I(\mu, \boldsymbol{\theta}_0)(\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

with the matrix $I(\mu, \boldsymbol{\theta}_0) = \boldsymbol{E}\nabla^2 \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$. This and the concentration property from Corollary 3.6 lead to the following bound on $\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0$:

**Corollary 3.9.** *Let (3.11) hold. Suppose that for some positive symmetric matrix $D$ and some $\mathfrak{r} > 0$, the function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ fulfills almost surely*

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top D^2(\boldsymbol{\theta} - \boldsymbol{\theta}_0), \qquad \boldsymbol{\theta} \in \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0), \qquad (3.15)$$

*Then for any $\mathfrak{z} \leq \mathfrak{r}$*

$$\boldsymbol{P}\big(\|D(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)\|^2 > \mathfrak{z}\big) \leq \mathfrak{Q}(\varrho)\mathrm{e}^{-\varrho[\mathfrak{z} - \mathfrak{b}(\mathfrak{z})]}.$$

*Proof.* It is obvious that

$$
\begin{aligned}
\big\{\|D(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)\|^2 > \mathfrak{z}\big\} \quad &\subseteq \quad \big\{\|D(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)\|^2 > \mathfrak{z}, \widetilde{\boldsymbol{\theta}} \in \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)\big\} \cup \big\{\widetilde{\boldsymbol{\theta}} \notin \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)\big\} \\
&\subseteq \quad \big\{\mathfrak{M}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) > \mathfrak{z}, \widetilde{\boldsymbol{\theta}} \in \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)\big\} \cup \big\{\overline{\mathfrak{M}}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) > \mathfrak{z}\big\} \\
&= \quad \big\{\overline{\mathfrak{M}}(\mu, \widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) > \mathfrak{z}\big\}
\end{aligned}
$$

and the result follows from Corollary 3.6. $\qquad\qquad\square$

In the next theorem we assume the lower bound (3.15) to be fulfilled on the whole parameter set $\Theta$. The general case can be reduced to this one by using once again the concentration property of Corollary 3.6.

**Theorem 3.10.** *Suppose $(E)$, $(ED)$ with $V(\boldsymbol{\theta}) \leq V^*$ for a matrix $V^*$. Let also for some $\mathfrak{a} > 0$ and a deterministic function $\overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ hold*

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq \mathfrak{a}^2(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top V^*(\boldsymbol{\theta} - \boldsymbol{\theta}_0), \qquad \boldsymbol{\theta} \in \Theta. \qquad (3.16)$$

*Fix some $\mathfrak{a}_1 \leq \mathfrak{a}$ and define $\mathrm{pen}(\boldsymbol{\theta})$ by*

$$\mathrm{pen}(\boldsymbol{\theta}) = \mathfrak{a}_1^2(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top V^*(\boldsymbol{\theta} - \boldsymbol{\theta}_0). \qquad (3.17)$$

*Then with $s = 1 - \mathfrak{a}_1^2/\mathfrak{a}^2$ it holds*

$$
\begin{aligned}
\mathfrak{Q}(\varrho, s) \quad &\stackrel{\mathrm{def}}{=} \quad \log \boldsymbol{E}\exp\big\{\varrho \sup_{\boldsymbol{\theta}}\big[\mu L(\boldsymbol{\theta}, \boldsymbol{\theta}_0) + \mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) - \mathrm{pen}(\boldsymbol{\theta})\big]\big\} \\
&\leq \quad 2\varrho + (1 - \varrho)\mathbb{Q}_p + \log\left(1 + \frac{\omega_p^{-1}\pi^{p/2}}{(1 - \varrho)^{p/2}\mathfrak{a}_1^p}\right) \\
&\leq \quad pC(\varrho) + p\log\big(|\mathfrak{a}^2(1 - s)(1 - \varrho)|^{-1/2}\big) \qquad (3.18)
\end{aligned}
$$

18

*for some fixed constant* $C(\varrho)$. *In addition,* $\mathfrak{b}(\mathfrak{r}) = 0$ *for all* $\mathfrak{r} \geq 0$ *yielding for any* $\mathfrak{z} > 0$ *the concentration property and confidence bound:*

$$
\begin{aligned}
\boldsymbol{P}\big(\widetilde{\boldsymbol{\theta}} \notin \mathcal{A}(\mathfrak{z}, \boldsymbol{\theta}_0)\big) &\leq \mathfrak{Q}(\varrho, s)\mathrm{e}^{-\varrho s \mathfrak{z}}, & \mathcal{A}(\mathfrak{z}, \boldsymbol{\theta}_0) &= \{\boldsymbol{\theta} : \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \leq \mathfrak{z}\}, \\
\boldsymbol{P}\big(\boldsymbol{\theta}_0 \notin \mathcal{E}(\mathfrak{z})\big) &\leq \mathfrak{Q}(\varrho, 0)\mathrm{e}^{-\varrho \mathfrak{z}}, & \mathcal{E}(\mathfrak{z}) &= \{\boldsymbol{\theta} : L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \leq \mathfrak{z}\}.
\end{aligned}
$$

*Proof.* We apply Theorem 3.4 with

$$
\varkappa(t) = \exp\big\{-(1 - \varrho)\mathfrak{a}_1^2(t - 1)_+^2\big\}
$$

leading for $\epsilon^2 = (1 - \varrho)/\varrho$ and $t = \epsilon^{-1}\big\|\sqrt{V^*}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)\big\|$ to the formula (3.17) for $\mathrm{pen}(\boldsymbol{\theta})$. By simple algebra

$$
\mathfrak{P}^* = \omega_p^{-1} \int_{I\!\!R^p} \varkappa(\|\boldsymbol{\theta}\|)d\boldsymbol{\theta} = 1 + \omega_p^{-1}\frac{\pi^{p/2}}{(1 - \varrho)^{p/2}\mathfrak{a}_1^p};
$$

cf. the bound (3.10) for $\mathfrak{P}^*$ with $\delta_1 = (1 - \varrho)\mathfrak{a}_1^2$. This implies the bound (3.18) for the $\mathfrak{Q}(\varrho)$ because $p^{-1}\mathbb{Q}_p$ and $p^{-1}\log\omega_p^{-1}$ are bounded by some fixed constants.

The inequality (3.16) ensures for $\mathfrak{r} = \overline{\mathfrak{M}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ that $\mathrm{pen}(\boldsymbol{\theta}) \leq \mathfrak{a}_1^2/\mathfrak{a}^2\mathfrak{r}$, i.e. $\mathfrak{b}(\mathfrak{r}) \leq \mathfrak{a}_1^2/\mathfrak{a}^2\mathfrak{r}$ and $\mathfrak{b} = \mathfrak{b}(0) = 0$. Finally, the concentration and coverage bounds follow from Corollaries 3.6 and 3.8. $\qquad\square$

## 3.4   Sub-ergodocity and root-n consistency

Consider for every $\boldsymbol{\theta} \in \Theta$ a $p \times p$ random matrix $B(\boldsymbol{\theta}) = \sum_t \mathfrak{n}_t^2(f_t(\boldsymbol{\theta}))\,\nabla f_t(\boldsymbol{\theta})\nabla f_t(\boldsymbol{\theta})^\top$; see (3.5). Condition (3.4) implies that for any $\boldsymbol{\gamma} \in \mathcal{S}^p$ and $|\lambda| \leq \lambda^*$

$$
\boldsymbol{E}\exp\big\{2\lambda\boldsymbol{\gamma}^\top\nabla\zeta(\boldsymbol{\theta}) - 2\lambda^2\boldsymbol{\gamma}^\top B(\boldsymbol{\theta})\boldsymbol{\gamma}\big\} \leq 1.
$$

The usual ergodicity condition for the sum $B(\boldsymbol{\theta})$ means that $n^{-1}B(\boldsymbol{\theta})$ converges to some deterministic matrix $b(\boldsymbol{\theta})$ for every $\boldsymbol{\theta}$ as $n$ grows. Sub-ergodicity can be understood in the sense that $n^{-1}B(\boldsymbol{\theta})$ is bounded by some deterministic matrix $v(\boldsymbol{\theta})$ with a high probability. We define $V(\boldsymbol{\theta}) = nv(\boldsymbol{\theta})$ and suppose that conditions $(ED)$ and $(V)$ are fulfilled for such defined $V(\boldsymbol{\theta})$.

Similarly the sub-ergodicity applied to the random quantity $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ means that there is a deterministic positive function $\overline{\mathfrak{m}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ such that $n^{-1}\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq \overline{\mathfrak{m}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ with a high probability. In this situation one can rewrite the main corollaries from Section 3.2 in terms of the functions $v(\boldsymbol{\theta})$ and $\overline{\mathfrak{m}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$. In particular, the concentration set $\mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)$ can be replaced by

$$
\overline{\mathcal{A}}(\mathfrak{r}, \boldsymbol{\theta}_0) \stackrel{\mathrm{def}}{=} \{\boldsymbol{\theta} : n\overline{\mathfrak{m}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \leq \mathfrak{r}\}.
$$

In addition we assume similarly to (3.15) that for some fixed symmetric positive matrix $D_1$ and some $\mathfrak{r} > 0$, it holds in the vicinity $\mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0)$ of the point $\boldsymbol{\theta}_0$:

$$\overline{\mathfrak{m}}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top D_1^2 (\boldsymbol{\theta} - \boldsymbol{\theta}_0), \qquad v(\boldsymbol{\theta}) \leq \mathfrak{a}^2 D_1^2 \quad \boldsymbol{\theta} \in \mathcal{A}(\mathfrak{r}, \boldsymbol{\theta}_0). \qquad (3.19)$$

**Corollary 3.11.** *Assume (3.11) and (3.19) for some* $\mathfrak{r} > 0$*. Then for any* $\mathfrak{z} \leq \mathfrak{r}$

$$\boldsymbol{P}\big(\|D_1(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)\|^2 > \mathfrak{z}/n\big) \leq \mathfrak{Q}(\varrho) \exp\{-\varrho(\mathfrak{z} - \mathfrak{b})\}.$$

# 4 Applications and examples

This section illustrates how the general results can be applied to some popular examples of parametric time series models which we already mentioned in Section 2.3. For all examples we assume that the two components of the parametric modeling are fixed: a parametric family $\mathcal{P} = (P_v, v \in \mathcal{U})$ and a family of parameter processes $\{f(\boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$. Moreover, we assume that both families are sufficiently regular, in particular, the functions $f_t(\boldsymbol{\theta})$ are differentiable w.r.t. $\boldsymbol{\theta}$. The corresponding gradient is denoted by $\nabla f_t(\boldsymbol{\theta}) = \partial f_t(\boldsymbol{\theta})/\partial \boldsymbol{\theta} \in I\!\!R^p$.

In all the examples, the real model dynamics is described by some predictable process $f_t$ accepting that the parametric assumption $f_t = f_t(\boldsymbol{\theta})$ is not precisely fulfilled whatever $\boldsymbol{\theta} \in \Theta$ is. By $\boldsymbol{\theta}_0$ we denote the parameter corresponding to the best parametric fit: $\boldsymbol{\theta}_0 = \operatorname{argmax}_{\boldsymbol{\theta}} \boldsymbol{E}L(\boldsymbol{\theta})$. Such defined vector parameter $\boldsymbol{\theta}_0$ is estimated by maximizing the objective function $L(\boldsymbol{\theta})$.

## 4.1 Linear autoregression

Assume the model $Y_t = f_t + \varepsilon_t$ in which the innovations $\varepsilon_t$ are martingale differences: $\boldsymbol{E}[\varepsilon_t|\mathcal{F}_{t-1}] = 0$ possible heterogeneous with bounded exponential moments: $\log \boldsymbol{E} \exp\{\lambda \varepsilon_t|\mathcal{F}_{t-1}\} \leq \varkappa_1 \lambda^2$ for $\lambda \leq \lambda^*$ and some fixed $\varkappa_1$.

With $L(\boldsymbol{\theta}) = -(2\sigma^2)^{-1} \sum_t \{Y_t - f_t(\boldsymbol{\theta})\}^2$, the least square estimate $\widetilde{\boldsymbol{\theta}}$ reads as

$$\widetilde{\boldsymbol{\theta}} = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname*{argmin}_{\boldsymbol{\theta}} \sum_t \{Y_t - f_t(\boldsymbol{\theta})\}^2.$$

Linear autoregression means the structural equation $f_t(\boldsymbol{\theta}) = \alpha_1 Y_{t-1} + \ldots + \alpha_p Y_{t-p}$ for $\boldsymbol{\theta} = (\alpha_1, \ldots, \alpha_p)^\top$ leading to maximization of the quadratic functional $\sum_t (Y_t - \Psi_t^\top \boldsymbol{\theta})^2$ with $\Psi_t = (Y_{t-1}, \ldots, Y_{t-p})^\top$ which admits a closed form solution:

$$\widetilde{\boldsymbol{\theta}} = \left(\sum_t \Psi_t \Psi_t^\top\right)^{-1} \sum_t Y_t \Psi_t = B^{-1} \sum_t Y_t \Psi_t$$

$$L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) = \frac{1}{2\sigma^2}(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^\top B(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0),$$

with $B \stackrel{\text{def}}{=} \sum_t \Psi_t \Psi_t^\top$ .

The matrix $B(\boldsymbol{\theta})$ from (3.5) does not depend on $\boldsymbol{\theta}$: $B(\boldsymbol{\theta}) \equiv B$. Next, the formulae for $\nabla\zeta(\boldsymbol{\theta})$ and $\nabla M(\boldsymbol{\theta})$ simplifies to

$$
\begin{aligned}
\nabla\zeta(\boldsymbol{\theta}) &= \sigma^{-2}\sum_t \{Y_t - f_t\}\Psi_t = \sigma^{-2}\sum_t \varepsilon_t \Psi_t \,, \\
\nabla M(\boldsymbol{\theta}) &= \sigma^{-2}\sum_t \{f_t - \Psi_t^\top \boldsymbol{\theta}\}\Psi_t \,.
\end{aligned}
$$

A natural choice for the matrix $V(\boldsymbol{\theta})$ is

$$
V(\boldsymbol{\theta}) = C_1 \sigma^{-2}\boldsymbol{E}B + C_2 \operatorname{Var}\nabla M(\boldsymbol{\theta})
$$

for some $C_1, C_2 \geq 1$. Such defined matrix $V(\boldsymbol{\theta})$ is a quadratic function of $\boldsymbol{\theta}$ and condition $(V)$ is straightforward. Condition $(ED)$ is also easy to check in typical situations.

The value $\boldsymbol{\theta}_0$ which maximizes $\boldsymbol{E}L(\boldsymbol{\theta})$ can be found by the following optimization problem:

$$
\begin{aligned}
\boldsymbol{\theta}_0 &= \operatorname*{argmax}_{\boldsymbol{\theta}} \boldsymbol{E}L(\boldsymbol{\theta}) \\
&= \operatorname*{argmin}_{\boldsymbol{\theta}} \boldsymbol{E}\sum_t \left(Y_t - \Psi_t \boldsymbol{\theta}\right)^2 \\
&= \operatorname*{argmin}_{\boldsymbol{\theta}} \boldsymbol{E}\sum_t \left(f_t - \Psi_t \boldsymbol{\theta}\right)^2 = \left(\boldsymbol{E}\sum_t \Psi_t \Psi_t^\top\right)^{-1}\boldsymbol{E}\sum_t f_t \Psi_t \,. \quad (4.1)
\end{aligned}
$$

Theorem 3.3 and Corollary 3.7 claim that $L(\widetilde{\boldsymbol{\theta}}, \boldsymbol{\theta}_0) = (2\sigma^2)^{-1}(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^\top B(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$ is stochastically bounded with finite polynomial and exponential moments:

$$
\boldsymbol{E}\big|(2\sigma^2)^{-1}(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^\top B(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)\big|^r \leq \mathfrak{R}(r).
$$

This also justifies the use of confidence sets in the form

$$
\mathcal{E}(\mathfrak{z}) \stackrel{\text{def}}{=} \big\{\boldsymbol{\theta} : (2\sigma^2)^{-1}(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta})^\top B(\widetilde{\boldsymbol{\theta}} - \boldsymbol{\theta}) \leq \mathfrak{z}\big\}.
$$

By Corollary 3.8, this set does not contain the target $\boldsymbol{\theta}_0$ with a probability exponentially decreasing in $\mathfrak{z}$. The result is valid simultaneously for stationary, stable (unit root) and explosive cases, also mixed structure is allowed. The only essential assumption is about exponential moments of the errors.

Finally we briefly discuss the concentration property of the estimate $\widetilde{\boldsymbol{\theta}}$. For applying Corollary 3.6 one has to compute or evaluate the rate function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$. Suppose for simplicity that the errors $\varepsilon_t$ are conditionally normal with zero mean and the conditional variance $\sigma_t^2$, that is, $Q_t = \mathcal{N}(0, \sigma_t^2)$. Then for any $v, v'$

$$
\begin{aligned}
\mathfrak{m}_t(\mu, v, v') &\stackrel{\text{def}}{=} -\log \boldsymbol{E}\big[\exp\{\mu\ell(Y_t, v) - \mu\ell(Y_t, v')\}\big|\mathcal{F}_{t-1}\big] \\
&= \frac{\mu}{2\sigma^2}\big\{(f_t - v)^2 - (f_t - v')^2\big\} - \frac{\mu^2 \sigma_t^2}{2\sigma^4}(v' - v)^2.
\end{aligned}
$$

21

Therefore,

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = \frac{\mu}{2\sigma^2} \sum_t \left\{ (f_t - \Psi_t \boldsymbol{\theta})^2 - (f_t - \Psi_t \boldsymbol{\theta}_0)^2 - \frac{\mu \sigma_t^2}{\sigma^2} (\Psi_t \boldsymbol{\theta} - \Psi_t \boldsymbol{\theta}_0)^2 \right\}.$$

In particular, under the parametric assumption $f_t = \Psi_t^\top \boldsymbol{\theta}_0$

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = \frac{\mu}{2\sigma^2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \sum_t (1 - \mu \sigma_t^2 / \sigma^2) \Psi_t \Psi_t^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0).$$

If, in addition, the variance $\sigma_t$ is homogeneous, $\sigma_t \equiv \sigma_0$, then

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \frac{\mu}{2\sigma^2} (1 - \mu \sigma_0^2 / \sigma^2) B (\boldsymbol{\theta} - \boldsymbol{\theta}_0).$$

Optimizing w.r.t. $\mu$ yields $\mu = \sigma^2 / (2\sigma_0^2)$ and $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top B (\boldsymbol{\theta} - \boldsymbol{\theta}_0) / (4\sigma_0^2)$. In the ergodic case, when the matrix $n^{-1} B$ is close to the stationary limit $b$, the result of Corollary 3.6 claims that $\widetilde{\boldsymbol{\theta}}$ concentrates in the root-$n$ neighborhood $\mathcal{A}(\mathfrak{z}, \boldsymbol{\theta}_0) = \{ \boldsymbol{\theta} : \| b^{1/2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \|^2 \le 4\sigma_0^2 \mathfrak{z}/n \}$ of $\boldsymbol{\theta}_0$.

If the parametric assumption $f_t = \Psi_t^\top \boldsymbol{\theta}$ is not fulfilled whatever $\boldsymbol{\theta}$ is, then the calculations become only slightly more complicated. Namely,

$$
\begin{aligned}
\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \;=\; & \frac{\mu}{2\sigma^2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \sum_t (1 - \mu \sigma_t^2 / \sigma^2) \Psi_t \Psi_t^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \\
& + \frac{\mu}{\sigma^2} \sum_t \left( f_t \Psi_t - \Psi_t \Psi_t^T \boldsymbol{\theta}_0 \right) (\boldsymbol{\theta} - \boldsymbol{\theta}_0).
\end{aligned}
$$

Due to (4.1), the expectation of the second sum is zero and this term is typically smaller in order than the first one. So, $\mathcal{A}(\mathfrak{z}, \boldsymbol{\theta}_0)$ remains to be a concentration set for $\widetilde{\boldsymbol{\theta}}$.

## 4.2   Categorical time series

Let $\mathcal{P} = (P_\upsilon, \upsilon \in \mathcal{U})$ be a given exponential family with the canonical parameter. The conditional distribution $Q_t = \mathcal{L}(Y_t | \mathcal{F}_{t-1})$ of every observation $Y_t$ given the past $\mathcal{F}_{t-1}$ is assumed to be in $\mathcal{P}$ and described by the varying stochastic parameter $f_t$: $Q_t = P_{f_t} \in \mathcal{P}$. The parametric assumption $Q_t = P_{f_t(\boldsymbol{\theta})}$ leads to the log-likelihood $L(\boldsymbol{\theta}) = \sum_t \ell(Y_t, f_t(\boldsymbol{\theta}))$ where $\ell(y, \upsilon) = y\upsilon - d(\upsilon)$ is the log-likelihood function for $\mathcal{P}$:

$$L(\boldsymbol{\theta}) = \sum_t \ell(Y_t, f_t(\boldsymbol{\theta})) = \sum_t \left\{ Y_t f_t(\boldsymbol{\theta}) - d(f_t(\boldsymbol{\theta})) \right\}. \tag{4.2}$$

We use the well known properties of the canonical exponential families: $E_\upsilon Y = \dot{d}(\upsilon)$ and $\log E_\upsilon \exp\{\mu Y\} = d(\upsilon + \mu) - d(\upsilon)$. This yields for every $\upsilon, \upsilon' \in \mathcal{U}$ and all $t$

$$
\begin{aligned}
m_t(\upsilon) &\stackrel{\text{def}}{=} \; E\big[\ell(Y_t, \upsilon) \big| \mathcal{F}_{t-1}\big] = \dot{d}(f_t)\upsilon - d(\upsilon), \\
\zeta_t(Y_t, \upsilon) &\stackrel{\text{def}}{=} \; E\big[\ell(Y_t, \upsilon) - m_t(\upsilon) \big| \mathcal{F}_{t-1}\big] = \{Y_t - \dot{d}(f_t)\}\upsilon, \\
\mathfrak{m}_t(\mu, \upsilon, \upsilon') &\stackrel{\text{def}}{=} \; -\log E\big[\exp\{\mu\ell(Y_t, \upsilon, \upsilon')\} \big| \mathcal{F}_{t-1}\big] \\
&= \; d(f_t) - d(f_t + \mu(\upsilon - \upsilon')) + \mu d(\upsilon) - \mu d(\upsilon').
\end{aligned}
$$

It is easy to see that condition $(E)$ is fulfilled if $f_t + \mu\{f_t(\boldsymbol{\theta}) - f_t(\boldsymbol{\theta}_0)\} \in \mathcal{U}$ for all $t$ and all $\boldsymbol{\theta} \in \Theta$.

Next, $\dot{\zeta}_t(Y_t, \upsilon) \overset{\text{def}}{=} \partial\zeta_t(Y_t, \upsilon)/\partial\upsilon = Y_t - \dot{d}(f_t)$. In particular, this expression does not depend on $\upsilon$. Let $\mathfrak{n}(\upsilon)$ be a function of $\upsilon$ which ensures for some fixed $\lambda_1^* > 0$ that

$$\log E_\upsilon \exp\left\{2\lambda\frac{Y - \dot{d}(\upsilon)}{\mathfrak{n}(\upsilon)}\right\} \le 2\lambda^2, \qquad \lambda \le \lambda_1^*$$

Define

$$M(\boldsymbol{\theta}) \overset{\text{def}}{=} \sum_t m_t\big(f_t(\boldsymbol{\theta})\big) = \sum_t \big\{\dot{d}(f_t)f_t(\boldsymbol{\theta}) - d(f_t(\boldsymbol{\theta}))\big\},$$

$$\zeta(\boldsymbol{\theta}) \overset{\text{def}}{=} \sum_t \zeta_t\big(Y_t, f_t(\boldsymbol{\theta})\big) = \sum_t \big\{Y_t - \dot{d}(f_t)\big\}f_t(\boldsymbol{\theta}),$$

$$B(\boldsymbol{\theta}) \overset{\text{def}}{=} \sum_t \mathfrak{n}^2(f_t)\,\nabla f_t(\boldsymbol{\theta})\big\{\nabla f_t(\boldsymbol{\theta})\big\}^\top.$$

Then under simple conditions on the parameter process $f_t(\boldsymbol{\theta})$, the derivatives

$$\nabla M(\boldsymbol{\theta}) \overset{\text{def}}{=} \sum_t \big\{\dot{d}(f_t) - \dot{d}(f_t(\boldsymbol{\theta}))\big\}\nabla f_t(\boldsymbol{\theta}),$$

$$\nabla\zeta(\boldsymbol{\theta}) \overset{\text{def}}{=} \sum_t \big\{Y_t - \dot{d}(f_t)\big\}\nabla f_t(\boldsymbol{\theta}),$$

fulfill the conditions $(ED)$ and $(V)$ with $V(\boldsymbol{\theta}) = C_1\boldsymbol{E}B(\boldsymbol{\theta}) + C_2\,\mathrm{Var}\,\nabla M(\boldsymbol{\theta})$ for some $C_1, C_2 \ge 1$.

Sub-ergodic properties are strictly related to the behavior of the parameter process $f_t(\boldsymbol{\theta})$ and its derivative $\nabla f_t(\boldsymbol{\theta})$. Usually it suffices to assume that the both processes remains bounded, at least with a dominating probability. Root-n consistency can be shown under the identifiability condition that

$$n^{-1}\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = n^{-1}\sum_t \mathfrak{m}_t\big(\mu, f_t(\boldsymbol{\theta}), f_t(\boldsymbol{\theta}_0)\big) \ge \overline{\mathfrak{m}}(\boldsymbol{\theta}, \boldsymbol{\theta}_0)$$

where $\overline{\mathfrak{m}}(\boldsymbol{\theta}, \boldsymbol{\theta}_0)$ is two times continuously differentiable and satisfies $\overline{\mathfrak{m}}(\boldsymbol{\theta}, \boldsymbol{\theta}_0) > 0$ for $\boldsymbol{\theta} \ne \boldsymbol{\theta}_0$. Then Corollary 3.9 ensures root-$n$ consistency of the estimate $\widetilde{\boldsymbol{\theta}}$.

## 4.3 Estimation for Generalized Linear time series

Here we consider the case when $\mathcal{P}$ is again an exponential family with the canonical parameter and the parametric function $f_t(\boldsymbol{\theta})$ is a transformation of another function which linearly depends on the parameter $\boldsymbol{\theta}$. However, in the contrary to the previous example, we admit that the true conditional data distribution is not in $\mathcal{P}$.

A generalized linear specification for this model is given by the following set of structural equations:

$$\mathcal{L}\big(Y_t\big|\mathcal{F}_{t-1}\big) = P_{f_t} \in \mathcal{P}, \qquad f_t = g(\Psi_t), \qquad \Psi_t = A(\boldsymbol{\theta})\Psi_{t-1} \tag{4.3}$$

where $\Psi_t$ is a predictable $R^d$-dimensional explanatory process, $g(\cdot)$ is a given mapping from $I\!\!R^d$ to $I\!\!R$, $A(\boldsymbol{\theta})$ is a given $d \times d$-matrix linearly depending on the parameter vector $\boldsymbol{\theta} \in \Theta \subset I\!\!R^p$.

Let $\boldsymbol{\theta}$ be the parameter vector. Then, given $\boldsymbol{\theta}$, the observations $Y_1, \ldots, Y_t$, and the prehistory $\Psi^0, Y^0$, one can uniquely reconstruct the process $\Psi_t = \Psi_t(\boldsymbol{\theta})$ and then $f_t(\boldsymbol{\theta}) = g(\Psi_t(\boldsymbol{\theta}))$ for $t \geq 1$ by recurrently applying the relation (4.3). This function $f_t(\boldsymbol{\theta})$ leads to the (quasi) log-likelihood $L(\boldsymbol{\theta}) = \sum_t \big\{ Y_t f_t(\boldsymbol{\theta}) - d\big(f_t(\boldsymbol{\theta})\big) \big\}$.

A useful feature of models of type (4.3) is that the gradient $\nabla f_t(\boldsymbol{\theta})$ also follows the linear structural equation:

$$
\begin{aligned}
\nabla f_t &= \nabla \Psi_t(\boldsymbol{\theta}) g'(\Psi_t(\boldsymbol{\theta})), \\
\nabla \Psi_t(\boldsymbol{\theta}) &= \nabla A \cdot X_{t-1}(\boldsymbol{\theta}) + A(\boldsymbol{\theta}) \nabla X_{t-1}(\boldsymbol{\theta})
\end{aligned}
$$

where $g'(\cdot)$ means the gradient of $g(\cdot)$, and $\nabla A$ is the gradient of of $A(\boldsymbol{\theta})$ which does not depend on $\boldsymbol{\theta}$ because $A(\boldsymbol{\theta})$ is linear in $\boldsymbol{\theta}$.

Define

$$
\begin{aligned}
b_t &\stackrel{\text{def}}{=} \boldsymbol{E}\big[Y_t\big|\mathcal{F}_{t-1}\big], \\
m_t(v) &\stackrel{\text{def}}{=} \boldsymbol{E}\big[\ell(Y_t, v)\big|\mathcal{F}_{t-1}\big] = b_t v - d(v), \\
\zeta_t(Y_t, v) &\stackrel{\text{def}}{=} \boldsymbol{E}\big[\ell(Y_t, v) - m_t(v)\big|\mathcal{F}_{t-1}\big] = (Y_t - b_t)v, \\
\mathfrak{m}_t(\mu, v, v') &\stackrel{\text{def}}{=} \log \boldsymbol{E}\big[\exp\{\mu\ell(Y_t, v, v')\}\big|\mathcal{F}_{t-1}\big]. 
\end{aligned} \tag{4.4}
$$

All these quantities are computed for the underlying data distribution $\boldsymbol{P}$ for which both the assumption on the conditional distribution $Q_t \in \mathcal{P}$ and the parametric dynamics $f_t = f_t(\boldsymbol{\theta}_0)$ can be violated.

We suppose the following condition to be satisfied:

**(Yt)**   *There exist a constant $\lambda_1^* > 0$ and a predictable process $\mathfrak{n}_t$ such that for all $t$ and $\lambda \leq \lambda_1^*$*

$$\log \boldsymbol{E} \exp\big\{2\lambda\mathfrak{n}_t^{-1}(Y_t - b_t)\big|\mathcal{F}_{t-1}\big\} \leq 2\lambda^2.$$

This condition ensures that for all $v, v' \in \mathcal{U}$ and $\mu > 0$ with $\mu(v - v') \leq 2\lambda_1^*/\mathfrak{n}_t$ the quantity $\mathfrak{m}_t(\mu, v, v')$ from (4.4) is well defined.

The description of the process $\nabla\Psi_t(\boldsymbol{\theta})$ also yields the representation for $\nabla\zeta(\boldsymbol{\theta})$, $\nabla M(\boldsymbol{\theta})$, $B(\boldsymbol{\theta})$; see (3.5):

$$\nabla\zeta(\boldsymbol{\theta}) = \sum_t (Y_t - b_t)\nabla f_t(\boldsymbol{\theta}) = \sum_t \{Y_t - b_t\}\nabla\Psi_t(\boldsymbol{\theta})g'(\Psi_t(\boldsymbol{\theta})),$$

$$\nabla M(\boldsymbol{\theta}) = \sum_t \{b_t - \dot{d}(f_t(\boldsymbol{\theta}))\}\nabla\Psi_t(\boldsymbol{\theta})g'(\Psi_t(\boldsymbol{\theta})).$$

$$B(\boldsymbol{\theta}) = \sum_t \mathfrak{n}_t^2\,\nabla\Psi_t(\boldsymbol{\theta})g'(\Psi_t(\boldsymbol{\theta}))\{\nabla\Psi_t(\boldsymbol{\theta})g'(\Psi_t(\boldsymbol{\theta}))\}^\top.$$

This suggests to take the matrix $V(\boldsymbol{\theta})$ in the form

$$V(\boldsymbol{\theta}) = C_1\boldsymbol{E}B(\boldsymbol{\theta}) + C_2\,\mathrm{Var}\,\nabla M(\boldsymbol{\theta}) \tag{4.5}$$

with some $C_1, C_2 \geq 1$. Checking the condition $(ED)$ is quite straightforward in the most of situations and all the results of Sections 3, 3.2 apply.

We now specify the above expressions for the important special case of $d = 1$ when the model is given by the equations $f_t = g(X_t)$ for a univariate link function $g$ and a univariate process $X_t$ following the linear dynamic equation

$$X_t = \omega + \alpha_1 Y_{t-1} + \ldots \alpha_p Y_{t-p} + \beta_1 X_{t-1} + \ldots \beta_q X_{t-q}.$$

Here $\boldsymbol{\theta} = (\beta_1, \ldots, \beta_q, \alpha_1, \ldots, \alpha_p, \omega)^\top$. One easily computes

$$\begin{aligned}
\nabla X_t(\boldsymbol{\theta}) &= (1, Y_{t-1}, \ldots Y_{t-p}, X_{t-1}(\boldsymbol{\theta}), \ldots, X_{t-q}(\boldsymbol{\theta})) \\
&\quad + \left(0, \ldots, 0, \beta_1\frac{\partial X_{t-1}(\boldsymbol{\theta})}{\partial \beta_1}, \ldots, \beta_q\frac{\partial X_{t-q}(\boldsymbol{\theta})}{\partial \beta_q}\right).
\end{aligned}$$

So, given the initial values $X_{1-q}, \ldots, X_0$, $\nabla X_{1-q}, \ldots, \nabla X_0$ and the observations $Y_1, \ldots, Y_n$, one can recurrently construct for every $\boldsymbol{\theta}$ the hidden process $X_t(\boldsymbol{\theta})$ and its gradient $\nabla X_t(\boldsymbol{\theta})$. This yields the representation $f_t(\boldsymbol{\theta}) = g(X_t(\boldsymbol{\theta}))$ and $\nabla f_t(\boldsymbol{\theta}) = g'(X_t(\boldsymbol{\theta}))\nabla X_t(\boldsymbol{\theta})$ for the parameter process $f_t(\boldsymbol{\theta})$. The formulae for $\nabla\zeta(\boldsymbol{\theta})$, $\nabla M(\boldsymbol{\theta})$, $B(\boldsymbol{\theta})$ can also be specified:

$$\nabla\zeta(\boldsymbol{\theta}) = \sum_t (Y_t - b_t)\nabla f_t(\boldsymbol{\theta}) = \sum_t \{Y_t - b_t\}g'(X_t(\boldsymbol{\theta}))\nabla X_t(\boldsymbol{\theta}),$$

$$\nabla M(\boldsymbol{\theta}) = \sum_t \{b_t - \dot{d}(f_t(\boldsymbol{\theta}))\}g'(X_t(\boldsymbol{\theta}))\nabla X_t(\boldsymbol{\theta}).$$

$$B(\boldsymbol{\theta}) = \sum_t \mathfrak{n}_t^2\,|g'(X_t(\boldsymbol{\theta}))|^2\nabla X_t(\boldsymbol{\theta})\{\nabla X_t(\boldsymbol{\theta})\}^\top.$$

Further one can proceed as in the general case with $d > 1$.

25

## 4.4  GARCH(1,1) estimation

GARCH-model with the parameter $\boldsymbol{\theta} = (\omega, \alpha, \beta)^{\top}$ can be described by structural equation:

$$X_t(\boldsymbol{\theta}) = \omega + \alpha R_{t-1} + \beta X_{t-1}(\boldsymbol{\theta}), \qquad t \geq 1, \qquad X_0(\boldsymbol{\theta}) = X_0, \tag{4.6}$$

yielding the (quasi) log likelihood

$$
\begin{aligned}
L(\boldsymbol{\theta}) &= -\frac{1}{2}\sum_t \big\{ \log(2\pi X_t(\boldsymbol{\theta})) + R_t/X_t(\boldsymbol{\theta}) \big\} \\
&= \frac{1}{2}\sum_t \big\{ \log(2\pi f_t(\boldsymbol{\theta})) - R_t f_t(\boldsymbol{\theta}) \big\}
\end{aligned}
$$

with $f_t(\boldsymbol{\theta}) = 1/X_t(\boldsymbol{\theta})$.

Below we assume that the innovations $\varepsilon_t^2$'s have bounded conditional exponential moments: with some $\lambda^* > 0$ and $\mathfrak{n}$, it holds almost surely

$$\log \boldsymbol{E} \exp\big\{ -\frac{\lambda}{\mathfrak{n}} (\varepsilon_t^2 - 1) \big| \mathcal{F}_{t-1} \big\} \leq 2\lambda^2, \qquad |\lambda| \leq \lambda^*.$$

If every $\varepsilon_t$ is conditionally standard normal then $\log \boldsymbol{E} \exp\big\{ -\lambda(\varepsilon_t^2 - 1) \big| \mathcal{F}_{t-1} \big\} = \lambda - \frac{1}{2}\log(1 + 2\lambda) \leq 2\lambda^2$ for $|\lambda| \leq 1/3$.

As before, $\boldsymbol{\theta}_0 = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \boldsymbol{E} L(\boldsymbol{\theta})$ denotes the "true" parameter vector. By some technical reason we assume that the constant term $\omega$ of the parameter vector $\boldsymbol{\theta}$ is not smaller than given value $\delta > 0$. Then the equation (4.6) ensures for every $\boldsymbol{\theta}$ and every $t \geq 1$ the lower bound $X_t(\boldsymbol{\theta}) \geq \delta$.

The gradient $\nabla X_t(\boldsymbol{\theta})$ of the process $X_t(\boldsymbol{\theta})$ satisfies the equation

$$\nabla X_t(\boldsymbol{\theta}) = (1, R_{t-1}, X_{t-1}(\boldsymbol{\theta}))^{\top} + \beta \nabla X_{t-1}(\boldsymbol{\theta}).$$

For the canonical parameter $f_t(\boldsymbol{\theta}) = 1/X_t(\boldsymbol{\theta})$ this yields

$$\nabla f_t(\boldsymbol{\theta}) = X_t^{-2}(\boldsymbol{\theta}) \nabla X_t(\boldsymbol{\theta}) = f_t^2(\boldsymbol{\theta}) \nabla X_t(\boldsymbol{\theta}).$$

Next, one easily computes

$$
\begin{aligned}
\nabla \zeta(\boldsymbol{\theta}) &= -\frac{1}{2}\sum_t (R_t - X_t) \nabla f_t(\boldsymbol{\theta}), \\
\nabla M(\boldsymbol{\theta}) &= \frac{1}{2}\sum_t \big\{ X_t - X_t(\boldsymbol{\theta}) \big\} \nabla f_t(\boldsymbol{\theta}).
\end{aligned}
$$

With $\mathfrak{n}_t \equiv \mathfrak{n}$ and $Y_t = -R_t/2$, one has for $\lambda \leq \lambda^*$

$$\log \boldsymbol{E} \big[ \exp\big\{ 2\lambda \mathfrak{n}_t^{-1}(Y_t - b_t) \big\} \big| \mathcal{F}_{t-1} \big] = \log \boldsymbol{E} \big[ \exp(-\lambda \varepsilon_t^2 + \lambda) \big| \mathcal{F}_{t-1} \big] \leq 2\lambda^2$$

for $|\lambda| \leq \lambda_1^*$, and condition $(Yt)$ is verified. Next, condition $(ED)$ can be checked with $V(\boldsymbol{\theta})$ from (4.5), (3.5). The results of Section 3 and 3.2 describe the accuracy of estimation in terms of the function

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \;=\; \sum_t \mathfrak{m}_t\{\mu, f_t(\boldsymbol{\theta}), f_t(\boldsymbol{\theta}_0)\},$$

where for $v, v'$ with $|\mu X_t^{-1}(v - v')|\mathfrak{n} \leq \lambda^*/2$ that

$$\begin{aligned}
\mathfrak{m}_t\{\mu, v, v'\} \;\; &\stackrel{\text{def}}{=} \;\; -\log \boldsymbol{E}\big[\exp\{\mu\ell(R_t, v, v')\}\big|\mathcal{F}_{t-1}\big] \\
&= \;\; \frac{\mu}{2}\log(v'/v) + \frac{\mu(v - v')}{2X_t} + \frac{1}{2}\log \boldsymbol{E}\exp\{\mu X_t^{-1}(v - v')(\varepsilon_t^2 - 1)\} \\
&\geq \;\; \frac{\mu}{2}\log(v'/v) + \frac{\mu(v - v')}{2X_t} - \frac{\mu^2(v - v')^2\mathfrak{n}^2}{4X_t^2}
\end{aligned}$$

yielding

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq \sum_t \left\{\frac{\mu}{2}\log\left(\frac{f_t(\boldsymbol{\theta}_0)}{f_t(\boldsymbol{\theta})}\right) + \frac{\mu\{f_t(\boldsymbol{\theta}) - f_t(\boldsymbol{\theta}_0)\}}{2X_t} - \frac{\mu^2|f_t(\boldsymbol{\theta}) - f_t(\boldsymbol{\theta}_0)|^2\mathfrak{n}^2}{4X_t^2}\right\}$$

In the parametric situation $X_t = 1/f_t(\boldsymbol{\theta}_0)$ and with $\delta_t(\boldsymbol{\theta}) = f_t(\boldsymbol{\theta})/f_t(\boldsymbol{\theta}_0) - 1$

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) \geq \sum_t \left\{\frac{\mu}{2}\log(1 + \delta_t(\boldsymbol{\theta})) + \frac{\mu\delta_t(\boldsymbol{\theta})}{2} - \frac{\mu^2\delta_t(\boldsymbol{\theta})^2\mathfrak{n}^2}{4}\right\}$$

So, $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ can be viewed as a kind of distance between two functions $f_t(\boldsymbol{\theta})$ and $f_t(\boldsymbol{\theta}_0)$. In the stationary case $\alpha + \beta < 1$ for any $\boldsymbol{\theta} = (\omega, \alpha, \beta)^\top \in \Theta$, the process $X_t(\boldsymbol{\theta})$ is ergodic and the normalized sum $n^{-1}\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ converges to the integral of every summand w.r.t. the stationary measure. One can easily seen that this integral is nearly quadratic in $\boldsymbol{\theta} - \boldsymbol{\theta}_0$ in a neighborhood of $\boldsymbol{\theta}_0$ yielding the root-n consistency of estimation by Corollary 3.9.

## 4.5    Median and quantile time series estimation

The median or more generally quantile estimation can be defined by maximizing the quasi log-likelihood

$$L(\boldsymbol{\theta}) = \sum_t \ell_\alpha\big(Y_t - f_t(\boldsymbol{\theta})\big)$$

with $\ell_\alpha(x) = (1 - \alpha)x_- - \alpha x_+$. Define

$$\begin{aligned}
m_t(\boldsymbol{\theta}) \;\; &\stackrel{\text{def}}{=} \;\; \boldsymbol{E}\big\{\ell_\alpha\big(Y_t - f_t(\boldsymbol{\theta})\big)\big|\mathcal{F}_{t-1}\big\}, \\
q_t(\boldsymbol{\theta}) \;\; &\stackrel{\text{def}}{=} \;\; \boldsymbol{P}\big(Y_t - f_t(\boldsymbol{\theta}) \leq 0\big|\mathcal{F}_{t-1}\big), \\
M(\boldsymbol{\theta}) \;\; &\stackrel{\text{def}}{=} \;\; \sum_t m_t(\boldsymbol{\theta}) = \sum_t \boldsymbol{E}\big\{\ell_\alpha\big(Y_t - f_t(\boldsymbol{\theta})\big)\big|\mathcal{F}_{t-1}\big\}, \\
\zeta(\boldsymbol{\theta}) \;\; &\stackrel{\text{def}}{=} \;\; \sum_t \big\{\ell_\alpha\big(Y_t - f_t(\boldsymbol{\theta})\big) - m_t(\boldsymbol{\theta})\big\}.
\end{aligned}$$

Simple derivations yield

$$\nabla\zeta(\boldsymbol{\theta}) \;=\; \sum_t \big\{\mathbf{1}(Y_t - f_t(\boldsymbol{\theta}) \leq 0) - q_t(\boldsymbol{\theta})\big\}\nabla f_t(\boldsymbol{\theta}),$$

$$\nabla M(\boldsymbol{\theta}) \;=\; \sum_t q_t(\boldsymbol{\theta})\nabla f_t(\boldsymbol{\theta}).$$

Here it is natural to set $\mathfrak{n}_t(\upsilon) \equiv 1$ leading to

$$B(\boldsymbol{\theta}) = \sum_t \nabla f_t(\boldsymbol{\theta})\nabla f_t(\boldsymbol{\theta})^\top.$$

As previously, the condition of Theorem 3.3 are easy to verify with $V(\boldsymbol{\theta}) = C_1 \boldsymbol{E} B(\boldsymbol{\theta}) + C_2 \operatorname{Var}\nabla M(\boldsymbol{\theta})$ for some $C_1, C_2 \geq 1$. The function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ from (3.1) describing the quality of estimation can be represented as

$$\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0) = \sum_t \mathfrak{m}_t(\mu, f_t(\boldsymbol{\theta}), f_t(\boldsymbol{\theta}_0))$$

with

$$\mathfrak{m}_t(\mu, \upsilon, \upsilon') = -\log \boldsymbol{E}\big[\exp\big\{\mu\ell_\alpha(Y_t - \upsilon) - \mu\ell_\alpha(Y_t - \upsilon')\big\}\big|\mathcal{F}_{t-1}\big]$$

## 4.6  Risk bounds. Summary

In this section we summarize what we obtained in all the previous examples and what can be stated in each particular case.

One can see that in every application there is a straightforward expression for the quasi likelihood function $L(\boldsymbol{\theta})$ and its components $\zeta(\boldsymbol{\theta})$ and $M(\boldsymbol{\theta})$ and for their gradient in $\boldsymbol{\theta}$ in terms of $f_t(\boldsymbol{\theta})$ and $\nabla f_t(\boldsymbol{\theta})$. Moreover, all the conditions required for the main results of Sections 3 and 3.2 can be easily verified with the matrix $V(\boldsymbol{\theta})$ given in (4.5). Therefore, all these theorems and corollaries are applicable. In particular, the concentration property and the structure of confidence sets is given.

To understand what the results exactly imply in particular cases, one has to bound the penalty function $\operatorname{pen}(\boldsymbol{\theta})$ from above and the rate function $\mathfrak{M}(\mu, \boldsymbol{\theta}, \boldsymbol{\theta}_0)$ from below. The penalty function can be defined in the most of cases via the linear ranking $\|\sqrt{V(\boldsymbol{\theta})}(\boldsymbol{\theta} - \boldsymbol{\theta})\|$. The rate admits typically a quadratic in $\boldsymbol{\theta} - \boldsymbol{\theta}$ lower bound which immediately yields the classical root-n accuracy by Corollary 3.11. The general case of Sections 3 and 3.2 apply in each of considered example without any significant change of formulation.

# 5 A penalized exponential bound for a random field

Let $(\mathcal{Y}(\boldsymbol{v}), \boldsymbol{v} \in \varUpsilon)$ be a random field on a probability space $(\varOmega, \mathcal{F}, \boldsymbol{P})$, where $\varUpsilon$ is a separable locally compact space. For any $\boldsymbol{v} \in \varUpsilon$ we assume the following exponential moment condition to be fulfilled:

$(\mathcal{E})$   *For every* $\boldsymbol{v} \in \varUpsilon$

$$\boldsymbol{E} \exp\{\mathcal{Y}(\boldsymbol{v})\} = 1.$$

The aim of this section is to establish a similar exponential bound for a supremum of $\mathcal{Y}(\boldsymbol{v})$ over $\boldsymbol{v} \in \varUpsilon$. A trivial corollary of the condition $(\mathcal{E})$ is that if the set $\varUpsilon$ is finite with $N = \#\varUpsilon$, then

$$\boldsymbol{E} \exp\Big\{\sup_{\boldsymbol{v} \in \varUpsilon} \mathcal{Y}(\boldsymbol{v})\Big\} \leq N.$$

Unfortunately, in the general case the supremum of $\mathcal{Y}(\boldsymbol{v})$ over $\boldsymbol{v}$ does not necessarily fulfill the condition of bounded exponential moments. We therefore, consider a penalized version of the process $\mathcal{Y}(\boldsymbol{v})$, that is, we try to bound the exponential moment of $\mathcal{Y}(\boldsymbol{v}) -$ pen$(\boldsymbol{v})$ for some *penalty* function pen$(\boldsymbol{v})$. The goal is to find a possibly minimal such function pen$(\boldsymbol{v})$ which provides

$$\boldsymbol{E} \exp\Big\{\sup_{\boldsymbol{v} \in \varUpsilon}\big[\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})\big]\Big\} \leq 1.$$

In the case of a finite set $\varUpsilon$, a natural candidate is pen$(\boldsymbol{v}) = \log(\#\varUpsilon)$. Below we show how this simple choice can be extended to the case of a general set $\varUpsilon$. There exists a number of results about a supremum of a centered random field which are heavily based on the theory of empirical processes. See e.g. the monographes van der Vaart and Wellner (1996), Van de Geer (2000), Massart (2007), and references therein. Our approach is a bit different. First the process $\mathcal{Y}(\boldsymbol{v})$ does not need to be centered, instead we use the normalization $\boldsymbol{E} \exp\{\mathcal{Y}(\boldsymbol{v})\} = 1$. Secondly we do not assume any particular structure of this process like independence of observations, so the methods of the empirical processes do not apply here. Finally, our analysis is focuses on the penalty function pen$(\cdot)$ rather then on the deviation probability of $\max_{\boldsymbol{v}} \mathcal{Y}(\boldsymbol{v})$.

## 5.1   A local bound

Define $\mathcal{M}(\boldsymbol{v}) = \boldsymbol{E}\mathcal{Y}(\boldsymbol{v})$, $\zeta(\boldsymbol{v}) = \mathcal{Y}(\boldsymbol{v}) - \boldsymbol{E}\mathcal{Y}(\boldsymbol{v})$, and denote $\zeta(\boldsymbol{v}, \boldsymbol{v}') = \zeta(\boldsymbol{v}) - \zeta(\boldsymbol{v}')$ for $\boldsymbol{v}, \boldsymbol{v}' \in \varUpsilon$. We assume a nonnegative symmetric function $\mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}')$ is given such that the following condition is fulfilled:

**(𝓔𝜖)** *There exist numbers* $\epsilon > 0$ *and* $\lambda^* > 0$ *, such that for any* $\lambda \leq \lambda^*$

$$\sup_{\boldsymbol{v},\boldsymbol{v}'\in\Upsilon\,:\,\mathfrak{D}(\boldsymbol{v},\boldsymbol{v}')\leq\epsilon}\log\boldsymbol{E}\exp\Big\{2\lambda\frac{\zeta(\boldsymbol{v},\boldsymbol{v}')}{\mathfrak{D}(\boldsymbol{v},\boldsymbol{v}')}\Big\}\leq 2\lambda^2.$$

Let $\epsilon > 0$ be shown in condition $(𝓔𝜖)$. Define for any point $\boldsymbol{v}^\circ \in \Upsilon$ the "ball"

$$\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)=\big\{\boldsymbol{v}:\mathfrak{D}(\boldsymbol{v},\boldsymbol{v}^\circ)\leq\epsilon\big\}.$$

To state the result, we have to introduce the notion of *local entropy*. We say that a discrete set $\mathcal{D}(\epsilon,\mathcal{C})$ is an $\epsilon$-net in $\mathcal{C} \subseteq \Upsilon$, if

$$\mathcal{C}\subset\bigcup_{\boldsymbol{v}^\circ\in\mathcal{D}(\epsilon,\mathcal{C})}\mathcal{B}(\epsilon,\boldsymbol{v}^\circ).\tag{5.1}$$

By $\mathbb{N}(\epsilon_0,\epsilon,\boldsymbol{v}^\circ)$ for $\epsilon_0 \leq \epsilon$ we denote the local covering number defined as the minimal number of sets $\mathcal{B}(\epsilon_0,\cdot)$ required to cover $\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)$. With this covering number we associate the local entropy

$$\mathbb{Q}(\epsilon,\boldsymbol{v}^\circ)\stackrel{\text{def}}{=}\sum_{k=1}^{\infty}2^{-k}\log\mathbb{N}(2^{-k}\epsilon,\epsilon,\boldsymbol{v}^\circ).$$

Assume that $\boldsymbol{v}^\circ \in \Upsilon$ is fixed. The following result controls the supremum in $\boldsymbol{v}$ of the penalized process $\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})$ over the ball $\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)$.

**Theorem 5.1.** *Assume* $(𝓔)$ *and* $(𝓔𝜖)$. *For any* $\varrho \in (0,1)$ *with* $\varrho\epsilon/(1-\varrho) \leq \lambda^*$ *, any* $\boldsymbol{v}^\circ \in \Upsilon$

$$\log\boldsymbol{E}\exp\Big\{\sup_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)}\varrho\big[\mathcal{Y}(\boldsymbol{v})-\mathrm{pen}(\boldsymbol{v})\big]\Big\}\leq\frac{2\epsilon^2\varrho^2}{1-\varrho}+(1-\varrho)\mathbb{Q}(\epsilon,\boldsymbol{v}^\circ)-\varrho\,\mathrm{pen}_\epsilon(\boldsymbol{v}^\circ)$$

*with*

$$\mathrm{pen}_\epsilon(\boldsymbol{v}^\circ)=\inf_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)}\mathrm{pen}(\boldsymbol{v}).$$

*Proof.* We begin with some result which bounds the stochastic component of the process $\mathcal{Y}(\boldsymbol{v})$ within the local ball $\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)$.

**Lemma 5.2.** *Assume that* $\zeta(\boldsymbol{v})$ *is a separable process satisfying condition* $(𝓔𝜖)$. *Then for any given* $\boldsymbol{v}^\circ \in \Upsilon$ *,* $\boldsymbol{v}^\sharp \in \mathcal{B}(\epsilon,\boldsymbol{v}^\circ)$ *, and* $\lambda \leq \lambda^*$

$$\log\boldsymbol{E}\exp\Big\{\frac{\lambda}{\epsilon}\sup_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)}\zeta(\boldsymbol{v},\boldsymbol{v}^\sharp)\Big\}\leq\mathbb{Q}(\epsilon,\boldsymbol{v}^\circ)+2\lambda^2.$$

*Proof.* The proof is based on the standard chaining argument; see e.g. van der Vaart and Wellner (1996). Without loss of generality, we assume that $\mathbb{Q}(\epsilon, \boldsymbol{v}^\circ) < \infty$. Then for any integer $k \geq 0$, there exists a $2^{-k}\epsilon$-net $\mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)$ in the local ball $\mathcal{B}(\epsilon, \boldsymbol{v}^\circ)$ having the cardinality $\mathbb{N}(2^{-k}\epsilon, \epsilon, \boldsymbol{v}^\circ)$. Using the nets $\mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)$ with $k = 1, \ldots, K-1$, one can construct a chain connecting an arbitrary point $\boldsymbol{v}$ in $\mathcal{D}_K(\epsilon, \boldsymbol{v}^\circ)$ and $\boldsymbol{v}^\sharp$. It means that one can find points $\boldsymbol{v}_k \in \mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)$, $k = 1, \ldots, K-1$, such that $\mathfrak{D}(\boldsymbol{v}_k, \boldsymbol{v}_{k-1}) \leq 2^{-k+1}\epsilon$ for $k = 1, \ldots, K$. Here $\boldsymbol{v}_K$ means $\boldsymbol{v}$ and $\boldsymbol{v}_0$ means $\boldsymbol{v}^\sharp$. Notice that $\boldsymbol{v}_k$ can be constructed recurrently: $\boldsymbol{v}_{k-1} = \tau_{k-1}(\boldsymbol{v}_k)$, $k = K, \ldots, 1$, where

$$\tau_{k-1}(\boldsymbol{v}) = \underset{\boldsymbol{v}' \in \mathcal{D}_{k-1}(\epsilon, \boldsymbol{v}^\circ)}{\operatorname{argmin}} \mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}').$$

It obviously holds

$$\zeta(\boldsymbol{v}, \boldsymbol{v}^\sharp) = \sum_{k=1}^{K} \zeta(\boldsymbol{v}_k, \boldsymbol{v}_{k-1}).$$

It holds for $\xi(\boldsymbol{v}_k, \boldsymbol{v}_{k-1}) = \zeta(\boldsymbol{v}_k, \boldsymbol{v}_{k-1})/\mathfrak{D}(\boldsymbol{v}_k, \boldsymbol{v}_{k-1})$ that

$$\zeta(\boldsymbol{v}_k, \boldsymbol{v}_{k-1}) = \mathfrak{D}(\boldsymbol{v}_k, \boldsymbol{v}_{k-1})\xi(\boldsymbol{v}_k, \boldsymbol{v}_{k-1}) = 2\epsilon\, c_k\, \xi(\boldsymbol{v}_k, \boldsymbol{v}_{k-1})$$

with $c_k = \mathfrak{D}(\boldsymbol{v}_k, \boldsymbol{v}_{k-1})/(2\epsilon) \leq 2^{-k}$. By condition $(\mathcal{E}\epsilon)$ $\log \boldsymbol{E} \exp\{2\lambda\xi(\boldsymbol{v}_k, \boldsymbol{v}_{k-1})\} \leq 2\lambda^2$. Next,

$$
\begin{aligned}
\sup_{\boldsymbol{v} \in \mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)} \zeta(\boldsymbol{v}, \boldsymbol{v}^\sharp) &\leq \sum_{k=1}^{K} \sup_{\boldsymbol{v}' \in \mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)} \zeta(\boldsymbol{v}', \tau_{k-1}(\boldsymbol{v}')) \\
&\leq 2\epsilon \sum_{k=1}^{K} \sup_{\boldsymbol{v}' \in \mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)} c_k \xi(\boldsymbol{v}', \tau_{k-1}(\boldsymbol{v}')). \quad (5.2)
\end{aligned}
$$

Since $c_k \leq 2^{-k}$, the Hölder inequality and condition $(\mathcal{E}\epsilon)$ imply

$$\log \boldsymbol{E} \exp\left\{\frac{\lambda}{\epsilon} \sup_{\boldsymbol{v} \in \mathcal{D}_K(\epsilon, \boldsymbol{v}^\circ)} \zeta(\boldsymbol{v}, \boldsymbol{v}^\sharp)\right\} \leq \log \boldsymbol{E} \exp\left\{2\lambda \sum_{k=1}^{K} \sup_{\boldsymbol{v}' \in \mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)} c_k \xi(\boldsymbol{v}', \tau_{k-1}(\boldsymbol{v}'))\right\}$$

$$\leq \sum_{k=1}^{K} 2^{-k} \log\left[\boldsymbol{E} \exp\left\{\sup_{\boldsymbol{v}' \in \mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)} 2^k c_k \times 2\lambda\xi(\boldsymbol{v}', \tau_{k-1}(\boldsymbol{v}'))\right\}\right]$$

$$\leq \sum_{k=1}^{K} 2^{-k} \log\left[\sum_{\boldsymbol{v}' \in \mathcal{D}_k(\epsilon, \boldsymbol{v}^\circ)} \boldsymbol{E} \exp\left\{2^k c_k \times 2\lambda\xi(\boldsymbol{v}', \tau_{k-1}(\boldsymbol{v}'))\right\}\right]$$

$$\leq \sum_{k=1}^{K} 2^{-k}\left\{\log \mathbb{N}(2^{-k}\epsilon, \epsilon, \boldsymbol{v}^\circ) + 2\lambda^2\right\}.$$

These inequalities and the separability of $\zeta(\boldsymbol{v}, \boldsymbol{v}^\sharp)$ yield

$$\log \boldsymbol{E} \exp\left\{\frac{\lambda}{\epsilon} \sup_{\boldsymbol{v} \in \mathcal{B}(\epsilon, \boldsymbol{v}^\circ)} \zeta(\boldsymbol{v}, \boldsymbol{v}^\sharp)\right\} = \lim_{K \to \infty} \log \boldsymbol{E} \exp\left\{\frac{\lambda}{\epsilon} \sup_{\boldsymbol{v} \in \mathcal{D}_K(\epsilon, \boldsymbol{v}^\circ)} \zeta(\boldsymbol{v}, \boldsymbol{v}^\sharp)\right\}$$

$$\leq \sum_{k=1}^{\infty} 2^{-k}\left\{2\lambda^2 + \log \mathbb{N}(2^{-k}\epsilon, \epsilon, \boldsymbol{v}^\circ)\right\} \leq 2\lambda^2 + \mathbb{Q}(\epsilon, \boldsymbol{v}^\circ)$$

which completes the proof of the lemma. $\qquad\square$

Now define for a fixed a point $\boldsymbol{v}^\circ$

$$\boldsymbol{v}^\sharp = \operatorname*{argmin}_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)} \{\mathcal{M}(\boldsymbol{v}) + \operatorname{pen}(\boldsymbol{v})\},$$

where $\mathcal{M}(\boldsymbol{v}) = -\boldsymbol{E}\mathcal{Y}(\boldsymbol{v})$. If there are many such points, then take any of them as $\boldsymbol{v}^\sharp$. Obviously

$$\sup_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)} \big\{\mathcal{Y}(\boldsymbol{v}) - \operatorname{pen}(\boldsymbol{v})\big\} \le \mathcal{Y}(\boldsymbol{v}^\sharp) - \operatorname{pen}(\boldsymbol{v}^\sharp) + \sup_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)} \zeta(\boldsymbol{v},\boldsymbol{v}^\sharp).$$

Therefore, by the Hölder inequality and Lemma 5.2 with $\lambda = \epsilon\varrho/(1-\varrho)$

$$\log \boldsymbol{E}\exp\Big\{\sup_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)} \varrho\big[\mathcal{Y}(\boldsymbol{v}) - \operatorname{pen}(\boldsymbol{v})\big]\Big\}$$

$$\le \varrho\log\boldsymbol{E}\exp\big\{\mathcal{Y}(\boldsymbol{v}^\sharp) - \operatorname{pen}(\boldsymbol{v}^\sharp)\big\} + (1-\varrho)\log\boldsymbol{E}\exp\Big\{\frac{\varrho}{1-\varrho}\sup_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)} \zeta(\boldsymbol{v},\boldsymbol{v}^\sharp)\Big\}$$

$$\le 2\epsilon^2\varrho^2/(1-\varrho) + (1-\varrho)\mathbb{Q}(\epsilon,\boldsymbol{v}^\circ) - \varrho\operatorname{pen}(\boldsymbol{v}^\sharp)$$

$$\le 2\epsilon^2\varrho^2/(1-\varrho) + (1-\varrho)\mathbb{Q}(\epsilon,\boldsymbol{v}^\circ) - \varrho\operatorname{pen}_\epsilon(\boldsymbol{v}^\circ).$$

which is the assertion of the theorem. $\qquad\square$

## 5.2 A global exponential bound for the penalized process

This section presents some sufficient conditions on the penalty function $\operatorname{pen}(\boldsymbol{v})$ which ensure the general exponential bound for the penalized process $\mathcal{Y}(\boldsymbol{v}) - \operatorname{pen}(\boldsymbol{v})$. For simplicity we assume that the local entropy numbers $\mathbb{Q}(\epsilon,\boldsymbol{v})$ are uniformly bounded by a constant $\mathbb{Q}^*(\Upsilon)$. Let also $\pi$ be a $\sigma$-finite measure on the space $\Upsilon$ and $\pi(A)$ stand for the $\pi$-measure of a set $A \subset \Upsilon$. The standard proposal for $\pi$ is the usual Lebesgue measure.

**Theorem 5.3.** *Assume* $(\mathcal{E})$ *and* $(\mathcal{E}_\epsilon)$ *with some fixed* $\epsilon$ *and* $\lambda^*$. *Let* $\varrho < 1$ *be such that* $\varrho\epsilon/(1-\varrho) \le \lambda^*$. *Let also* $\mathbb{Q}(\epsilon,\boldsymbol{v}) \le \mathbb{Q}^*(\Upsilon)$ *for all* $\boldsymbol{v} \in \Upsilon$. *Let a* $\sigma$-*finite measure* $\pi$ *on* $\Upsilon$ *be such that for some* $\nu \ge 1$

$$\sup_{\boldsymbol{v},\boldsymbol{v}':\mathfrak{D}(\boldsymbol{v},\boldsymbol{v}')\le\epsilon} \frac{\pi(\mathcal{B}(\epsilon,\boldsymbol{v}))}{\pi(\mathcal{B}(\epsilon,\boldsymbol{v}'))} \le \nu. \tag{5.3}$$

*Finally, let a function* $\operatorname{pen}(\boldsymbol{v})$ *satisfy*

$$\mathfrak{H}_\epsilon(\varrho) \stackrel{\text{def}}{=} \log\int_\Upsilon \frac{1}{\pi(\mathcal{B}(\epsilon,\boldsymbol{v}^\circ))}\exp\big\{-\varrho\operatorname{pen}_\epsilon(\boldsymbol{v}^\circ)\big\}d\pi(\boldsymbol{v}^\circ) < \infty$$

with $\mathrm{pen}_\epsilon(\boldsymbol{v}^\circ) = \inf_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)} \mathrm{pen}(\boldsymbol{v})$. *Then*

$$\boldsymbol{E}\exp\Big\{\sup_{\boldsymbol{v}\in\varUpsilon}\varrho\big[\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})\big]\Big\} \le \mathfrak{Q}(\varrho,\epsilon), \tag{5.4}$$

*where*

$$\log\mathfrak{Q}(\varrho,\epsilon) = \frac{2\epsilon^2\varrho^2}{1-\varrho} + (1-\varrho)\mathbb{Q}^*(\varUpsilon) + \log\nu + \mathfrak{H}_\epsilon(\varrho). \tag{5.5}$$

*Proof.* We begin with a simple technical result which bounds the maximum of a given function via the weighted integral of the local maxima.

**Lemma 5.4.** *Let* $f(\boldsymbol{v})$ *be a nonnegative function on* $\varUpsilon \subset I\!\!R^p$ *and let for every point* $\boldsymbol{v} \in \varUpsilon$ *a vicinity* $A(\boldsymbol{v})$ *be fixed such that* $\boldsymbol{v}' \in A(\boldsymbol{v})$ *implies* $\boldsymbol{v} \in A(\boldsymbol{v}')$. *Let also the measure* $\pi\big(A(\boldsymbol{v})\big)$ *of the set* $A(\boldsymbol{v})$ *fulfill for every* $\boldsymbol{v}^\circ \in \varUpsilon$

$$\sup_{\boldsymbol{v}\in A(\boldsymbol{v}^\circ)} \frac{\pi\big(A(\boldsymbol{v})\big)}{\pi\big(A(\boldsymbol{v}^\circ)\big)} \le \nu. \tag{5.6}$$

*Then*

$$\sup_{\boldsymbol{v}\in\varUpsilon} f(\boldsymbol{v}) \le \nu \int_\varUpsilon f^*(\boldsymbol{v})\frac{1}{\pi\big(A(\boldsymbol{v})\big)}d\pi(\boldsymbol{v})$$

*with*

$$f^*(\boldsymbol{v}) \stackrel{\mathrm{def}}{=} \sup_{\boldsymbol{v}'\in A(\boldsymbol{v})} f(\boldsymbol{v}').$$

*Proof.* For every $\boldsymbol{v}^\circ \in \varUpsilon$

$$\begin{aligned}
\int_\varUpsilon f^*(\boldsymbol{v})\frac{1}{\pi\big(A(\boldsymbol{v})\big)}d\pi(\boldsymbol{v}) &\ge \int_{A(\boldsymbol{v}^\circ)} f^*(\boldsymbol{v})\frac{1}{\pi\big(A(\boldsymbol{v})\big)}d\pi(\boldsymbol{v}) \\
&\ge f(\boldsymbol{v}^\circ)\int_{A(\boldsymbol{v}^\circ)} \frac{1}{\pi\big(A(\boldsymbol{v})\big)}d\pi(\boldsymbol{v})
\end{aligned}$$

because $\boldsymbol{v} \in A(\boldsymbol{v}^\circ)$ implies $\boldsymbol{v}^\circ \in A(\boldsymbol{v})$ and hence, $f(\boldsymbol{v}^\circ) \le f^*(\boldsymbol{v})$. Now by (5.6)

$$\int_\varUpsilon f^*(\boldsymbol{v})\frac{1}{\pi\big(A(\boldsymbol{v})\big)}d\pi(\boldsymbol{v}) \ge \frac{f(\boldsymbol{v}^\circ)}{\nu}\int_{A(\boldsymbol{v}^\circ)} \frac{1}{\pi\big(A(\boldsymbol{v}^\circ)\big)}d\pi(\boldsymbol{v}) = f(\boldsymbol{v}^\circ)/\nu$$

as required. $\qquad\square$

This result applied to $f(\boldsymbol{v}) = \exp\big\{\varrho\big[\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})\big]\big\}$ and $A(\boldsymbol{v}) = \mathcal{B}(\epsilon,\boldsymbol{v})$ implies

$$\sup_{\boldsymbol{v}\in\varUpsilon}\exp\big\{\varrho\big[\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})\big]\big\} \le \nu\int_\varUpsilon \sup_{\boldsymbol{v}\in\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)}\exp\big\{\varrho\big[\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})\big]\big\}\frac{d\pi(\boldsymbol{v}^\circ)}{\pi\big(\mathcal{B}(\epsilon,\boldsymbol{v}^\circ)\big)}.$$

This implies by Theorem 5.1

$$\log \boldsymbol{E} \sup_{\boldsymbol{v} \in \Upsilon} \exp\Big\{ \varrho\big[\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})\big] \Big\}$$

$$\leq \frac{2\epsilon^2 \varrho^2}{1 - \varrho} + (1 - \varrho)\mathbb{Q}^*(\Upsilon) + \log\Big\{ \nu \int_\Upsilon \exp\{-\varrho \, \mathrm{pen}_\epsilon(\boldsymbol{v}^\circ)\} \frac{d\pi(\boldsymbol{v}^\circ)}{\pi\big(\mathcal{B}(\epsilon, \boldsymbol{v}^\circ)\big)} \Big\}$$

$$\leq \frac{2\epsilon^2 \varrho^2}{1 - \varrho} + (1 - \varrho)\mathbb{Q}^*(\Upsilon) + \log(\nu) + \mathfrak{H}_\epsilon(\varrho)$$

and the assertion follows. □

## 5.3 Smooth case

Here we discuss the special case when $\Upsilon \subset {I\!\!R}^p$, the process $\mathcal{Y}(\boldsymbol{v})$ and its stochastic component $\zeta(\boldsymbol{v}) \stackrel{\mathrm{def}}{=} d\zeta(\boldsymbol{v})/d\boldsymbol{v}$ has bounded exponential moments. We also assume that $\pi$ is the Lebesgue measure on $\Upsilon$. Suppose the following condition is fulfilled:

**($\mathcal{E}D$)** *There exist $\lambda^* > 0$ and for each $\boldsymbol{v} \in \Upsilon$, a symmetric non-negative matrix $H(\boldsymbol{v})$ such that for any $\lambda \leq \lambda^*$*

$$\sup_{\boldsymbol{v} \in \Upsilon} \sup_{\boldsymbol{\gamma} \in S^p} \log \boldsymbol{E} \exp\Big\{ 2\lambda \frac{\boldsymbol{\gamma}^\top \nabla \zeta(\boldsymbol{v})}{\|H(\boldsymbol{v})\boldsymbol{\gamma}\|} \Big\} \leq 2\lambda^2.$$

The matrix function $H(\boldsymbol{v})$ can be used for defining a natural topology in $\Upsilon$. Namely, for any $\boldsymbol{v}, \boldsymbol{v}' \in \Upsilon$ define $\mathfrak{d} = \|\boldsymbol{v} - \boldsymbol{v}'\|$, $\boldsymbol{\gamma} = (\boldsymbol{v} - \boldsymbol{v}')/\mathfrak{d}$ and

$$\mathfrak{D}^2(\boldsymbol{v}, \boldsymbol{v}') \stackrel{\mathrm{def}}{=} \|\boldsymbol{v} - \boldsymbol{v}'\|^2 \int_0^1 \boldsymbol{\gamma}^\top H^2(\boldsymbol{v} + t\mathfrak{d}\boldsymbol{\gamma})\boldsymbol{\gamma} \, dt.$$

Next, introduce for each $\boldsymbol{v}^\circ \in \Upsilon$ and $\epsilon > 0$ the set

$$\mathcal{B}(\epsilon, \boldsymbol{v}^\circ) \stackrel{\mathrm{def}}{=} \{\boldsymbol{v} : \mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}^\circ) \leq \epsilon\}$$

To state the result, we need one more condition on the uniform continuity of the matrix $H(\boldsymbol{v})$ in $\boldsymbol{v}$.

**($H$)** *There exist constants $\epsilon > 0$ and $\nu_1 \geq 1$ such that*

$$\sup_{\boldsymbol{v}, \boldsymbol{v}' : \mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}') \leq \epsilon} \sup_{\boldsymbol{\gamma} \in S^p} \frac{\boldsymbol{\gamma}^\top H^2(\boldsymbol{v})\boldsymbol{\gamma}}{\boldsymbol{\gamma}^\top H^2(\boldsymbol{v}')\boldsymbol{\gamma}} \leq \nu_1 \, .$$

**Theorem 5.5.** *Let $(\mathcal{E})$ be satisfied. Suppose that $(\mathcal{E}D)$ holds with some $\lambda^*$ and a matrix function $H(\boldsymbol{v})$ which fulfills $(H)$. If for some $\varrho \in (0, 1)$ and $\epsilon > 0$ with $\varrho\epsilon/(1-\varrho) \leq \lambda^*$, the penalty function $\mathrm{pen}(\boldsymbol{v})$ fulfills*

$$\mathfrak{H}_\epsilon(\varrho) \stackrel{\mathrm{def}}{=} \log\Big\{ \omega_p^{-1} \epsilon^{-p} \int_\Upsilon \det(H(\boldsymbol{v}^\circ)) \exp\{-\varrho \, \mathrm{pen}_\epsilon(\boldsymbol{v}^\circ)\} d\boldsymbol{v}^\circ \Big\} < \infty$$

*with* $\operatorname{pen}_\epsilon(\boldsymbol{v}^\circ) = \inf_{\boldsymbol{v} \in \mathcal{B}(\epsilon, \boldsymbol{v}^\circ)} \operatorname{pen}(\boldsymbol{v})$ *, then*

$$\boldsymbol{E} \exp\left\{\sup_{\boldsymbol{v} \in \Upsilon} \varrho\big[\mathcal{Y}(\boldsymbol{v}) - \operatorname{pen}(\boldsymbol{v})\big]\right\} \leq \mathfrak{Q}(\varrho, \epsilon) \tag{5.7}$$

*where*

$$\log \mathfrak{Q}(\varrho, \epsilon) \;=\; \frac{2\epsilon^2 \varrho^2}{1 - \varrho} + (1 - \varrho)\mathbb{Q}_p + \mathfrak{H}_\epsilon(\varrho) + p\log(\nu_1)$$

*with* $\mathbb{Q}_p$ *being the usual entropy number for the Euclidean ball in* $\mathrm{I\!R}^p$ *.*

*Proof.* First we show that the differentiability condition $(\mathcal{E}D)$ implies the local moment condition $(\mathcal{E}\epsilon)$.

**Lemma 5.6.** *Assume that* $(\mathcal{E}D)$ *holds with some* $\lambda^*$ *. Then for any* $\boldsymbol{v}^\circ \in \Upsilon$ *and any* $\lambda$ *with* $|\lambda| \leq \lambda^*/\nu_1^{1/2}$ *, it holds*

$$\sup_{\boldsymbol{v} \in \mathcal{B}(\epsilon, \boldsymbol{v}^\circ)} \log \boldsymbol{E} \exp\left\{2\lambda \frac{\zeta(\boldsymbol{v}, \boldsymbol{v}^\circ)}{\mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}^\circ)}\right\} \leq 2\lambda^2. \tag{5.8}$$

*Proof.* For $\boldsymbol{v} \in \mathcal{B}(\epsilon, \boldsymbol{v}^\circ)$, denote $\mathfrak{d} = \|\boldsymbol{v} - \boldsymbol{v}^\circ\|$, $\boldsymbol{\gamma} = (\boldsymbol{v} - \boldsymbol{v}^\circ)/\mathfrak{d}$. With this notation

$$\zeta(\boldsymbol{v}, \boldsymbol{v}^\circ) = \mathfrak{d}\boldsymbol{\gamma}^\top \int_0^1 \nabla \zeta(\boldsymbol{v}^\circ + t\mathfrak{d}\boldsymbol{\gamma}) dt.$$

The condition $(H)$ implies for every $t \in [0, 1]$ that

$$\lambda \frac{u\|H(\boldsymbol{v}^\circ + t\mathfrak{d}\boldsymbol{\gamma})\boldsymbol{\gamma}\|}{\mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}^\circ)} \leq \lambda \nu_1^{1/2} \leq \lambda^*.$$

Now the Hölder inequality and $(\mathcal{E}D)$ yield

$$
\begin{aligned}
&\log \boldsymbol{E} \exp\left\{2\lambda \frac{\zeta(\boldsymbol{v}, \boldsymbol{v}^\circ)}{\mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}^\circ)} - \lambda^2\right\} \\
&= \log \boldsymbol{E} \exp\left\{\int_0^1 \boldsymbol{\gamma}^\top \left[\frac{2\lambda\mathfrak{d}}{\mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}^\circ)} \nabla \zeta(\boldsymbol{v}^\circ + t\mathfrak{d}\boldsymbol{\gamma}) - \frac{2\lambda^2\mathfrak{d}^2}{\mathfrak{D}^2(\boldsymbol{v}, \boldsymbol{v}^\circ)} H^2(\boldsymbol{v}^\circ + t\mathfrak{d}\boldsymbol{\gamma})\boldsymbol{\gamma}\right] dt\right\} \\
&\leq \int_0^1 \log \boldsymbol{E} \exp\left\{\boldsymbol{\gamma}^\top \left[\frac{2\lambda\mathfrak{d}}{\mathfrak{D}(\boldsymbol{v}, \boldsymbol{v}^\circ)} \nabla \zeta(\boldsymbol{v}^\circ + t\mathfrak{d}\boldsymbol{\gamma}) - \frac{2\lambda^2\mathfrak{d}^2}{\mathfrak{D}^2(\boldsymbol{v}, \boldsymbol{v}^\circ)} H^2(\boldsymbol{v}^\circ + t\mathfrak{d}\boldsymbol{\gamma})\boldsymbol{\gamma}\right]\right\} dt \\
&\leq 0
\end{aligned}
$$

as required. $\qquad\square$

Next we show that condition $(H)$ implies (5.3). Consider for every $\boldsymbol{v}^\circ \in \Upsilon$ an elliptic neighborhood $\mathcal{B}'(\epsilon, \boldsymbol{v}^\circ) = \{\boldsymbol{v} : \|H(\boldsymbol{v}^\circ)(\boldsymbol{v} - \boldsymbol{v}^\circ)\| \leq \epsilon\}$.

**Lemma 5.7.** *Assume* $(H)$ *. Then*

*1. for any $\epsilon > 0$ and any $\boldsymbol{v} \in \Upsilon$*

$$\begin{aligned}
\mathcal{B}'(\nu_1^{-1/2}\epsilon, \boldsymbol{v}) &\subset \mathcal{B}(\epsilon, \boldsymbol{v}) \subset \mathcal{B}'(\nu_1^{1/2}\epsilon, \boldsymbol{v}), \\
\mathcal{B}(\nu_1^{-1/2}\epsilon, \boldsymbol{v}) &\subset \mathcal{B}'(\epsilon, \boldsymbol{v}) \subset \mathcal{B}(\nu_1^{1/2}\epsilon, \boldsymbol{v}).
\end{aligned} \tag{5.9}$$

*2. For every $\boldsymbol{v} \in \Upsilon$,*

$$\nu_1^{-p/2} \leq \epsilon^{-p} \pi(\mathcal{B}(\epsilon, \boldsymbol{v})) \det(H(\boldsymbol{v}))/\omega_p \leq \nu_1^{p/2}, \tag{5.10}$$

*where $\omega_p$ is the Lebesgue measure of the unit ball in $\mathbb{R}^p$.*

*3. condition (5.3) holds with $\nu = \nu_1^p$.*

*Proof.* Condition $(H)$ implies that for any $\boldsymbol{v}^\circ \in \Upsilon$ and $\boldsymbol{v} \in \mathcal{B}(\epsilon, \boldsymbol{v}^\circ)$ that

$$\nu_1^{-1} \boldsymbol{\gamma}^\top H^2(\boldsymbol{v}^\circ) \boldsymbol{\gamma} \leq \int_0^1 \boldsymbol{\gamma}^\top H^2(\boldsymbol{v}^\circ + t\mathfrak{d}\boldsymbol{\gamma}) \boldsymbol{\gamma} \, dt \leq \nu_1 \boldsymbol{\gamma}^\top H^2(\boldsymbol{v}^\circ) \boldsymbol{\gamma}$$

with $\mathfrak{d} = \|\boldsymbol{v} - \boldsymbol{v}^\circ\|$ and $\boldsymbol{\gamma} = (\boldsymbol{v} - \boldsymbol{v}^\circ)/\mathfrak{d}$, which yields the first assertion of the lemma.

The Lebesgue measure of the ellipsoid $\mathcal{B}'(\epsilon, \boldsymbol{v})$ is equal to $\omega_p \epsilon^p / \det(H(\boldsymbol{v}))$. This and (5.9) imply the second assertion. This, in turns, implies (5.3) in view of $(H)$. $\qquad\square$

The next result claims that in the smooth case the local entropy number $\mathbb{Q}(\epsilon, \boldsymbol{v}^\circ)$ is similar to the usual Euclidean situation.

**Lemma 5.8.** *Assume $(H)$. Then $\sup_{\boldsymbol{v} \in \Theta} \mathbb{Q}(\epsilon, \boldsymbol{v}) \leq \mathbb{Q}_p + p\log(\nu_1)$.*

*Proof.* Fix any $\boldsymbol{v}^\circ \in \Upsilon$. Linear transformation with the matrix $H^{-1}(\boldsymbol{v}^\circ)$ reduces the situation to the case when $H(\boldsymbol{v}^\circ) \equiv I$ and $\mathcal{B}'(\epsilon_0, \boldsymbol{v}^\circ)$ is a usual Euclidean ball for any $\epsilon_0 \leq \epsilon$. Moreover, by $(H)$, each elliptic set $\mathcal{B}'(\epsilon_0, \boldsymbol{v})$ for $\boldsymbol{v} \in \mathcal{B}(\epsilon, \boldsymbol{v}^\circ)$ is nearly an Euclidean ball in the sense that the ratio of its largest and smallest axes (which is the ratio of the largest and smallest eigenvalues of $H^{-1}(\boldsymbol{v}^\circ) H^2(\boldsymbol{v}) H^{-1}(\boldsymbol{v}^\circ)$) is bounded by $\nu_1$. Therefore, for any $\epsilon_0 \leq \epsilon$, a Euclidean net $\mathcal{D}^e(\epsilon_0/\nu_1)$ with the step $\epsilon_0/\nu_1$ ensures a covering of $\mathcal{B}(\epsilon, \boldsymbol{v}^\circ)$ by the sets $\mathcal{B}(\epsilon_0, \boldsymbol{v}^\circ)$, $\boldsymbol{v}^\circ \in \mathcal{D}^e(\epsilon_0/\nu_1)$. Therefore, the corresponding covering number is bounded by $(\nu_1 \epsilon/\epsilon_0)^p$ yielding the claimed bound for the local entropy. $\qquad\square$

Now the result of theorem 5.5 is reduced to the statement of Theorem 5.3. $\qquad\square$

Computing of the penalty simplifies a lot when the matrix $H(\boldsymbol{v})$ is uniformly bounded by a matrix $H^*$, or, equivalently, condition $(H)$ is fulfilled for $H(\boldsymbol{v}) \equiv H^*$. Then one can define $\operatorname{pen}(\boldsymbol{v})$ as a function of the norm $\|H^*(\boldsymbol{v} - \boldsymbol{v}_0)\|$ for a fixed $\boldsymbol{v}_0$.

**Theorem 5.9.** *Assume additionally to the conditions of Theorem 5.5 that $H(\boldsymbol{v}) \leq H^*$ for a symmetric matrix $H^*$. Suppose that $\varkappa(t)$ is a monotonously decreasing positive function on $[0, +\infty)$ satisfying*

$$\mathfrak{P}^* \stackrel{\text{def}}{=} \omega_p^{-1} \int_{I\!\!R^p} \varkappa(\|\boldsymbol{u}\|) d\boldsymbol{u} = p \int_0^\infty \varkappa(t) t^{p-1} dt < \infty. \tag{5.11}$$

*Define*

$$\mathrm{pen}(\boldsymbol{v}) = -\varrho^{-1} \log \varkappa\big(\epsilon^{-1}\|H^*(\boldsymbol{v} - \boldsymbol{v}_0)\| + 1\big)$$

*Then*

$$\boldsymbol{E} \exp\Big\{\sup_{\boldsymbol{v} \in \varUpsilon} \varrho\big[\mathcal{Y}(\boldsymbol{v}) - \mathrm{pen}(\boldsymbol{v})\big]\Big\} \leq \mathfrak{Q}(\varrho, \epsilon) \tag{5.12}$$

*with*

$$\log \mathfrak{Q}(\varrho, \epsilon) \quad = \quad \frac{2\epsilon^2 \varrho^2}{1 - \varrho} + (1 - \varrho)\mathbb{Q}_p + \log(\mathfrak{P}^*),$$

*where $\omega_p$ is the volume of the unit ball in $I\!\!R^p$.*

*Proof.* Let us fix $\boldsymbol{v}^\circ \in \varUpsilon$. Definition of the semi-metric $\mathfrak{D}$ and condition $(H)$ imply for every $\boldsymbol{v} \in \mathcal{B}(\epsilon, \boldsymbol{v}^\circ)$ that

$$\|H^*(\boldsymbol{v}^\circ - \boldsymbol{v})\| \leq \epsilon.$$

The triangle inequality and $(H)$ now imply for this $\boldsymbol{v}$ that

$$\epsilon^{-1}\|H^*(\boldsymbol{v} - \boldsymbol{v}_0)\| + 1 \geq \epsilon^{-1}\|H^*(\boldsymbol{v}^\circ - \boldsymbol{v}_0)\|$$

and $\mathrm{pen}_\epsilon(\boldsymbol{v}^\circ) \geq -\varrho^{-1} \log \varkappa\big(\epsilon^{-1}\|H^*(\boldsymbol{v}^\circ - \boldsymbol{v}_0)\|\big)$. Therefore, it follows by change of variables $\boldsymbol{u} = \epsilon H^*(\boldsymbol{v} - \boldsymbol{v}_0)$ that

$$\omega_p^{-1} \epsilon^{-p} \int_{\varUpsilon} \det(H^*) \exp\{-\varrho \, \mathrm{pen}_\epsilon(\boldsymbol{v})\} d\boldsymbol{v} \quad \leq \quad \omega_p^{-1} \int_{I\!\!R^p} \varkappa(\|\boldsymbol{u}\|) d\boldsymbol{u}$$

$$\leq \quad p \int_0^\infty \varkappa(t) t^{p-1} dt = \mathfrak{P}^*,$$

and the result follows from Theorem 5.5. $\qquad \square$

Natural candidates for the function $\varkappa(\cdot)$ and the corresponding $\mathfrak{P}^*$-values are:

$$\varkappa_1(t) \quad = \quad e^{-\delta_1(t-1)_+^2}, \quad \mathfrak{P}_1^* \quad = \quad 1 + \omega_p^{-1}(\pi/\delta_1)^{p/2},$$

$$\varkappa_2(t) \quad = \quad \|1 + t\|^{-p-\delta_2}, \quad \mathfrak{P}_2^* \quad = \quad p/\delta_2,$$

where $\delta_1, \delta_2 > 0$ are some constants. The result of Theorem 5.9 yields

**Corollary 5.10.** *Under conditions of Theorem 5.9, the bound (5.12) holds with*

$$
\begin{aligned}
\mathrm{pen}_1(\boldsymbol{v}) &= \varrho^{-1}\delta_1\,\epsilon^{-2}\|H^*(\boldsymbol{v}-\boldsymbol{v}_0)\|^2, \\
\log\mathfrak{Q}_2(\varrho,\epsilon) &= \frac{2\epsilon^2\varrho^2}{1-\varrho} + (1-\varrho)\mathbb{Q}_p + \log(1+\omega_p^{-1}|\pi/\delta_1|^{p/2}). \\
\mathrm{pen}_1(\boldsymbol{v}) &= -\varrho^{-1}(p+\delta_2)\log\big(\epsilon^{-1}\|H^*(\boldsymbol{v}-\boldsymbol{v}_0)\|+2\big), \\
\log\mathfrak{Q}_1(\varrho,\epsilon) &= \frac{2\epsilon^2\varrho^2}{1-\varrho} + (1-\varrho)\mathbb{Q}_p + \log(p/\delta_2),
\end{aligned}
$$

Sometimes it is useful to combine the functions $\varkappa_1(\cdot)$ and $\varkappa_2(\cdot)$ in the form

$$
\varkappa(t) = \varkappa_1(t)\mathbf{1}(t\geq r) + \varkappa_2(t)\mathbf{1}(t\leq r) \tag{5.13}
$$

for a properly selected $r$ which still ensures (5.11) with

$$
\mathfrak{P}^* \leq \omega_p^{-1}|\pi/\delta_1|^{p/2} + pr^{-\delta_2}/\delta_2.
$$

# References

Anderson, T. (1994). *The statistical analysis of time series.* Wiley Classics Library. Chichester: John Wiley &amp; Sons Ltd. xiv, 704 p. .

Basawa, I. V. and Brockwell, P. J. (1984). Asymptotic conditional inference for regular nonergodic models with an application to autoregressive processes. *Ann. Statist.*, 12(1):161–171.

Bollerslev, T. (1986). Generalized autoregressive conditional heteroscedasticity. *J. Econ.*, 31:307–327.

Brockwell, P. J. and Davis, R. A. (1991). *Time series: theory and methods. 2nd ed.* Springer Series in Statistics. Berlin etc.: Springer-Verlag.

Chan, N. H. and Wei, C. Z. (1988). Limiting distributions of least squares estimates of unstable autoregressive processes. *Ann. Statist.*, 16(1):367–401.

Chen, H. (1995). Asymptotically efficient estimation in semiparametric generalized linear models. *Ann. Statist.*, 23(4):1102–1129.

Chen, K., Hu, I., and Ying, Z. (1999). Strong consistency of maximum quasi-likelihood estimators in generalized linear models with fixed and adaptive designs. *Ann. Statist.*, 27(4):1155–1163.

Chen, Y. and Spokoiny, V. (2007). Modeling and estimation for nonstationary time series with applications to robust risk management. *Econometrics Journal.*

Cox, D. D. and Llatas, I. (1991). Maximum likelihood type estimation for nearly nonstationary autoregressive time series. *Ann. Statist.*, 19(3):1109–1128.

Dickey, D. A. and Fuller, W. A. (1981). Likelihood ratio statistics for autoregressive time series with a unit root. *Econometrika*, 49(4):1057–1072.

Fan, J. and Yao, Q. (2003). *Nonlinear time series. Nonparametric and parametric methods.* Springer Series in Statistics. New York, Springer .

Fokianos, K. and Kedem, B. (2003). Regression theory for categorical time series. *Stat. Sci.*, 18(3):357–376.

Fountis, N. G. and Dickey, D. A. (1989). Testing for a unit root nonstationarity in multivariate autoregressive time series. *Ann. Statist.*, 17(1):419–428.

Francq, C. and Zakoian, J.-M. (2007). Quasi-maximum likelihood estimation in GARCH processes when some coefficients are equal to zero. *Stochastic Process. Appl.*, 117(9):1265–1284.

Giacomini, E., Härdle, W., and Spokoiny, V. (2007). Inhomogeneous dependency modelling with time varying copulae. *Journal of Business and Economic Statistics.*

Golubev, Y. and Spokoiny, V. (2009). Preprint, WIAS, Berlin.

Green, P. J. and Silverman, B. (1994). *Nonparametric regression and generalized linear models: a roughness penalty approach.* London: Chapman & Hall. .

Johansen, S. (1995). *Likelihood-based inference in cointegrated vector autoregressive models.* Oxford: Oxford Univ. Press.

Johansen, S. (2002). A small sample correction for the test of cointegrating rank in the vector autoregressive model. *Econometrica*, 70(5):1929–1961.

Kedem, B. and Fokianos, K. (2002). *Regression models for time series analysis.* Hoboken, Wiley. .

Koenker, R. (2005). *Quantile regression.* Cambridge University Press.

Koenker, R. and Xiao, Z. (2006). Quantile autoregression. *J. Am. Stat. Assoc.*, 101(475):980–990.

Konev, V. and Pergamenshchikov, S. (1997). On guaranteed estimation of the mean of an autoregressive process. *Ann. Statist.*, 25(5):2127–2163.

Koul, H. L. and Saleh, A. K. M. E. (1993). $R$-estimation of the parameters of autoregressive [AR($p$)] models. *Ann. Statist.*, 21(1):534–551.

Lai, T. L. and Siegmund, D. (1983). Fixed accuracy estimation of an autoregressive parameter. *Ann. Statist.*, 11(2):478–485.

Lee, S.-W. and Hansen, B. E. (1994). Asymptotic theory for the GARCH(1, 1) quasi-maximum likelihood estimator. *Econometric Theory*, 10(1):29–52.

Massart, P. (2007). *Concentration inequalities and model selection. Ecole d'Eté de Probabilités de Saint-Flour XXXIII-2003.* Lecture Notes in Mathematics. Springer.

McCullagh, P. and Nelder, J. (1989). *Generalized linear models. 2nd ed.* Monographs on Statistics and Applied Probability. 37. London etc.: Chapman and Hall. xix, 511 p. .

Nelson, D. B. (1991). Conditional heteroskedasticity in asset returns: A new approach. *Econometrica*, 59(2):347–370.

Phillips, P. C. and Xu, K.-L. (2006). Inference in autoregression under heteroscedasticity. *J. Time Ser. Anal.*, 27(2):289–308.

Sentana, E. (1995). Quadratic ARCH models. *Rev. Econ. Stud.*, 62(4):639–661.

Shiryaev, A. N. and Spokoiny, V. G. (1997). On sequential estimation of an autoregressive parameter. *Stochastics Stochastics Rep.*, 60(3-4):219–240.

Spokoiny, V. (2007). Multiscale local change point detection with with applications to value-at-risk. *Ann. Statist.*

Sriram, T. N. (1987). Sequential estimation of the mean of a first-order stationary autoregressive process. *Ann. Statist.*, 15(3):1079–1090.

Sun, J., Loader, C., and McCormick, W. P. (2000). Confidence bands in generalized linear models. *Ann. Statist.*, 28(2):429–460.

Sun, Y. and Stengos, T. (2006). Semiparametric efficient adaptive estimation of asymmetric GARCH models. *J. Econometrics*, 133(1):373–386.

Van de Geer, S. A. (2000). *Applications of empirical process theory.* Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press.

van der Vaart, A. and Wellner, J. A. (1996). *Weak convergence and empirical processes. With applications to statistics.* Springer Series in Statistics. New York, Springer.

Čížek, P., Härdle, W., and Spokoiny, V. (2007). Adaptive pointwise estimation in time-inhomogeneous time-series models. *Econometrics Journal.*

Wang, C. W. H. (1986). A minimum distance estimator for first-order autoregressive processes. *Ann. Statist.*, 14(3):1180–1193.