



# Model selection for affine causal processes

Kare Kamila

## ► To cite this version:

Kare Kamila. Model selection for affine causal processes. Statistics [math.ST]. Université Paris 1 - Panthéon-Sorbonne, 2021. English. ⟨NNT : ⟩. ⟨tel-03500833⟩

**HAL Id: tel-03500833**

**<https://hal.science/tel-03500833v1>**

Submitted on 22 Dec 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

UNIVERSITÉ PARIS 1 PANTHÉON-SORBONNE

École Doctorale Sciences Mathématiques de Paris Centre (ED 386)

Laboratoire : Statistique, Analyse et Modélisation Multidisciplinaire, EA 4543

## Sélection de modèles pour les séries affines causales

KAMILA KARE

Thèse de Doctorat de Mathématiques Appliquées

Dirigée par

Jean-Marc BARDET et William KENGNE

Thèse présentée et soutenue le : 12 Juillet 2021

Après avis des rapporteurs : CHRISTIAN FRANCO (ENSAE)  
PASCAL MASSART (Université Paris-Saclay)

JURY DE SOUTENANCE :	JEAN-MARC BARDET	(Université Paris 1) Directeur de thèse
	FABIENNE COMTE	(Université de Paris) Examinatrice
	CHRISTIAN FRANCO	(ENSAE) Rapporteur
	WILLIAM KENGNE	(Cergy Paris Université) Codirecteur
	BÉATRICE LAURENT-BONNEAU	(INSA) Examinatrice
	PASCAL MASSART	(Université Paris-Saclay) Rapporteur
	OLIVIER WINTENBERGER	(Sorbonne Université) Examineur



## Résumé

L'analyse des séries temporelles est un sujet de recherche très actif en Statistique et en Data Science. L'abondance de ce type de données a créé d'énormes besoins de méthodologies efficaces et précises. C'est ainsi que plusieurs familles de modèles ont vu le jour. Étant donnée cette multitude de modèles, comment en choisir un pour modéliser une série temporelle ? L'objet de cette thèse est de proposer et d'étudier des critères de sélection de modèles pour une grande famille de modèles contenant les séries temporelles autorégressives telles les ARMA ainsi que les séries temporelles conditionnellement hétéroscédastiques telles les GARCH.

Nous commençons par une brève présentation des séries temporelles, en particulier de la famille des modèles affines causaux tout en rappelant quelques précédents résultats utiles pour cette thèse. Nous décrivons ensuite quelques critères classiques de sélection de modèles obtenus pour les séries temporelles et terminons par un succinct résumé de nos principales contributions.

Dans la suite, nous allons présenter les quatre contributions originales de cette thèse. Le chapitre 2 donne des conditions suffisantes sur la pénalité en fonction de la régularité de la dépendance du processus par rapport à son passé afin d'obtenir un critère consistant en probabilité. Nous proposons également un test d'adéquation du modèle sélectionné basé sur l'autocorrélation du carré des résidus du modèle. Les simulations numériques ont montré des résultats satisfaisants.

Au chapitre 3, nous proposons une généralisation du critère de Hannan et Quinn à la classe des séries affines causales. Cette généralisation induit une certaine constante connue pour les modèles classiques (type ARMA, GARCH ou APARCH) et pour les modèles complexes tels les ARMA-GARCH, la constante est inconnue mais peut être estimée de manière adaptative via l'heuristique de pente. Là également, quelques études de simulation ont attesté de la qualité des critères obtenus.

Dans la troisième contribution, nous construisons des critères asymptotiquement efficaces. Nous proposons une généralisation du critère AIC d'Akaike se basant sur la pénalité dite idéale. Le comportement asymptotique de cette pénalité idéale nous a suggéré un terme de pénalité qui vaut exactement  $2D_m$  comme dans l'AIC pour des modèles assez simples, et pour des modèles complexes, nous avons donné une formule moins explicite. A la suite de Schwartz, nous dérivons également le critère BIC qui s'appuie sur la maximisation de la probabilité a posteriori de choisir le vrai modèle.

Au chapitre 5, nous nous sommes restreints à l'étude non asymptotique d'un processus particulier de la classe des modèles affines causaux. Un estimateur des moindres carrés pénalisé est construit à partir d'un critère de sélection adaptatif et la sélection est opérée parmi une collection de modèles linéaires. Nous avons montré que l'estimateur final est presque aussi performant que le meilleur sur la collection considérée, *i.e.* qu'il réalise, à une constante près, le compromis biais-variance. La pénalité obtenue généralise celle de Mallows et dépend d'une constante que l'on pourrait estimer avec des algorithmes de calibration adaptative.

Enfin, nous donnons quelques pistes de recherche dans la Conclusion générale du travail.

**Mots-Clefs :** Modèles affines causaux, sélection de modèles, estimation par quasi-vraisemblance, critères consistants, critères efficients, inégalité oracle, heuristique de pente.

## Abstract

Time series analysis is a very active research topic in Statistics and Data Science. The abundance of this type of data has created a huge need for efficient and accurate methodologies. Thus, several families of models have emerged. Given this multitude of models, how do you choose one to model a time series? The purpose of this thesis is to propose and study model selection criteria for a large family of models containing autoregressive time series such as ARMA and conditionally heteroscedastic time series such as GARCH.

We start with a brief presentation of time series, in particular of the family of causal affine models, while recalling some previous results useful for this thesis. We then describe some classical model selection criteria obtained for time series and end with a brief summary of our main contributions.

The rest of our work presents our four contributions. The first chapter gives sufficient conditions on the penalty depending on the regularity of the dependence of the process on its past in order to obtain a consistent criterion. We also propose a goodness-of-fit test of the selected model based on the autocorrelation of the square of the model residuals. The numerical simulations have shown satisfactory results.

In Chapter 3, we propose a generalization of the Hannan and Quinn criterion to the class of causal affine series. This generalization induces a certain constant known for classical models (ARMA, GARCH or APARCH type) and can be data-driven estimated for complex models like ARMA-GARCH. Here again, some simulation studies have attested to the quality of the criteria obtained.

In the third contribution, we construct asymptotically efficient criteria. We propose a generalization of Akaike's AIC criterion based on the so-called ideal penalty. The asymptotic behavior of this ideal penalty suggested a penalty term which is exactly  $2D_m$  as in the AIC for simple models, and for complex models, we give a less explicit formula. Following Schwartz, we also derive the BIC criterion based on the maximization of the a posteriori probability of choosing the true model.

In Chapter 5, we restricted ourselves to the non-asymptotic study of a particular process of the class of causal affines. A penalized least-squares estimator is built on a data driven selected model among a collection of linear models. We showed that the final estimator performs almost as well as the best over the considered collection, *i.e.* it achieves, up to a constant, the bias-variance tradeoff. The penalty obtained generalizes Mallows' penalty and depends on a constant that is estimated with data-driven calibration algorithms.

Finally, we give some research directions in the General Conclusion of the work.

**Keywords:** Causal affine models, model selection, quasi-likelihood estimation, data-driven, consistent criteria, efficient criteria, oracle inequality, slope heuristic.



## Remerciements

La réalisation de ce travail de longue haleine a été possible grâce au concours de plusieurs personnes à qui je voudrais témoigner toute ma gratitude.

Je voudrais tout d'abord adresser toute ma reconnaissance à mon directeur de thèse, Jean-Marc BARDET, pour sa patience, sa disponibilité, sa rigueur scientifique et surtout ses précieux conseils, qui ont contribué à alimenter ma réflexion. Tu m'as fait confiance dès nos premiers échanges et tu m'as soutenu dans toutes les étapes du concours COFUND Math In Paris qui m'a permis de décrocher une bonne source de financement. JM, 1000 Mercis!

Je tiens à exprimer toute ma gratitude à William KENGNE, co-directeur de cette thèse. Ton dynamisme, ta disponibilité, ta rigueur, tes conseils m'ont été précieux dans l'avancée de mes travaux.

Je suis extrêmement reconnaissant envers Christian FRANCO et Pascal MASSART pour l'honneur qu'ils m'ont fait en acceptant de rapporter cette thèse. Merci pour cette lecture minutieuse et vos rapports.

Mes remerciements vont à l'endroit de Fabienne COMTE, Béatrice LAURENT-BONNEAU et Olivier WINTENBERGER pour leur participation à mon jury de thèse.

Je remercie la Fondation des Sciences Mathématiques de Paris pour avoir mis sur pied le programme doctoral COFUND Math In Paris. Je dis aussi merci à l'UE pour le financement accordé au COFUND. Je profite pour remercier Ariela pour sa disponibilité et l'accompagnement constant durant les trois années de financement.

J'adresse mes sincères remerciements à Paul DOUKHAN pour les réponses à mes questions.

J'ai une pensée pour mes enseignants du Master MASS POP d'Aix-Marseille, notamment à Marie-Christine ROUBAUD, Nicolas PECH, Thomas WILLER et Pierre PUDLO, pour leur accompagnement durant le concours COFUND. Marie-Christine, c'est aussi grâce à tes paroles motivantes et captivantes qui raisonnent encore aujourd'hui dans mes oreilles que j'ai été lauréat COFUND 2018.

J'exprime également mes remerciements à la chaleureuse équipe du SAMM et plus particulièrement à Nicolas, Branda, Xavier, Dafnis, Florian, Alexandre, Diem, Clara, Alice, Marie, Antoine, Béchir, Alain, Julien, Joseph, la liste étant loin d'être exhaustive.

Une mention spéciale à tous les amis qui n'ont cessé de m'apporter leur soutien: Merci à toi Dolores, Tara (Moussa), NDAM, Marcel, Alida, Ismaël, Ousman.

Je remercie ma famille, ma source de motivation: Maman, chacune de nos discussions me donnait plein d'énergie pour surmonter les moments de doute et travailler sans relâche. Papa, Mami, Digora, Fidel, Bayer merci pour votre soutien.





## Notations

$\mathbb{R}^\infty$	Set of sequences of real numbers with a finite number of non-zero terms
$ m $ or $D_m$	is the dimension of the model $m$
p.s.	means almost surely (in Chapter 1)
$X_n \xrightarrow[n \rightarrow +\infty]{a.s.} X$	$\mathbb{P}(X_n \xrightarrow[n \rightarrow \infty]{} X) = 1$
$X_n = o_{a.s.}(g_n)$	$\frac{X_n}{g_n} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$
$X_n \xrightarrow[n \rightarrow \infty]{\mathcal{P}} X$	$\forall \epsilon > 0, \mathbb{P}( X_n - X  > \epsilon) \xrightarrow[n \rightarrow \infty]{} 0$
$X_n = o_P(g_n)$	$\frac{X_n}{g_n} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} 0$
$X_n = O_P(g_n)$	$\forall \epsilon > 0, \exists M_\epsilon > 0$ such that $\mathbb{P}( X_n  \geq M_\epsilon g_n) \leq \epsilon$
$X_n \xrightarrow[n \rightarrow \infty]{\mathcal{D}} X$	$F_{X_n}(x) \xrightarrow[n \rightarrow \infty]{} F_X(x)$ for every number $x \in \mathbb{R}$ at which $F$ is continuous where $F_{X_n}$ and $F_X$ are the cumulative distribution functions of random variables $X_n$ and $X$ , respectively
$X \sim \text{SG}(\sigma^2)$	$\mathbb{E}[e^{\lambda X}] \leq e^{\lambda^2 \sigma^2 / 2}$ for every $\lambda > 0$ provided that $\mathbb{E}[X] = 0$
$X \sim \text{SE}(\sigma^2, \alpha)$	$\mathbb{E}[e^{\lambda X}] \leq e^{\lambda^2 \sigma^2 / 2}, \forall \lambda :  \lambda  < \frac{1}{\alpha}$ provided that $\mathbb{E}[X] = 0$
$\ \cdot\ $	denotes the usual Euclidean norm on $\mathbb{R}^\nu$ , with $\nu \geq 1$
$\ \cdot\ $ $\ A\  = \sup_{v \neq 0} \frac{\ Av\ }{\ v\ }$	for a matrix $A$ , denote the subordinate matrix norm such that
$\ \cdot\ _r$ $(\mathbb{E}[\ X\ ^r])^{1/r} \in [0, \infty]$	if $X$ is a $\mathbb{R}^\nu$ -random variable and $r \geq 1$ , we set $\ X\ _r = (\mathbb{E}[\ X\ ^r])^{1/r} \in [0, \infty]$
$\ \cdot\ _\Theta$	for $\theta \in \Theta \subset \mathbb{R}^d$ , if $\Psi_\theta : \mathbb{R}^\infty \rightarrow E$ where $E = \mathbb{R}^\nu$ or $E$ is a set of square matrix, denote $\ \Psi_\theta(\cdot)\ _\Theta = \sup_{\theta \in \Theta} \{\ \Psi_\theta(\cdot)\ \}$
$\partial_\theta$ .	for $\theta \in \Theta \subset \mathbb{R}^d$ , if $\Psi_\theta : \mathbb{R}^\infty \rightarrow \mathbb{R}$ is a $\mathcal{C}^2(\Theta \times \mathbb{R}^\infty)$ function, we will denote $\partial_\theta \Psi_\theta(\cdot) = \left( \frac{\partial}{\partial \theta_i} \Psi_\theta(\cdot) \right)_{1 \leq i \leq d} = (\partial_{\theta_i} \Psi_\theta(\cdot))_{1 \leq i \leq d}$
$\partial_{\theta^2}^2$ .	$\partial_{\theta^2}^2 \Psi_\theta(\cdot) = \left( \frac{\partial^2}{\partial \theta_i \partial \theta_j} \Psi_\theta(\cdot) \right)_{1 \leq i, j \leq d}$

# Contents

1	Révue de Littérature et Synthèse des Travaux	1
1.1	Modèles Affines Causaux	2
1.2	Sélection de modèles	7
1.3	Synthèse des travaux	15
2	Consistent model selection criteria and goodness-of-fit test for common time series models	29
2.1	Introduction	30
2.2	General Framework	33
2.3	Asymptotic results	36
2.4	Examples	38
2.5	Portmanteau test	41
2.6	Numerical Results	43
2.7	Proofs	49
3	General Hannan and Quinn Criterion for Common Time Series	65
3.1	INTRODUCTION	66
3.2	MODEL SELECTION CONSISTENCY	68
3.3	NUMERICAL EXPERIMENTS	74
3.4	Proofs	76
4	Efficient and consistent data-driven model selection for time series	85
4.1	Introduction	86
4.2	Model selection framework	87
4.3	Asymptotic behavior of the QMLE	91
4.4	Efficient model selection result	93
4.5	From a Bayesian model selection to a data-driven consistent model selection	95
4.6	Computations of the ideal penalties	97
4.7	Numerical Studies	100
4.8	Proofs	101
5	Data driven model selection for same-realization predictions in autoregressive processes	113
5.1	INTRODUCTION	114
5.2	MODEL SELECTION APPROACH AND PRELIMINARY RESULTS	115
5.3	Bias-Variance Result	121
5.4	PROOFS	121
5.5	THEORETICAL TOOLS	133
6	Conclusion Générale et Perspectives	135

# 1

## Révue de Littérature et Synthèse des Travaux

---

<b>1.1</b>	<b>Modèles Affines Causaux</b>	<b>2</b>
1.1.1	Quelques préliminaires sur les séries temporelles	2
1.1.2	Modèles affines causaux	3
1.1.2.a	Modèle général et propriétés	3
1.1.2.b	Exemples de processus affines causaux	4
1.1.2.c	QMLE et quelques résultats	5
<b>1.2</b>	<b>Sélection de modèles</b>	<b>7</b>
1.2.1	Pourquoi sélectionner un modèle?	7
1.2.2	Objectifs de la Sélection de modèles	8
1.2.2.a	Identification	8
1.2.2.b	Prédiction	8
1.2.3	Stratégies de Sélection de modèles	9
1.2.3.a	Tests Statistiques	9
1.2.3.b	Validation croisée	10
1.2.3.c	Pénalisation	10
1.2.4	Dérivation de quelques critères dans le cadre des séries temporelles	11
1.2.4.a	FPE <a href="#">Akaike (1969)</a>	11
1.2.4.b	AIC <a href="#">Akaike (1973)</a>	12
1.2.4.c	BIC <a href="#">Schwarz (1978)</a>	14
1.2.4.d	HQ <a href="#">Hannan and Quinn (1979)</a>	15
1.2.5	Calibration de la constante multiplicative	15
<b>1.3</b>	<b>Synthèse des travaux</b>	<b>15</b>
1.3.1	Critères de Sélection de modèles consistants et test d'ajustement pour les modèles affines causaux	15
1.3.1.a	Hypothèses	17
1.3.1.b	Résultats théoriques	18
1.3.1.c	Test Portmanteau	18
1.3.2	Généralisation du critère de Hannan et Quinn pour les séries affines causales	19
1.3.2.a	Résultats théoriques	20
1.3.3	Sélection de modèles efficiente et consistante pour les séries temporelles	21
1.3.3.a	Hypothèses	22

1.3.3.b	Résultats théoriques	23
1.3.3.c	Dérivation du BIC pour les modèles affines causaux	24
1.3.3.d	Quelques exemples de calcul du terme trace dans la pénalité idéale	25
1.3.4	Sélection de modèles data driven pour la prédiction d'un processus linéaire autorégressif	25
1.3.4.a	Hypothèses	26
1.3.4.b	Résultats théoriques	27

---

Ce chapitre introductif a pour but de présenter tour à tour le modèle statistique étudié, quelques critères de sélection de modèles et nos contributions sur la question de la sélection de modèles pour les séries affines causales.

## 1.1 Modèles Affines Causaux

### 1.1.1 Quelques préliminaires sur les séries temporelles

Une série temporelle est une suite d'observations, observées sur une période de temps  $T$ . Lorsque la période  $T$  est un sous ensemble de  $\mathbb{R}$ , on parlera de série à temps continu et de série à temps discret quand  $T$  est plutôt une partie de  $\mathbb{Z}$ . Dans toute cette thèse, l'on ne considérera que les séries à temps discret.

Les observations temporelles apparaissent dans de nombreux domaines clés du monde réel, allant de la finance à la biologie, en passant par la météorologie, le traitement du langage naturel, la détection des anomalies (dans les systèmes de contrôle), l'audio et le traitement vidéo, pour n'en citer que quelques-uns. Compte tenu de leur abondance, l'analyse des séries chronologiques est devenue un domaine clé de la Statistique et aussi du Machine Learning.

Dans l'analyse des séries temporelles, l'on s'intéresse à la compréhension et à la modélisation des relations entre les observations d'une même variable: la variable d'intérêt est identique aux variables explicatives à la seule différence qu'elles sont observées à des instants distincts. Donnons une définition plus formelle et générale des séries chronologiques.

**Définition 1.1.** Une série temporelle  $(X_1, X_2, \dots, X_n)$  est une réalisation du processus stochastique  $(X_t)_{t \in \mathbb{Z}}$  défini sur un espace probabilisé  $(\Omega, \mathcal{F}, P)$ .

À partir d'une seule réalisation  $(X_1, X_2, \dots, X_n)$ , l'on aimerait faire de l'inférence sur les paramètres associés au processus  $(X_t)_{t \in \mathbb{Z}}$ , notamment l'espérance, la variance etc. En Statistique classique, l'on sait que la moyenne empirique des réalisations **indépendantes** d'une même variable aléatoire est un bon estimateur de l'espérance. Et l'on sait aussi que l'hypothèse d'indépendance est fondamentale car si les réalisations sont fortement dépendantes alors il est probable que la moyenne empirique ne recouvre pas l'espace de probabilité tout entier et qu'elle ne converge pas vers l'espérance.

Ainsi, peut-on obtenir de bons estimateurs des paramètres lorsqu'on a affaire à une série temporelle? La réponse est oui mais au prix de quelques conditions sur le processus. Par exemple, l'on peut supposer que le processus est de moyenne constante et qu'il est de mémoire courte de sorte que l'effet de  $X_t$  sur  $X_{t+k}$  ne soit pas trop important pour  $k$  assez grand.

Introduisons plus explicitement la notion de stationnarité qui est un concept assez intuitif et suppose que la structure du processus sous-jacent n'évolue pas avec le temps.

**Définition 1.2.** (*forte stationnarité*) Un processus aléatoire  $X = (X_t)_{t \in \mathbb{Z}}$  est dit (*fortement*) stationnaire s'il est invariant en distribution par toute translation du temps, i.e.  $\forall k \in \mathbb{N}^*, \forall (t_1, \dots, t_k) \in \mathbb{Z}^k, \forall c \in \mathbb{Z}$ , les vecteurs  $(X_{t_1}, \dots, X_{t_k})$  et  $(X_{t_1+c}, \dots, X_{t_k+c})$  ont la même distribution.

Ce type de stationnarité comme l'indique son nom, est plutôt une condition forte et est souvent difficile à vérifier en pratique. On lui préfère la stationnarité faible qui ne nécessite que l'invariance des moments jusqu'à l'ordre deux.

**Définition 1.3.** (*faible stationnarité*) Un processus aléatoire  $X = (X_t)_{t \in \mathbb{Z}}$  est dit (*faiblement*) stationnaire si son espérance est constante et si pour tout  $t, k \in \mathbb{Z}$ , la covariance entre  $X_t$  et  $X_{t+k}$  ne dépend que de  $k$ .

Un autre concept qui accompagne très souvent la notion de stationnarité est l'ergodicité. Il permet de généraliser la loi des grands nombres aux variables aléatoires dépendantes.

**Définition 1.4.** (*Ergodicité pour processus stationnaire*) Le processus fortement stationnaire  $X$  est dit ergodique si et seulement si pour tout borélien  $B$  et pour tout entier  $k$ ,

$$n^{-1} \sum_{t=1}^n \mathbb{I}_B(X_t, X_{t+1}, \dots, X_{t+k}) \xrightarrow[n \rightarrow +\infty]{a.s.} \mathbb{P}\{(X_1, \dots, X_{1+k}) \in B\}.$$

### 1.1.2 Modèles affines causaux

Pour décrire les relations entre les observations, l'on recourt à la modélisation. Dans le cadre des séries temporelles, trois grandes familles de modèles se distinguent:

- Les séries linéaires: importantes pour modéliser la structure de covariance de la série. Cette famille comprend les modèles autorégressifs (AR), les modèles de moyenne mobile (MA) et leur combinaison (ARMA);
- Les modèles non linéaires pour prendre en compte les irrptions soudaines et irrationnelles que l'on observe dans les données, notamment financières. Il s'agit des modèles GARCH et associés;
- Les combinaisons des deux premières familles i.e. les ARMA-GARCH.

La classe de modèles affines causaux se propose d'unifier l'écriture de ces trois familles afin de les traiter simultanément dans un cadre identique.

#### 1.1.2.a Modèle général et propriétés

**Définition 1.5.** Soit  $(\xi_t)_{t \in \mathbb{Z}}$  une suite de variables aléatoires iid telle que  $\mathbb{E}[|\xi_0|^r] < \infty$ . Le processus  $(X_t)_{t \in \mathbb{Z}}$  est dit affine causal s'il existe deux fonctions mesurables  $M, f : \mathbb{R}^\infty \rightarrow \mathbb{R}$  telles que

$$X_t = M((X_{t-i})_{i \in \mathbb{N}^*}) \xi_t + f((X_{t-i})_{i \in \mathbb{N}^*}) \text{ pour tout } t \in \mathbb{Z}. \quad (1.1.1)$$

L'on notera  $X = (X_t)_{t \in \mathbb{Z}} \in \mathcal{AC}(M, f)$ . Dans toute cette thèse, l'on se restreindra au cas où les fonctions  $M$  et  $f$  sont connus mais dépendent d'un paramètre inconnu  $\theta^*$  appartenant à une certaine région admissible de paramètres  $\Theta$  de  $\mathbb{R}^d$ . Nos travaux s'inscrivent dans un cadre sémi-paramétrique puisque l'on ne fera aucune hypothèse sur la distribution du bruit  $(\xi_t)_{t \in \mathbb{Z}}$ .

Une fois le modèle posé, il est fondamental avant l'étape d'inférence, de montrer qu'il est bien défini i.e. qu'il admet de solutions ayant un intérêt pratique (finitude, stationarité, moments, etc). Ainsi, sous quelles conditions (1.1.1) admet-il une solution stationnaire? L'étude de l'existence de solution du modèle (1.1.1) et de certaines propriétés a été effectuée par [Doukhan and Wintenberger \(2008\)](#) en termes de coefficients de contraction des fonctions  $M_\theta$  et  $f_\theta$  dits coefficients de Lipschitz.

Nous énonçons la principale condition de [Doukhan and Wintenberger \(2008\)](#) avec  $\Psi_\theta = f_\theta, M_\theta, H_\theta = M_\theta^2$  où  $\theta \in \Theta$  ( $\Theta$  supposé compact) et  $i = 0, 1, 2$ .

**A<sub>i</sub>( $\Psi_\theta, \Theta$ ):** Supposons  $\|\Psi_\theta(0)\|_\Theta < \infty$  et qu'il existe une suite de nombres positifs  $(\alpha_k^i(\Psi_\theta, \Theta))_{k \geq 1}$  telle que  $\sum_{k=1}^\infty \alpha_k^i(\Psi_\theta, \Theta) < \infty$  vérifiant:

$$\left\| \frac{\partial^i \Psi_\theta(x)}{\partial \theta^i} - \frac{\partial^i \Psi_\theta(y)}{\partial \theta^i} \right\|_\Theta \leq \sum_{k=1}^\infty \alpha_k^i(\Psi_\theta, \Theta) |x_k - y_k| \quad \text{pour tout } x, y \in \mathbb{R}^\infty.$$

Pour les modèles de type GARCH admettant une espérance conditionnelle nulle, il est plus optimale de considérer

**A<sub>i</sub>( $H_\theta, \Theta$ ):** Supposons  $\|H_\theta(0)\|_\Theta < \infty$  et qu'il existe une suite de nombres positifs  $(\alpha_k^i(H_\theta, \Theta))_{k \geq 1}$  telle que  $\sum_{k=1}^\infty \alpha_k^i(H_\theta, \Theta) < \infty$  vérifiant:

$$\left\| \frac{\partial^i H_\theta(x)}{\partial \theta^i} - \frac{\partial^i H_\theta(y)}{\partial \theta^i} \right\|_\Theta \leq \sum_{k=1}^\infty \alpha_k^i(H_\theta, \Theta) |x_k^2 - y_k^2| \quad \text{pour tout } x, y \in \mathbb{R}^\infty.$$

L'idée de [Doukhan and Wintenberger \(2008\)](#) consiste à restreindre l'espace des paramètres de sorte qu'une condition sur le moment du bruit entraîne l'existence d'une solution stationnaire ergodique admettant des moments. Ainsi, si  $\xi_0$  admet des moments d'ordre  $r \geq 1$ , considérons

$$\Theta(r) = \left\{ \theta \in \mathbb{R}^d, A_0(f_\theta, \{\theta\}) \text{ et } A_0(M_\theta, \{\theta\}) \text{ vraies avec } \sum_{k=1}^\infty \alpha_k^0(f_\theta, \{\theta\}) + \|\xi_0\|_r \sum_{k=1}^\infty \alpha_k^0(M_\theta, \{\theta\}) < 1 \right\} \quad (1.1.2)$$

ou bien pour les modèles de type GARCH

$$\Theta(r) = \left\{ \theta \in \mathbb{R}^d, A_0(H_\theta, \{\theta\}) \text{ vraie avec } \|\xi_0\|_r \sum_{k=1}^\infty \alpha_k^0(H_\theta, \{\theta\}) < 1 \right\}. \quad (1.1.3)$$

Nous pouvons rappeler le résultat de [Doukhan and Wintenberger \(2008\)](#).

**Proposition 1.1.** *Pour tout  $\theta \in \Theta(r)$  avec  $r \geq 1$ , (1.1.1) admet une unique solution  $(X_t)_{t \in \mathbb{Z}}$  stationnaire, ergodique, faiblement dépendante avec  $\|X_0\|_r < \infty$ .*

### 1.1.2.b Exemples de processus affines causaux

#### 1. Modèle AR( $\infty$ )

Soit  $(\psi_k(\theta))_{k \in \mathbb{N}}$  une suite de réels dépendante de  $\theta \in \Theta \subset \mathbb{R}^d$ . Considérons l'AR( $\infty$ ) définie ainsi qu'il suit:

$$X_t = \sum_{k \geq 1} \psi_k(\theta) X_{t-k} + \sigma \xi_t \quad \text{pour tout } t \in \mathbb{Z}, \quad (1.1.4)$$

où  $(\xi_t)_{t \in \mathbb{Z}}$  est une suite de variable aléatoire iid admettant des moments d'ordre  $r \geq 1$ , et  $\sigma > 0$ . Ce processus correspond à (1.1.1) avec  $f_\theta((x_i)_{i \geq 1}) = \sum_{k \geq 1} \psi_k(\theta) x_k$  et  $M_\theta \equiv \sigma$  pour tout  $\theta \in \Theta$ . Les coefficients de Lipschitz associés à  $f_\theta$  sont  $\alpha_k^0(f_\theta) = \|\psi_k(\theta)\|_\Theta$ . La Proposition 1.1 est vraie pour tout paramètre  $\theta$  vérifiant  $\sum_{k=1}^\infty \|\psi_k(\theta)\|_{\{\theta\}} < 1$ .

Ainsi soit l'ARMA( $p, q$ ) suivant

$$X_t - \sum_{i=1}^p \phi_i X_{t-i} = \xi_t + \sum_{i=1}^p \theta_i \xi_{t-i}, \quad (1.1.5)$$

tel que toutes les racines du polynôme  $\theta(z) = \sum_{i=1}^p \theta_i z^i$  soient de module strictement plus grand que 1. Alors le processus  $(X_t)_{t \in \mathbb{Z}}$  est inversible et (1.1.5) admet une représentation affine causale (1.1.4).

## 2. Modèle ARCH( $\infty$ )

Considérons  $(\psi_k(\theta))_{k \in \mathbb{N}}$  une suite de réels dépendante de  $\theta \in \Theta \subset \mathbb{R}^d$ . Soit l'ARCH( $\infty$ ) suivant

$$X_t = \left( \psi_0(\theta) + \sum_{k=1}^\infty \psi_k(\theta) X_{t-k}^2 \right)^{1/2} \xi_t \quad \text{pour tout } t \in \mathbb{Z}, \quad (1.1.6)$$

où  $(\xi_t)_t$  admet des moments d'ordre  $r \geq 1$ . Nous avons là un exemple de (1.1.1) avec  $f_\theta = 0$  et  $M_\theta((x_i)_{i \geq 1}) = \left( \psi_0(\theta) + \sum_{k=1}^\infty \psi_k(\theta) x_{t-k}^2 \right)^{1/2}$ . Ainsi les coefficients de Lipschitz (relativement à la série  $\{x_t^2\}$ ) sont  $\alpha_k^0(H_\theta) = \|\psi_k(\theta)\|_\Theta$  et tout paramètre  $\theta$  vérifiant  $\|\xi_0\|_r \sum_{k=1}^\infty \|\psi_k(\theta)\|_{\{\theta\}} < 1$  assure que le processus  $(X_t)$  est stationnaire, ergodique d'après la Proposition 1.1.

## 3. Modèle TARCH( $\infty$ )

Ce modèle est défini ainsi qu'il suit

$$X_t = \left( \psi_0(\theta) + \sum_{k=1}^\infty \psi_k^+(\theta) \max(X_{t-k}, 0) - \psi_k^-(\theta) \min(X_{t-k}, 0) \right) \xi_t \quad \text{pour tout } t \in \mathbb{Z}, \quad (1.1.7)$$

avec  $(\xi_t)_t$  admet des moments d'ordre  $r \geq 1$ . On a encore un cas de (1.1.1) avec  $f_\theta = 0$  et  $M_\theta((x_i)_{i \geq 1}) = \psi_0(\theta) + \sum_{k=1}^\infty \psi_k^+(\theta) \max(X_{t-k}, 0) - \psi_k^-(\theta) \min(X_{t-k}, 0)$ . On obtient que  $\alpha_k^0(M_\theta) = \max(\|\psi_k^+(\theta)\|_\Theta, \|\psi_k^-(\theta)\|_\Theta)$  de sorte que tout  $\theta$  vérifiant  $\|\xi_0\|_r \sum_{k=1}^\infty \alpha_k^0(M_\theta) < 1$  assure toutes les propriétés de la Proposition 1.1.

### 1.1.2.c QMLE et quelques résultats

À présent que nous avons des garanties théoriques sur l'existence de la solution stationnaire du processus (1.1.1), une question naturelle qui vient à l'esprit serait celle de l'estimation du paramètre  $\theta^*$ .

Supposons que nous avons observé  $(X_1, X_2, \dots, X_n)$  où les  $X_t$  proviennent de (1.1.1) avec  $\theta^* \in \Theta(r)$  inconnu pour un  $r \geq 1$ . Nous allons construire un estimateur de  $\theta^*$  en utilisant un contraste basé sur la vraisemblance.



Dans l'étude des séries chronologiques, très souvent la distribution des innovations est inconnue, c'est pourquoi nous aurons recours à la méthode de quasi-vraisemblance. Elle consiste à supposer dans un premier temps que le bruit est gaussien et à trouver l'expression de la densité jointe des observations. Cette expression qui donne un contraste sera utilisée même si le bruit n'est pas gaussien. Par ailleurs, il est aussi important de rappeler qu'ici, la densité jointe des observations est en pratique impossible à maximiser. L'on se contente d'une version conditionnelle.

Plus formellement, supposons que  $(\xi_t)_{t \in \mathbb{Z}}$  soit un bruit blanc gaussien (standard). À partir de (1.1.1), l'on déduit que la log densité de  $X_t$  sachant  $\sigma(X_i, i < t)$  est

$$-\frac{1}{2} \left[ \frac{(X_t - f_{\theta^*}^t)^2}{H_{\theta^*}^t} + \log(H_{\theta^*}^t) \right].$$

Ainsi, la log densité conditionnelle de  $(X_1, \dots, X_n)$  sachant  $\sigma(X_t, t \leq 0)$  est

$$-\frac{1}{2} \sum_{t=1}^n \left[ \frac{(X_t - f_{\theta^*}^t)^2}{H_{\theta^*}^t} + \log(H_{\theta^*}^t) \right].$$

Desormais, l'hypothèse de gaussianité du bruit est abandonnée et la log-densité conditionnelle permet de définir pour tout  $\theta \in \Theta$

$$L_n(\theta) := -\frac{1}{2} \sum_{t=1}^n q_t(\theta), \text{ avec } q_t(\theta) := \frac{(X_t - f_{\theta}^t)^2}{H_{\theta}^t} + \log(H_{\theta}^t) \quad (1.1.8)$$

où  $f_{\theta}^t := f_{\theta}(X_{t-1}, X_{t-2}, \dots)$ ,  $M_{\theta}^t := M_{\theta}(X_{t-1}, X_{t-2}, \dots)$  et  $H_{\theta}^t = (M_{\theta}^t)^2$ . L'on remarque tout de suite que la fonction de vraisemblance  $L_n$  n'est pas évaluable puisqu'elle dépend du passé  $(X_{-j})_{j \in \mathbb{N}}$  qui est inconnu. C'est pourquoi, nous considérerons une approximation observable de  $L_n$  notée  $\hat{L}_n$  et définie ainsi qu'il suit

$$\hat{L}_n(\theta) := -\frac{1}{2} \sum_{t=1}^n \hat{q}_t(\theta), \text{ avec } \hat{q}_t(\theta) := \frac{(X_t - \hat{f}_{\theta}^t)^2}{\hat{H}_{\theta}^t} + \log(\hat{H}_{\theta}^t) \quad (1.1.9)$$

où  $\hat{f}_{\theta}^t := f_{\theta}(X_{t-1}, X_{t-2}, \dots, X_1, 0, \dots, 0)$ ,  $\hat{M}_{\theta}^t := M_{\theta}(X_{t-1}, X_{t-2}, \dots, X_1, 0, \dots, 0)$  et  $\hat{H}_{\theta}^t = (\hat{M}_{\theta}^t)^2$ . L'estimateur de quasi-vraisemblance (QMLE) de  $\theta^*$  est donc:

$$\hat{\theta}_n = \operatorname{argmax}_{\theta \in \Theta} \hat{L}_n(\theta). \quad (1.1.10)$$

Notons que les propriétés (consistance et normalité asymptotique) de cet estimateur ont été étudiées par [Bardet and Wintenberger \(2009\)](#) sous certaines hypothèses de régularité qui seront utilisées tout au long de cette thèse.

La première hypothèse porte sur l'identifiabilité du "vrai" paramètre  $\theta^*$ . Il induit que deux paramètres qui conduisent à une même valeur de la vraisemblance sont identiques.

**A1:** Pour tout  $\theta, \theta' \in \Theta$ ,  $(f_{\theta}^0 = f_{\theta'}^0 \text{ et } M_{\theta}^0 = M_{\theta'}^0) \text{ p.s.} \implies \theta = \theta'$ .

Notons que cette hypothèse est vérifiée pour tous les modèles affines causaux classiques lorsque le bruit est non dégénéré.

Aussi, dans les définitions des fonctions de vraisemblance et quasi-vraisemblance conditionnelles, il apparaît un dénominateur qui ne devrait pas s'annuler. Ainsi, nous supposons

**A2:**  $\exists \underline{h} > 0$  tel que  $\inf_{\theta \in \Theta} (H_\theta(x)) \geq \underline{h}$  pour tout  $x \in \mathbb{R}^\infty$ .

La condition suivante permet d'assurer l'inversibilité de la matrice hessienne de la fonction  $L_n$  et est très importante pour prouver la normalité asymptotique de l'estimateur QMLE.

**A3:** Une des familles  $(\partial f_\theta^t / \partial \theta^{(i)})_{1 \leq i \leq d}$  ou  $(\partial H_\theta^t / \partial \theta^{(i)})_{1 \leq i \leq d}$  est linéairement indépendante presque sûrement.

**Théorème 1.1.** (*Bardet and Wintenberger (2009)*) Sous les conditions **A1-A3** et si

$$\alpha_k^0(f_\theta, \Theta) + \alpha_k^0(M_\theta, \Theta) = O(k^{-\gamma}) \quad \text{avec } \gamma > 3/2,$$

Alors le QMLE est fortement consistant i.e.

$$\hat{\theta}_n \xrightarrow[n \rightarrow +\infty]{a.s.} \theta^*. \quad (1.1.11)$$

Si par ailleurs,

$$\alpha_k^1(f_\theta, \Theta) + \alpha_k^1(M_\theta, \Theta) = O(k^{-\gamma}) \quad \text{avec } \gamma > 3/2,$$

Alors  $\hat{\theta}_n$  est asymptotiquement normal i.e.

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}\left(0, (F(\theta^*))^{-1} G(\theta^*) (F(\theta^*))^{-1}\right) \quad (1.1.12)$$

avec  $G(\theta^*) = \frac{1}{4} \mathbb{E}[\partial_\theta q_0(\theta^*) \partial_\theta q_0(\theta^*)^\top]$  et  $F(\theta^*) = -\frac{1}{2} \mathbb{E}[\partial_{\theta^2}^2 q_0(\theta^*)]$ .

## 1.2 Sélection de modèles

### 1.2.1 Pourquoi sélectionner un modèle?

En pratique, pour une analyse, le statisticien dispose des données. À partir de ces données, il aimerait tirer des informations pour atteindre l'objectif de l'étude. L'extraction de ces informations passe par une phase de modélisation. Quel est le modèle le plus approprié pour ces données? Rappelons qu'une mauvaise modélisation conduira à des conclusions gravement trompeuses et donc à une prise de décision catastrophique puisque l'intérêt pratique de la Statistique est d'éclairer les décideurs. Comme il n'existe pas de modèle qui soit universellement adapté à toutes les données, le statisticien explore plusieurs modèles/algorithmes afin de mieux comprendre le processus ayant généré ces données ou d'obtenir une meilleure performance de prédiction.

Ainsi, il est beaucoup plus judicieux de considérer une grande famille de modèles pour réduire le risque de se tromper et de choisir le modèle le plus approprié de la collection.

La sélection de modèles est donc le processus de sélection d'un modèle dans une famille convenablement choisie à cet effet étant donné un jeu de données. Elle est fondamentale à toute analyse de données pour des fins d'inférence ou de prédiction fiable. Elle intervient donc dans tous les domaines utilisant les données notamment l'économie, l'ingénierie, l'écologie, la finance, les sciences politiques, la biologie, l'épidémiologie, etc.

### 1.2.2 Objectifs de la Sélection de modèles

L'objectif d'une procédure de sélection de modèles est en réalité le but de la modélisation. On en distingue principalement deux:

- L'identification ou l'interprétation: ici, l'on aimerait comprendre le processus de génération des données, interpréter la nature des données ou appuyer un modèle physique. L'on préférera donc les modèles assez simples et interprétables aux modèles complexes;
- La prédiction: Il ne s'agit plus de découvrir ou d'inférer des modèles beaucoup plus réels mais plutôt de répondre à des questions pratiques, comme: quelles températures il fera demain, quel sera le cours du bitcoin dans deux semaines, comment détecter une fraude lors d'une transaction, etc.

#### 1.2.2.a Identification

Considérons  $\mathcal{M}$ , la famille de modèles candidats. Puisque l'objectif est de sélectionner le vrai modèle (notée  $m^*$  dans tout ce travail) ayant généré les données, l'on fait l'hypothèse que  $\mathcal{M}$  contient  $m^*$ .

Une procédure de sélection de modèles vise à sélectionner le meilleur modèle de la famille noté  $\hat{m} := \hat{m}_n$  dans le but d'estimer  $m^*$ .

**Définition 1.6.** Une telle procédure sera dite **consistante** si elle arrive à retrouver  $m^*$  avec probabilité approchant 1 asymptotiquement. Plus formellement,

$$\mathbb{P}(\hat{m} = m^*) \xrightarrow[n \rightarrow \infty]{} 1. \quad (1.2.1)$$

Il existe aussi une version de la consistance avec une convergence presque sûre. Elle est dite **consistance forte**:

$$\hat{m} \xrightarrow[n \rightarrow +\infty]{a.s.} m^*. \quad (1.2.2)$$

La première est beaucoup plus classique dans la littérature et est facilement vérifiable à travers des études de Monte Carlo.

#### 1.2.2.b Prédiction

Bien que la consistance d'une procédure de sélection de modèles est une propriété mathématique convaincante, en pratique l'on ne connaît pas le vrai modèle et supposer que  $m^* \in \mathcal{M}$  est souvent irréaliste. Le vrai modèle peut être de dimension infinie et il est impossible d'identifier pareil modèle. Dans ce cas, ce qui pourrait être fait est une approximation de la réalité par des modèles de dimension finie qui sont évidemment tous faux. Et le but serait donc de sélectionner le modèle qui fait le moins possible d'erreur quand il s'agit de prédire. Ainsi, l'on va considérer une collection de modèles au plus dénombrable  $\mathcal{M}_n$  de cardinal dépendant de  $n$ . Pour un certain prédictor  $\theta$ , l'on mesure son erreur de prédiction par

$$R(\theta) := \mathbb{P}\gamma(\theta) = \mathbb{E}[\gamma(\theta, X_1)]$$

$$\text{avec } \gamma(\theta, X_t) := \frac{(X_t - f_\theta^t)^2}{H_\theta^t} + \log(H_\theta^t) \quad (1.2.3)$$

Le contraste  $\gamma$  est -2 fois la log-densité conditionnelle de  $X_t$  sachant le passé (en ayant considéré un bruit gaussien). D'après [Bardet and Wintenberger \(2009\)](#), la fonction  $R$

atteint son minimum en  $\theta^*$  qui est souvent non estimable directement soit parcequ'il est de dimension infinie ou parce que la famille de modèle candidat est mal spécifiée. Ainsi, il est naturel de définir un risque de prédiction relativement à  $\theta^*$

$$\ell(\theta, \theta^*) := R(\theta) - R(\theta^*) \geq 0. \quad (1.2.4)$$

Le meilleur choix possible dans la collection  $\mathcal{M}_n$  est appelé *oracle* (que nous noterons aussi  $m^*$  mais qui ne signifie plus vrai modèle) défini par

$$m^* \in \arg \inf_{m \in \mathcal{M}_n} \ell(\hat{\theta}_m, \theta^*). \quad (1.2.5)$$

L'oracle  $m^*$  est un idéal et n'est pas atteignable puisqu'il dépend de  $\theta^*$  et de la distribution du processus qui sont inconnus. Ainsi, une procédure de sélection de modèles peut avoir pour but de sélectionner un modèle  $\hat{m}$  qui imitera  $m^*$  en matière de risque.

**Définition 1.7.** Une procédure de sélection de modèles est dite *asymptotiquement efficiente* si le rapport de risque de  $\hat{m}$  et de  $m^*$  converge en probabilité vers 1 quand  $n \rightarrow \infty$ ; i.e.,

$$\frac{\ell(\hat{\theta}_{\hat{m}}, \theta^*)}{\ell(\hat{\theta}_{m^*}, \theta^*)} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} 1. \quad (1.2.6)$$

L'estimateur  $\hat{\theta}_m$  n'est plus le QMLE mais le minimiseur du risque empirique associée au contraste  $\gamma$ :

$$\hat{\theta}_m = \operatorname{argmin}_{\theta \in \Theta_m} \gamma_n(\theta) \quad \text{avec} \quad \gamma_n(\theta) = \frac{1}{n} \sum_{t=1}^n \gamma(\theta, X_t)$$

Notons que la version non asymptotique de cette définition est la propriété de plus en plus recherchée chez les théoriciens de la sélection de modèles: elle est appelée **inégalité oracle**

$$\ell(\hat{\theta}_{\hat{m}}, \theta^*) \leq C_1 \inf_{m \in \mathcal{M}_n} \{\ell(\hat{\theta}_m, \theta^*)\} + \frac{C_2}{n} \quad (1.2.7)$$

avec  $C_1$  proche de 1 et  $C_2 > 0$ . Le but est d'obtenir (1.2.7) sur un évènement de grande probabilité ou en espérance.

### 1.2.3 Stratégies de Sélection de modèles

Nous avons vu l'intérêt d'un recours à la sélection de modèles, les différents objectifs. Mais comment est opérée cette sélection? Cette sous-section se propose de répondre à cette question.

#### 1.2.3.a Tests Statistiques

Le problème de sélection de modèles peut être traité en utilisant un test d'hypothèses. Dans cette approche, l'on compare les modèles candidats deux à deux si ceux-ci sont imbriqués, sinon l'on pourrait préalablement les plonger dans un modèle plus grand les contenant. Ainsi, un test de Fisher (F-test) peut être réalisé pour départager les deux modèles. Pour faire simple, supposons que nous voulons sélectionner entre un AR(1) et un AR(2) ayant observé une trajectoire. Comme hypothèses du test, l'on peut considérer pour hypothèse nulle " $\theta^* \in \Theta_1 \subset ]-1, 1[$ " vs " $\theta^* \in \Theta_2 \subset ]-1, 1[\times]-1, 1[$ " l'hypothèse alternative. Un rejet de l'hypothèse nulle pourrait conduire à retenir l'AR(2) bien que cela

ne signifie pas que  $H_1$  est vraie. Ainsi, cela pourrait servir pour progressivement éliminer certains modèles.

Pour conserver un sous-ensemble de paramètres importants, l'on dispose également des tests comme le test de Wald, de student ou du rapport de vraisemblance.

Cette approche de sélection de modèles par test présente un certain nombre de défauts: avant d'effectuer le test, il faut se fixer un niveau de significativité; ainsi, deux analystes différents avec exactement les mêmes modèles et données peuvent arriver à des conclusions différentes basées sur des niveaux de significativité différents. Aussi, le nombre de F-test a effectué croît très vite.

### 1.2.3.b Validation croisée

La validation croisée est très utilisée notamment en apprentissage statistique/ Machine Learning quand il s'agit d'un objectif de prédiction. Plusieurs approches de la validation croisée existent ([Stone \(1974\)](#), [Allen \(1974\)](#)). L'on distingue principalement deux :

- La **testset validation** ou **hold out**: Elle consiste à diviser les données en deux: la première partie appelée données d'entraînement servira à calculer les estimations des paramètres des modèles et la deuxième appelée données test permettra d'évaluer les capacités prédictives des estimateurs obtenus lors de la phase d'apprentissage. Le choix de la proportion de la partie qui servira de test (généralement moins de 40%) influence beaucoup le résultat final;
- La  **$k$  fold validation** qui consiste à diviser les données en  $k$  blocs et d'utiliser  $k - 1$  blocs pour entraîner les estimateurs et le bloc restant pour le test. Cet entraînement est répété  $k$  fois de sorte que chaque bloc serve exactement une fois d'ensemble test. Contrairement au hold-out, le choix de  $k$  est certes important mais n'affectera pas sensiblement le résultat final.

En pratique, il est connu que le hold-out est relativement instable par rapport au  $k$  fold. Une question assez ouverte est de savoir comment choisir  $k$  et quelle est son influence sur les performances statistiques de la méthode? A cet effet, [Arlot \(2007\)](#) a montré qu'il y'a surpénalisation lorsque  $k$  est petit et d'un point de vue asymptotique, il est indispensable de faire tendre  $k$  vers l'infini avec  $n$  pour obtenir une procédure optimale.

Plusieurs autres études ont abordé cette approche de sélection de modèles dans le cadre de l'estimation de la densité (voir par exemple [Celisse et al. \(2014\)](#), [Celisse \(2008\)](#)).

### 1.2.3.c Pénalisation

Une approche assez populaire pour opérer une sélection de modèles est la pénalisation d'un certain critère empirique qui peut être -2 fois la log-vraisemblance ou le contraste des moindres carrés. Puisque le critère empirique est généralement décroissant en fonction de la complexité du modèle, le minimiser fournira le modèle le plus complexe qui sera difficilement interprétable et en plus n'aura pas nécessairement les meilleures capacités de prédiction en raison de sa grande variabilité due à l'estimation d'un très grand nombre de paramètres.

L'idée de pénaliser vient donc résoudre ce problème; en ajoutant un terme (croissant en fonction de la complexité), l'on défavorise les très gros modèles. A l'inverse, si ce terme est prépondérant devant le contraste empirique, l'on va privilégier les modèles beaucoup plus simples présentant un biais beaucoup plus important.

La pénalisation remonte aux années 1970 avec les travaux de [Akaike \(1969\)](#), [Mallows \(1973\)](#) et [Akaike \(1973\)](#). Bien qu'il soit probable que cette idée ait déjà existé dans

d'autres contextes tels que la sélection de sous-ensembles par [Beale et al. \(1967\)](#), [Hocking and Leslie \(1967\)](#), et la régression ridge par [Hoerl and Kennard \(1970\)](#).

En utilisant les moindres carrés ordinaires dans le cadre de régression, Mallows a obtenu le critère  $C_p$ . Parallèlement, Akaike a dérivé l'AIC pour l'estimation de la densité en utilisant le contraste log-vraisemblance. Quelques années plus tard, à la suite d'Akaike, [Schwarz \(1978\)](#) a proposé une approche alternative pour l'estimation de la densité et a dérivé le critère d'information bayésien (BIC).

Le terme de pénalité de ces critères est proportionnel à la dimension du modèle. Au cours des dernières décennies, différentes approches de pénalisation ont émergé telles que la pénalisation  $\mathbb{L}^2$  pour la régularisation ridge [Hoerl and Kennard \(1970\)](#), la pénalisation  $\mathbb{L}^1$  pour la procédure LASSO de [Tibshirani \(1996\)](#) et l'elastic net qui combine les normes  $\mathbb{L}^1$  et  $\mathbb{L}^2$  [Zou and Hastie \(2005\)](#).

Formellement, en considérant comme contraste empirique la quasi-vraisemblance, l'on définit un critère de pénalisation  $\hat{C}$  ainsi qu'il suit:

$$\hat{C}(m) = -2 \hat{L}_n(\hat{\theta}(m)) + \kappa_n(m), \quad (1.2.8)$$

pour tout modèle candidat  $m \in \mathcal{M}$ , avec  $\kappa_n$  est une suite croissante de la taille de modèle. Puisqu'il existe une multitude de choix pour  $\kappa_n$ , toute la difficulté de la sélection de modèles par pénalisation réside dans la spécification de ce terme de pénalité. Les choix les plus classiques sont entre autres:

- $\kappa_n(m) = 2c D_m \log \log n$  avec  $c > 1$ , on retrouve le critère de Hannan et Quinn [Hannan and Quinn \(1979\)](#);
- $\kappa_n = D_m \log n$ ,  $\hat{C}$  est le critère BIC [Schwarz \(1978\)](#);
- $\kappa_n = 2 D_m$ ,  $\hat{C}$  est l'AIC [Akaike \(1973\)](#).

Toutes ces pénalités ont été obtenues soit pour la prédiction ou l'interprétation. Rappelons aussi qu'il existe aussi des critères de pénalisation où le terme de pénalité est considéré multiplicativement par rapport au risque empirique; c'est l'exemple du FPE de Akaike.

Compte tenu de cette richesse de choix, comment un statisticien peut-il décider du critère à utiliser? C'est la raison pour laquelle les approches qui ne peuvent être ni mises en œuvre ni comprises par la communauté scientifique ne sont pas acceptées. Cela implique qu'il faut au moins une méthode qui puisse être implémentée facilement et donnent des résultats qui peuvent être interprétés par les utilisateurs. D'un point de vue statistique, les approches suffisamment générales pour traiter une grande variété de problèmes et vérifiant une propriété de sélection de modèles sont les plus prisées.

## 1.2.4 Dérivation de quelques critères dans le cadre des séries temporelles

### 1.2.4.a FPE Akaike (1969)

C'est l'un des tout premiers critères de pénalisation à être obtenus du moins pour les séries chronologiques (notamment pour un  $\text{AR}(p)$ ). Supposons avoir observé une trajectoire  $(X_1, \dots, X_n)$  d'un  $\text{AR}(p)$   $X_t = \sum_{i=1}^p \phi_i^* X_{t-i} + \sigma \xi_t$ . Puisque l'on ignore la valeur de  $p$ , alors quel ordre d'auto-régression choisir pour ajuster les données? L'idée d'Akaike est de choisir l'ordre qui minimise l'erreur quadratique moyenne lorsque le modèle estimé (au moyen de  $(X_1, \dots, X_n)$ ) est utilisé pour prédire la valeur  $Y_{n+1}$  d'une suite d'observation  $(Y_1, \dots, Y_n)$  provenant du même  $\text{AR}(p)$  mais indépendant de  $(X_1, \dots, X_n)$ .

Soit  $\hat{\phi} := (\hat{\phi}_1(X_1, \dots, X_n), \dots, \hat{\phi}_p(X_1, \dots, X_n))^\top$ , l'estimateur du QMLE de  $\phi^* = (\phi_1^*, \dots, \phi_p^*)^\top$ . D'après le TCL (1.1.12), on a

$$\sqrt{n}(\hat{\phi} - \phi^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, (G(\phi^*))^{-1}) \quad (1.2.9)$$

puisque  $G(\phi^*) = -F(\phi^*)$ . Pour  $i, j = 1, \dots, p$

$$\begin{aligned} (G(\phi^*))_{i,j} &= \frac{1}{4} \mathbb{E} \left[ \frac{\partial q_0(\phi^*)}{\partial \phi_i} \frac{\partial q_0(\phi^*)}{\partial \phi_j} \right] \\ &= \frac{1}{4} \mathbb{E} \left[ \frac{4}{\sigma^4} X_{-i} X_{-j} \xi_0^2 \right] = \frac{1}{\sigma^2} \mathbb{E}[X_i X_j]. \end{aligned}$$

Avec  $X = (X_1, \dots, X_n)^\top$  et  $Y = (Y_1, \dots, Y_n)^\top$ , l'erreur quadratique est donc:

$$\begin{aligned} \mathbb{E}[(Y_{n+1} - \hat{\phi}_1 Y_n - \dots - \hat{\phi}_p Y_{n+1-p})^2] &= \sigma^2 \mathbb{E}[\xi_{n+1}^2] + \mathbb{E} \left[ \left( \sum_{i=1}^p (\hat{\phi}_i - \phi_i^*) Y_{n+1-i} \right)^2 \right] \\ &= \sigma^2 + \mathbb{E}[(\hat{\phi} - \phi^*)^\top \mathbb{E}[Y Y^\top | X] (\hat{\phi} - \phi^*)] \\ &= \sigma^2 + \sigma^2 \mathbb{E}[(\hat{\phi} - \phi^*)^\top G(\phi^*) (\hat{\phi} - \phi^*)] \\ &\approx \sigma^2 + \sigma^2 \frac{1}{n} \text{Trace}(G(\phi^*) (G(\phi^*))^{-1}) \\ &= \sigma^2 \left( 1 + \frac{p}{n} \right) \end{aligned} \quad (1.2.10)$$

où l'approximation a été rendue possible grâce à (1.2.9). Puisque  $\sigma$  est en général inconnu, l'idée est de le substituer par un estimateur consistant. D'après [Brockwell and Davis \(1991\)](#), l'estimateur MLE de  $\sigma^2$  vérifie

$$n \frac{\hat{\sigma}^2}{\sigma^2} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \chi(n-p)$$

avec  $\hat{\sigma}^2 = \frac{1}{n} \sum_{t=1}^n \left( X_t - \sum_{i=1}^p \hat{\phi}_i X_{t-i} \right)^2$ . En remplaçant  $\sigma^2$  par l'estimateur  $\frac{n}{n-p} \hat{\sigma}^2$  dans la forme développée (1.2.10), Akaike définit le FPE pour un modèle candidat d'ordre  $p$  quelconque à être

$$\text{FPE}_p = \hat{\sigma}^2 \frac{n+p}{n-p}.$$

[Ing and Wei \(2005\)](#) ont démontré que le FPE possède la propriété d'efficacité asymptotique.

#### 1.2.4.b AIC Akaike (1973)

L'AIC a été obtenu par [Akaike \(1973\)](#). Il fut dérivé en recherchant un estimateur non biaisé de la divergence de Kullback entre le vrai modèle et les modèles candidats. Supposons observé  $X_1, \dots, X_n$  d'un ARMA( $p, q$ )

$$X_t - \sum_{i=1}^p \phi_i^* X_{t-i} = \xi_t + \sum_{i=1}^q \beta_i^* \xi_{t-i}$$

où  $(\xi_t)$  est un bruit blanc gaussien. Posons  $\theta^* = (\phi^*, \beta^*)$ , la divergence de Kullback-Leibler entre la loi de  $X_0$  (sachant tout le passé) suivant le vrai modèle (représenté par  $\theta^*$ ) et un modèle candidat est

$$\mathcal{K}(\theta^* | \theta) = \mathbb{E} \left[ -2 \ln \left( \frac{g(X_0, \theta)}{g(X_0, \theta^*)} \right) \right]$$

où l'espérance est prise sous  $\theta^*$  et  $g(X_0, \theta)$  représente la densité de  $X_0$  sachant  $X_{-1}, X_{-2}, \dots$  sous  $\theta$ . L'objectif est de construire un estimateur non biaisé de  $\mathcal{K}(\theta^*|\theta)$ . Puisque cette pseudo distance est positive et s'annule uniquement en  $\theta = \theta^*$ , minimiser cet estimateur non biaisé nous conduira au modèle le "plus proche" du vrai.

Pour simplifier, supposons que les modèles compétitifs sont aussi gaussiens. Puisque  $\mathbb{E}[-2 \ln(g(X_0, \theta^*))]$  est indépendant du modèle candidat choisi, nous ignorerons ce terme dans l'expression de  $\mathcal{K}(\theta^*|\theta)$ . Donc, un estimateur naturel de cette approximation sera :

$$-\frac{2}{n} \sum_{t=1}^n \ln(g(X_t, \theta)) = -\frac{2}{n} L_n(\theta).$$

Considérons  $\hat{\theta}$  l'estimateur MLE de  $\theta^*$ . Sous l'hypothèse de différentiabilité de  $L_n$  et en prenant  $\hat{\theta}$  comme maximum local, un développement de Taylor au second ordre nous donne

$$L_n(\theta^*) \approx L_n(\hat{\theta}) + \frac{1}{2} (\theta^* - \hat{\theta})^\top \partial_{\theta^2}^2 L_n(\hat{\theta}) (\theta^* - \hat{\theta}). \quad (1.2.11)$$

En vertu du théorème ergodique et de la continuité uniforme (voir [Bardet and Wintenberger \(2009\)](#)), l'on a:

$$\frac{1}{n} \partial_{\theta^2}^2 L_n(\hat{\theta}) \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta^*) = -\frac{1}{2} \mathbb{E}[\partial_{\theta^2}^2 q_0(\theta^*)]$$

avec  $q_t(\theta) = -2 \ln(g(X_t, \theta))$ . Le TCL (1.1.12) s'écrit ici

$$\sqrt{n}(\hat{\theta} - \theta^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}\left(0, 2(\mathbb{E}[\nabla_{\theta}^2 q_0(\theta^*)])^{-1}\right) \quad (1.2.12)$$

puisque  $G(\theta^*) = \frac{1}{4} \mathbb{E}[\partial_{\theta} q_0(\theta^*) \partial_{\theta} q_0(\theta^*)^\top] = -F(\theta^*)$ .

Ainsi, l'on obtient en prenant l'espérance de chaque terme en (1.2.11)

$$\begin{aligned} \mathbb{E}[L_n(\hat{\theta})] &\approx \mathbb{E}[L_n(\theta^*)] - \frac{1}{2} \text{Trace}\left(F(\theta^*) 2(\mathbb{E}[\nabla_{\theta}^2 q_0(\theta^*)])^{-1}\right) \\ &= \mathbb{E}[L_n(\theta^*)] + \frac{1}{2} (p + q). \end{aligned} \quad (1.2.13)$$

En effectuant également un développement de  $\mathbb{E}[q_0(\theta)]$  autour de  $\theta^*$ , il vient que

$$\mathbb{E}[q_0(\hat{\theta})] \approx \mathbb{E}[q_0(\theta^*)] + \frac{1}{2} (\hat{\theta} - \theta^*)^\top \mathbb{E}[\nabla_{\theta}^2 q_0(\theta^*)] (\hat{\theta} - \theta^*),$$

le terme de premier ordre étant nul puisque  $\theta^*$  est l'unique minimum de la divergence de Kullback. Ainsi,

$$\mathbb{E}[\mathbb{E}[q_0(\hat{\theta})]] \approx \mathbb{E}[q_0(\theta^*)] + \frac{p + q}{n}. \quad (1.2.14)$$

En prenant l'espérance de la différence entre les termes de (1.2.13) et (1.2.14), on a

$$\begin{aligned} \mathbb{E}\left[-\frac{2}{n} L_n(\hat{\theta}) - \mathbb{E}[q_0(\hat{\theta})]\right] &\approx \mathbb{E}\left[-\frac{2}{n} L_n(\theta^*)\right] - \mathbb{E}[q_0(\theta^*)] - \frac{p + q}{n} - \frac{p + q}{n} \\ &= -2 \frac{(p + q)}{n} \end{aligned}$$

de sorte qu'un estimateur sans biais (non normalisé) de  $\mathbb{E}[q_0(\hat{\theta})]$  est

$$AIC = -2 L_n(\hat{\theta}) + 2(p + q).$$

L'optimalité asymptotique de l'AIC a été prouvée par [Shibata \(1980\)](#) et aussi par [Ing and Wei \(2005\)](#).



#### 1.2.4.c BIC Schwarz (1978)

Le critère BIC est l'un des plus connus en sélection de modèles. Contrairement aux autres critères, sa dérivation s'est faite dans un cadre bayésien: le paramètre  $\theta^*$  associé au vrai modèle  $m^*$  est supposé aléatoire. Etant donné une famille de modèle  $\mathcal{M}$ , un modèle  $m \in \mathcal{M}$  est tiré selon une distribution a priori  $\pi_m$ . Très souvent cet a priori est non informatif car pris comme distribution uniforme ou bien ignoré dans la formule d'approximation du BIC. Ainsi, conditionnellement au modèle  $m$  tiré, un paramètre  $\theta \in \Theta_m$  est tiré suivant une autre distribution  $\mu_m$ . Le choix du meilleur modèle dans la famille  $\mathcal{M}$  est opéré en considérant le modèle maximisant la probabilité a posteriori i.e.,

$$\begin{aligned}\hat{m} &= \operatorname{argmax}_{m \in \mathcal{M}} \{ \mathbb{P}(m | X_1, X_2, \dots, X_n) \} \\ &= \operatorname{argmax}_{m \in \mathcal{M}} \left\{ \frac{\pi_m \mathbb{P}(X_1, X_2, \dots, X_n | m)}{\mathbb{P}(X_1, X_2, \dots, X_n)} \right\}.\end{aligned}$$

Le terme prépondérant dans cette formule est bien  $\mathbb{P}(X_1, X_2, \dots, X_n | m)$  qui peut encore s'écrire en intégrant sur l'espace des paramètres

$$\begin{aligned}\mathbb{P}(X_1, X_2, \dots, X_n | m) &= \int_{\Theta_m} \mathbb{P}(X_1, X_2, \dots, X_n | \theta, m) d\mu_m(\theta) \\ &= \int_{\Theta_m} \exp(L_n(\theta)) d\mu_m(\theta).\end{aligned}$$

Cette intégrale est généralement impossible à calculer et est approchée au moyen de la formule de l'approximation de Laplace (Lebarbier and Mary-Huard (2004)). Ainsi, l'on obtient

$$\log \mathbb{P}(X_1, X_2, \dots, X_n | m) \approx L_n(\hat{\theta}_m) - \frac{\log(n)}{2} |m| + \frac{\log(2\pi)}{2} |m| - \frac{1}{2} \log(\det(-\partial_{\theta}^2 L_n(\hat{\theta}))) + O(n^{-1}).$$

Et en prenant  $\pi_m = 1/|\mathcal{M}|$ , il vient que

$$\begin{aligned}\log(\pi_m \mathbb{P}(X_1, X_2, \dots, X_n | m)) &\approx L_n(\hat{\theta}_m) - \frac{\log(n)}{2} |m| \\ &\quad \underbrace{- \log(|\mathcal{M}|) + \frac{\log(2\pi)}{2} |m| - \frac{1}{2} \log(\det(-\partial_{\theta}^2 L_n(\hat{\theta}))) + O(n^{-1})}_{O_P(1)}.\end{aligned}$$

En négligeant tous ces termes qui restent bornés (et le terme  $\mathbb{P}(X_1, X_2, \dots, X_n)$  indépendant du choix de  $m$ ), il vient que maximiser la probabilité a posteriori revient à maximiser  $L_n(\hat{\theta}_m) - \frac{\log(n)}{2} |m|$ . Ce qui conduit Schwartz à définir le BIC à être

$$BIC_m = -2 L_n(\hat{\theta}_m) + \log(n) |m|$$

et donc à choisir le modèle  $\hat{m}$  comme celui qui minimise le BIC.

Plusieurs études ont montré que si le vrai modèle  $m^*$  est inclus dans la famille des modèles candidats  $\mathcal{M}$ , alors  $\hat{m}$  obtenu par minimisation du BIC est consistant (voir par exemple Csiszár and Shields (1999), Csiszár et al. (2000), Garivier (2006), Bardet et al. (2012)).

#### 1.2.4.d HQ Hannan and Quinn (1979)

Il a été obtenu par Hannan and Quinn (1979) dans le cadre de l'estimation de l'ordre à considérer lorsqu'un modèle linéaire autorégressif est ajusté aux données. Le terme de pénalité en  $\log \log n$  provient d'une application de la loi des logarithmes itérés (LIL) à la série des autocorrélations partielles. Il est défini par

$$HQ = -2 L_n(\hat{\theta}) + 2 c p \log \log n$$

où  $c > 1$ . En effet, Hannan and Quinn (1979) ont argué que le  $\log n$  dans la pénalité BIC qui assure la consistance dans la sélection du vrai ordre, n'est pas la suite qui croît la plus lentement possible et ont proposé le  $\log \log n$  qui vérifie la propriété de forte consistance Hannan and Quinn (1979), Hannan and Deistler (2012). D'après eux, le BIC peut très souvent choisir des modèles très simples, éventuellement erronés, pour de petits échantillons.

#### 1.2.5 Calibration de la constante multiplicative

Dans la littérature de la sélection de modèles non asymptotique, il est fréquent de rencontrer les pénalités dépendantes d'une certaine constante universelle inconnue:

$$\text{pen}(m) = \kappa \text{pen}_{\text{shape}}(m)$$

où  $\text{pen}_{\text{shape}}$  est une fonction de la dimension du modèle

- $D_m$  pour les modèles linéaires gaussiens Birgé and Massart (2001) ;
- $D_m (1 + \sqrt{2 L_m})$  Birgé and Massart (2007b) dans un cadre gaussien plus général;
- $D_m (1 + c \log(n/D_m))$  pour la détection de ruptures Lebarbier (2005);
- et bien d'autres.

Déterminer la valeur de la constante  $\kappa$  est fondamentale avant la mise en œuvre de la procédure de sélection. Massart (2007) a proposé une méthode efficace pour calibrer  $\kappa$  au moyen des données uniquement. Il s'agit de l'*heuristique de pente* qui se décline en deux algorithmes classiques:

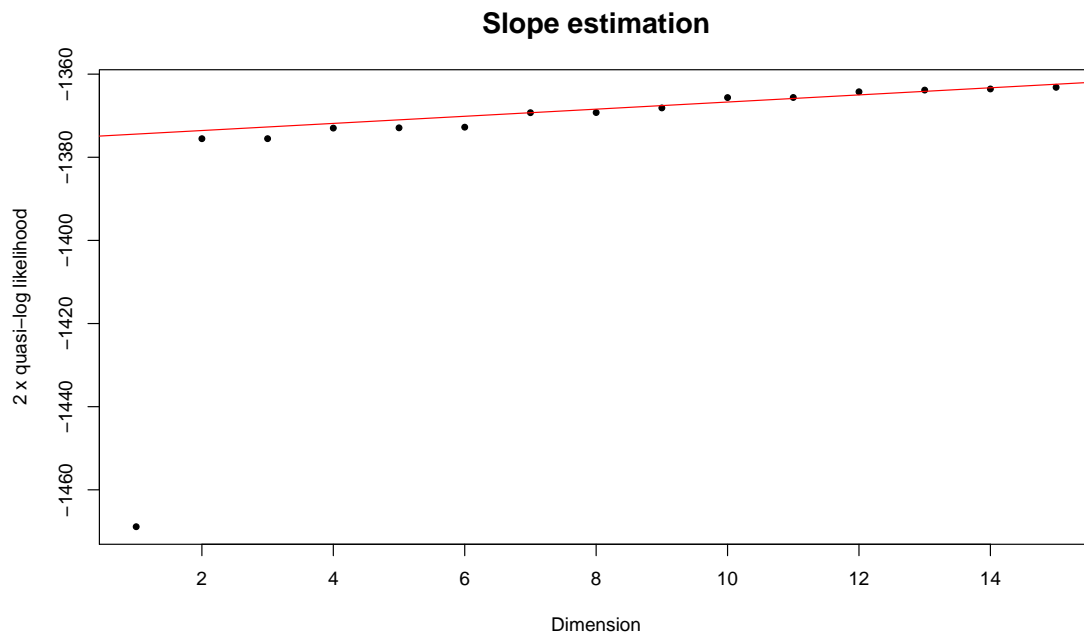
- *slope estimation*: elle consiste à prendre comme une estimation de  $\kappa$  deux fois la valeur de la pente de croissance du critère empirique en fonction de la dimension  $D_m$  pour les modèles complexes (voir Figure 1.1);
- *saut de dimension* qui consiste à considérer plutôt 2 fois la valeur de  $\kappa$  (en ayant préalablement choisi une grille de valeurs pour  $\kappa$ ) qui donne le plus grand saut de la complexité (voir Figure 1.2).

### 1.3 Synthèse des travaux

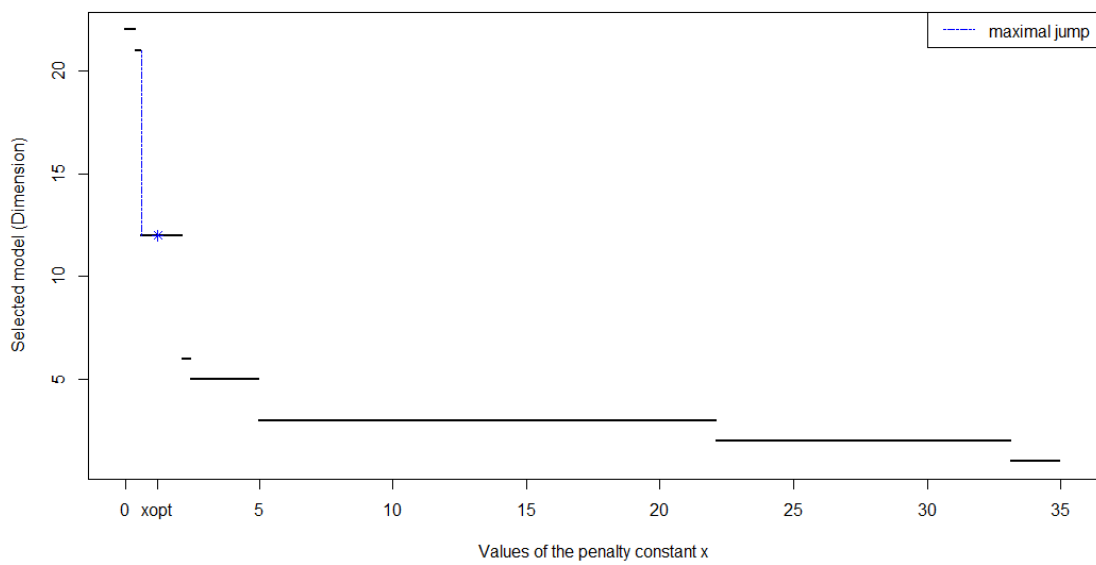
#### 1.3.1 Critères de Sélection de modèles consistants et test d'ajustement pour les modèles affines causaux

Supposons avoir observé une trajectoire  $(X_1, X_2, \dots, X_n)$  du processus affine causal

$$X_t = M_{\theta^*}((X_{t-i})_{i \in \mathbb{N}^*}) \xi_t + f_{\theta^*}((X_{t-i})_{i \in \mathbb{N}^*}) \quad \text{pour tout } t \in \mathbb{Z} \quad (1.3.1)$$



**Figure 1.1:** Courbe de la quasi log-vraisemblance d'un AR(2) en fonction de la dimension; la pente de la droite oblique vaut  $\hat{a} = 0.86$ .



**Figure 1.2:** Dimension Jump

où  $\theta^* \in \mathbb{R}^d$ . Le modèle (1.3.1) ayant généré les données sera noté  $m^*$  et appartient à une certaine famille finie de modèles affines causaux  $\mathcal{M}$ . Notre objectif est de construire des critères de sélection vérifiant la propriété de consistance (1.2.1) et de tester la qualité d'ajustement du modèle sélectionné aux données. Pour ce faire, nous construisons un estimateur  $\hat{m}$  de  $m^*$  par pénalisation de la quasi-vraisemblance:

$$\hat{m} = \underset{m \in \mathcal{M}}{\operatorname{argmin}} \hat{C}(m) \quad \text{avec} \quad \hat{C}(m) = -2 \hat{L}_n(\hat{\theta}(m)) + |m| \kappa_n, \quad (1.3.2)$$

A la suite de Bardet et al. (2012), l'idée est de trouver des conditions sur la suite  $(\kappa_n)$  afin de garantir la convergence (1.2.1); lesquelles conditions devraient dépendre des coefficients de Lipschitz des fonctions  $M_\theta$  et  $f_\theta$ .

Avant d'énoncer les principaux résultats obtenus, nous donnons un résultat intermédiaire pratique qui dit qu'il est toujours possible de plonger deux modèles affines causaux paramétriques dans un plus grand modèle affine causal.

**Proposition 1.2.** *Soit  $d_1, d_2 \in \mathbb{N}$ ,  $\Theta_1 \subset \mathbb{R}^{d_1}$  et  $\Theta_2 \subset \mathbb{R}^{d_2}$ , et pour  $i = 1, 2$ , définissons  $f_{\theta_i}^{(i)}, M_{\theta_i}^{(i)} : \mathbb{R}^\infty \rightarrow \mathbb{R}$  et pour  $\theta_i \in \Theta_i$ . Alors, il existe  $\max(d_1, d_2) \leq d \leq d_1 + d_2$ ,  $\Theta \subset \mathbb{R}^d$ , et une famille de fonctions  $f_\theta : \mathbb{R}^\infty \rightarrow \mathbb{R}$  et  $M_\theta : \mathbb{R}^\infty \rightarrow [0, \infty)$  avec  $\theta \in \Theta$ , tels que pour tout  $\theta_1 \in \Theta_1$  et  $\theta_2 \in \Theta_2$ , il existe  $\theta \in \Theta$  vérifiant*

$$\mathcal{AC}(M_{\theta_1}^{(1)}, f_{\theta_1}^{(1)}) \cup \mathcal{AC}(M_{\theta_2}^{(2)}, f_{\theta_2}^{(2)}) \subset \mathcal{AC}(M_\theta, f_\theta).$$

Ainsi dans toute la suite,  $\Theta$  sera l'espace des paramètres contenant tous les  $\Theta_m$  avec  $m \in \mathcal{M}$ .

### 1.3.1.a Hypothèses

Outre les hypothèses de régularité ayant permis d'obtenir le TCL (1.1.12), nous considérerons la condition suivante nous garantissant une relation entre la suite  $\kappa_n$  et la vitesse de décroissance des coefficients de Lipschitz de  $M_\theta$  et  $f_\theta$ .

**Hypothèse  $K(\Theta)$ :** *Sous les hypothèses  $A_i(\Psi_\theta, \Theta)$ , avec  $i = 0, 1$ ,  $\Psi_\theta = f_\theta, M_\theta$  et s'il existe  $r \geq 2$  tel que  $\theta^* \in \Theta(r)$ . Par ailleurs, avec  $s = \min(1, r/3)$ , supposons que la suite  $(\kappa_n)_{n \in \mathbb{N}}$  satisfasse*

$$\sum_{k \geq 1} \left( \frac{1}{\kappa_k} \right)^s \left( \sum_{j \geq k} \alpha_j^0(f_\theta, \Theta) + \alpha_j^0(M_\theta, \Theta) + \alpha_j^1(f_\theta, \Theta) + \alpha_j^1(M_\theta, \Theta) \right)^s < \infty.$$

Pour les modèles de type GARCH, nous considérerons plutôt

**Hypothèse  $\widetilde{K}(\Theta)$ :** *Sous les hypothèses  $A_i(\Psi_\theta, \Theta)$ , avec  $i = 0, 1$ ,  $\Psi_\theta = f_\theta, M_\theta$  et supposons qu'il existe  $r \geq 2$  tel que  $\theta^* \in \Theta(r)$ . Par ailleurs, avec  $s = \min(1, r/4)$ , supposons que la suite  $(\kappa_n)_{n \in \mathbb{N}}$  satisfasse*

$$\sum_{k \geq 1} \left( \frac{1}{\kappa_k} \right)^s \left( \sum_{j \geq k} \alpha_j^0(\widetilde{H}_\theta, \Theta) + \alpha_j^1(\widetilde{H}_\theta, \Theta) \right)^s < \infty.$$

**Remarque 1.1.** *Ces conditions sur  $(\kappa_n)_{n \in \mathbb{N}}$  ne sont pas restrictives : par exemple, si les coefficients de Lipschitz de  $f_\theta, M_\theta$  (le cas utilisant  $\widetilde{H}_\theta$  peut être traité de manière similaire) et leurs dérivées sont bornées par une décroissance géométrique ou riemannienne, l'on a:*

1. *Cas géométrique:*  $\alpha_j^0(f_\theta, \Theta) + \alpha_j^0(M_\theta, \Theta) + \alpha_j^1(f_\theta, \Theta) + \alpha_j^1(M_\theta, \Theta) = O(a^j)$  avec  $0 \leq a < 1$ , alors tout  $(\kappa_n)$  telle  $1/\kappa_n = o(1)$  peut être choisie; par exemple  $\kappa_n = \log n$  ou  $\log(\log n)$ ; C'est le cas pour les ARMA, GARCH, APARCH or ARMA-GARCH.

2. *Cas Riemannien:*  $\alpha_j^0(f_\theta, \Theta) + \alpha_j^0(M_\theta, \Theta) + \alpha_j^1(f_\theta, \Theta) + \alpha_j^1(M_\theta, \Theta) = O(j^{-\gamma})$  avec  $\gamma > 1$ :

- si  $r \geq 3$  alors
  - si  $\gamma > 2$  alors toute suite  $(\kappa_n)$  vérifiant  $1/\kappa_n = o(1)$  peut être choisie;
  - si  $1 < \gamma < 2$ , toute suite vérifiant  $(\kappa_n)$  telle que  $\kappa_n = O(n^\delta)$  avec  $\delta > 2 - \gamma$  peut être choisie.
- si  $1 \leq r < 3$ 
  - si  $\gamma > (r + 3)/r$  alors toute suite  $(\kappa_n)$  satisfaisant  $1/\kappa_n = o(1)$  peut être choisie;
  - si  $1 < \gamma < (r + 3)/r$  alors toute suite  $(\kappa_n)$  vérifiant  $\kappa_n = n^\delta$  avec  $\delta > (r + 3)/r - \gamma$  peut être choisie.

Dans le dernier cas de ces deux conditions sur  $r$ , nous pouvons voir que le choix habituel du BIC,  $\kappa_n = \log n$  ne vérifie pas l'hypothèse en général.

### 1.3.1.b Résultats théoriques

**Théorème 1.2.** *Considérons  $(X_1, \dots, X_n)$  une trajectoire observée d'un processus affine causal  $X$  appartenant à  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$  (ou  $\widehat{\mathcal{AC}}(\tilde{H}_{\theta^*})$ ) où  $\theta^*$  est un paramètre inconnu du compact  $\Theta$  contenu dans  $\Theta(r) \subset \mathbb{R}^d$  (ou  $\tilde{\Theta}(r) \subset \mathbb{R}^d$ ) avec  $r \geq 4$ . Si les hypothèses **A1-A3** et  $K(\Theta)$  (ou  $\tilde{K}(\Theta)$ ),  $A_2(f_\theta, \Theta)$  et  $A_2(M_\theta, \Theta)$  (ou  $A_2(\tilde{H}_\theta, \Theta)$ ) sont satisfaites, alors*

$$\mathbb{P}(\hat{m} = m^*) \xrightarrow[n \rightarrow \infty]{} 1 \quad \text{et} \quad \hat{\theta}(\hat{m}) \xrightarrow[n \rightarrow \infty]{\mathcal{P}} \theta^*. \quad (1.3.3)$$

Nous avons aussi établi la normalité asymptotique du QMLE du modèle choisi.

**Théorème 1.3.** *Sous les hypothèses du Théorème 1.2 et si  $\theta^* \in \overset{\circ}{\Theta}$ , alors*

$$\sqrt{n} \left( (\hat{\theta}(\hat{m}))_i - (\theta^*)_i \right)_{i \in m^*} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}_{|m^*|}(0, F(\theta^*, m^*)^{-1} G(\theta^*, m^*) F(\theta^*, m^*)^{-1}) \quad (1.3.4)$$

où  $(F(\theta^*, m^*))_{i,j} = -\frac{1}{2} \mathbb{E} \left[ \frac{\partial^2 q_0(\theta^*)}{\partial \theta_i \partial \theta_j} \right]$  et  $(G(\theta^*, m^*))_{i,j} = \frac{1}{4} \mathbb{E} \left[ \frac{\partial q_0(\theta^*)}{\partial \theta_i} \frac{\partial q_0(\theta^*)}{\partial \theta_j} \right]$  pour  $i, j \in m^*$ .

### 1.3.1.c Test Portmanteau

Bien que les pénalités obtenues nous garantissent la convergence des critères, il est également important de vérifier si le modèle choisi est approprié. Cette section tente de répondre à cette question en construisant un test portmanteau comme outil de diagnostic basé sur la séquence des carrés des résidus du modèle choisi. Pour tout  $m \in \mathcal{M}$ , et  $K$  un entier naturel non nul, soit le vecteur corrélogramme des carrés des résidus suivant

$$\hat{\rho}(m) := (\hat{\rho}_1(m), \dots, \hat{\rho}_K(m))',$$

où pour  $k = 1, \dots, K$ ,  $\hat{\rho}_k(m) := \frac{\hat{\gamma}_k(m)}{\hat{\gamma}_0(m)}$  avec

$$\hat{\gamma}_k(m) := \frac{1}{n} \sum_{t=k+1}^n (\hat{e}_t^2(m) - 1)(\hat{e}_{t-k}^2(m) - 1) \quad \text{et} \quad \hat{e}_t(m) := (\widehat{M}_{\hat{\theta}(m)}^t)^{-1} (X_t - \hat{f}_{\hat{\theta}(m)}^t).$$

Le résultat suivant fournit un TCL pour  $\widehat{\rho}(m^*)$  et  $\widehat{\rho}(\widehat{m})$  et établit la distribution asymptotique de la statistique de test.

**Théorème 1.4.** *Sous les hypothèses du Théorème 1.3, avec en outre*

- $\mathbb{E}[\xi_0^3] = 0$ ;

- 

$$\sum_{t=1}^{\infty} t^{-1/4} \left( \sum_{j \geq t} \alpha_j^0(f_\theta, \Theta) + \alpha_j^0(M_\theta, \Theta) \right)^{1/2} < \infty$$

$$\text{ou} \quad \sum_{t=1}^{\infty} t^{-1/4} \left( \sum_{j \geq t} \alpha_j^0(\widetilde{H}_\theta, \Theta) \right)^{1/2} < \infty,$$

alors,

1. avec  $V(\theta^*, m^*)$  définie en (1.3.7), on a

$$\sqrt{n} \widehat{\rho}(m^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}_K(0, V(\theta^*, m^*)). \quad (1.3.5)$$

2. avec  $\widehat{Q}_K(m^*) := n \widehat{\rho}(m^*)' (V(\widehat{\theta}(m^*), m^*))^{-1} \widehat{\rho}(m^*)$ , nous avons

$$\widehat{Q}_K(m^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \chi^2(K). \quad (1.3.6)$$

3. Les points 1. et 2. demeurent vraies lorsque  $m^*$  est remplacé par  $\widehat{m}$ .

$$\begin{aligned} V(\theta^*, m^*) := & I_K + (\mu_4 - 1)^{-2} J_K(m^*) F(\theta^*, m^*)^{-1} G(\theta^*, m^*) F(\theta^*, m^*)^{-1} J'_K(m^*) \\ & - 2(\mu_4 - 1)^{-1} J_K(m^*) F(\theta^*, m^*)^{-1} J'_K(m^*). \end{aligned} \quad (1.3.7)$$

Le Théorème 1.4 permet de tester asymptotiquement:

$$\begin{cases} H_0 : \exists m^* \in \mathcal{M}, \text{ tel que } (X_1, \dots, X_n) \text{ est une trajectoire de } X \in \mathcal{AC}(M_{\theta^*}, f_{\theta^*}) \\ H_1 : \nexists m^* \in \mathcal{M}, \text{ tel que } (X_1, \dots, X_n) \text{ est une trajectoire de } X \in \mathcal{AC}(M_{\theta^*}, f_{\theta^*}) \end{cases}$$

avec  $\theta^* \in \Theta(m^*)$  dans les deux cas.

### 1.3.2 Généralisation du critère de Hannan et Quinn pour les séries affines causales

Nous considérons toujours le problème de sélection de modèles (1.3.2). Nous avons vu précédemment qu'il est important de faire dépendre la suite  $(\kappa_n)$  des coefficients de Lipschitz de  $M_\theta$  et  $f_\theta$ . Une limite à cela est d'intérêt pratique: puisque nous donnons un large éventail de choix de la pénalité, comment l'analyste choisit-il sa pénalité? En décroissance exponentielle par exemple, parmi  $\log(\log n)$ ,  $\log n$  ou même  $\sqrt{n}$ , laquelle est la plus optimale?

Dans cette nouvelle contribution, nous établissons pour les processus à mémoire pas trop longue typiquement

$$\alpha_j^0(f_\theta, \Theta) + \alpha_j^0(M_\theta, \Theta) + \alpha_j^1(f_\theta, \Theta) + \alpha_j^1(M_\theta, \Theta) = O(j^{-\gamma}) \quad \text{avec} \quad \gamma > 2, \quad (1.3.8)$$

que l'on peut toujours trouver une constante  $c$  telle que avec  $\kappa_n = 2c \log(\log n)$ , l'on ait la propriété de forte consistance (1.2.2). La constante  $c$  est connue pour les modèles classiques ARMA, GARCH, APARCH mais pour les modèles complexes i.e. combinaison des ARMA et GARCH, il est difficile de trouver la valeur théorique de  $c$ . Cependant, les algorithmes de calibration de la constante multiplicative peuvent être utilisés pour une estimation adaptative de  $c$  dans les cas complexes.

Ainsi, la condition  $\mathbf{K}(\Theta)$  devient

**Condition  $\mathbf{K}(\Theta)$ :**

$$\sum_{k \geq e} \frac{1}{\log \log k} \sum_{j \geq k} \alpha_j^0(f_\theta, \Theta) + \alpha_j^0(M_\theta, \Theta) + \alpha_j^1(f_\theta, \Theta) + \alpha_j^1(M_\theta, \Theta) < \infty.$$

Et pour les modèles de type GARCH, l'on considèrera

**Condition  $\tilde{\mathbf{K}}(\Theta)$ :**

$$\sum_{k \geq e} \frac{1}{\log \log k} \sum_{j \geq k} \alpha_j^0(H_\theta, \Theta) + \alpha_j^1(H_\theta, \Theta) < \infty.$$

Il est facile de voir que tout processus satisfaisant (1.3.8) vérifie  $\mathbf{K}(\Theta)$ .

Pour obtenir un résultat assez général (avec les ARMA-GARCH inclus), il nous a été commode de supposer une relation entre la matrice d'information de Fisher  $G(\theta_m^*)$  et la matrice limite de la matrice hessienne de la log-likelihood  $F(\theta_m^*)$

$$(F(\theta_m^*))_{i,j} = -\frac{1}{2} \mathbb{E} \left[ \frac{\partial^2 q_0(\theta_m^*)}{\partial \theta_i \partial \theta_j} \right] \quad \text{et} \quad (G(\theta_m^*))_{i,j} = \frac{1}{4} \mathbb{E} \left[ \frac{\partial q_0(\theta_m^*)}{\partial \theta_i} \frac{\partial q_0(\theta_m^*)}{\partial \theta_j} \right],$$

avec  $\theta_m^* := (\theta^*, 0, \dots, 0)^\top \in \Theta(m)$ .

**A4:** Il existe des constantes  $\alpha_1$  et  $\alpha_2$  telles que pour tout  $m \in \mathcal{M}$  vérifiant  $m^* \subset m$ ,

$$\mathbf{1}_m^\top \Sigma_{\theta_m^*} \mathbf{1}_m = \alpha_1 D_m^1 + \alpha_2 D_m^2 \quad (1.3.9)$$

où  $D_m^1$  et  $D_m^2$  sont tels que  $D_m^1 + D_m^2 = D_m$ ,  $\mathbf{1}_m := (1, 1, \dots, 1)^\top \in \mathbb{R}^{D_m}$  et  $\Sigma_{\theta_m^*} := G(\theta_m^*)^{1/2} F(\theta_m^*)^{-1} G(\theta_m^*)^{1/2}$ .

Nous verrons que **A4** est vérifiée pour tous les modèles qui ne sont pas combinaison de ARMA et GARCH.

### 1.3.2.a Résultats théoriques

La proposition suivante suggère l'existence d'un terme qui sera capital pour notre résultat principal.

**Proposition 1.3.** *Considérons  $(X_1, X_2, \dots, X_n)$  une trajectoire observée d'un processus affine causal  $m^*$ . Pour tout modèle  $m$  vérifiant  $\theta_m^* \in \widetilde{\Theta(m)}$ , et si **A1-A4** sont vérifiées, alors on a*

$$\limsup_{n \rightarrow \infty} \frac{\widehat{L}_n(\widehat{\theta}(m)) - \widehat{L}_n(\theta_m^*)}{2 \log \log n} = \frac{1}{4} (\alpha_1 D_m^1 + \alpha_2 D_m^2) \quad p.s. \quad (1.3.10)$$

Soit  $m \in \mathcal{M}$ , désignons par  $c_{\min}(m)$  la quantité

$$c_{\min}(m) := \frac{1}{4} (\alpha_1 D_m^1 + \alpha_2 D_m^2) \quad (1.3.11)$$

Nous énonçons maintenant un résultat qui donne les valeurs de  $\alpha_1$  et de  $\alpha_2$  pour la plupart des modèles causaux affines classiques.

**Proposition 1.4.** *Sous les hypothèses et notation de la Proposition 1.3, nous avons*

- Si  $\mu_4 = \mathbb{E}[\xi_0^4] = 3$  (bruit gaussien par exemple), alors  $\alpha_1 = 2$ ,  $\alpha_2 = 2$  et  $c_{\min}(m) = \frac{1}{2} D_m$ ;
- Si le paramètre  $\theta$  identifiant un modèle affine causal  $X_t = M_\theta^t \xi_t + f_\theta^t$  peut s'écrire  $\theta = (\theta_1, \theta_2)'$  avec  $f_\theta^t = \tilde{f}_{\theta_1}^t$  et  $M_\theta^t = \tilde{M}_{\theta_2}^t$ , alors  $\alpha_1 = 2$ ,  $\alpha_2 = \mu_4 - 1$  et

$$c_{\min}(m) = \frac{1}{2} D_m^1 + \frac{\mu_4 - 1}{4} D_m^2$$

La seconde configuration de la Proposition 1.4 contient la plupart des séries temporelles classiques

- modèles GARCH( $p, q$ ), APARCH( $\delta, p, q$ ) et associés,  $c_{\min}(m) = \frac{\mu_4 - 1}{4} D_m$ ;
- modèles ARMA( $p, q$ ),  $c_{\min}(m) = \frac{D_m}{2}$  si la variance du bruit est connue et  $c_{\min}(m) = \frac{D_m - 1}{2} + \frac{\mu_4 - 1}{4}$  sinon.

Le principal résultat de notre contribution peut donc être énoncé.

**Théorème 1.5.** *Considérons  $(X_1, \dots, X_n)$  une trajectoire observée d'un processus affine causal  $X$  appartenant à  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$  (ou  $\tilde{\mathcal{AC}}(\tilde{H}_\theta)$ ) où  $\theta^*$  est un paramètre inconnu du compact  $\Theta$  contenu dans  $\Theta(r) \subset \mathbb{R}^d$  (ou  $\tilde{\Theta}(r) \subset \mathbb{R}^d$ ) avec  $r \geq 8$ . Supposons les hypothèses **A1-A3**,  $K(\Theta)$  (ou  $\tilde{K}(\Theta)$ ),  $A_2(f_\theta, \Theta)$  et  $A_2(M_\theta, \Theta)$  (ou  $A_2(\tilde{H}_\theta, \Theta)$ ) satisfaites. Soit  $\mathcal{M}$  une famille finie de modèles contenant  $m^*$  et vérifiant **A4**. Alors, avec  $c_{\min} := \max(\frac{\alpha_1}{2}, \frac{\alpha_2}{2})$ , on a que*

pour tout  $\kappa_n(m) = 2c D_m \log \log n$  avec

$$c \geq c_{\min}, \tag{1.3.12}$$

le modèle sélectionné  $\hat{m}$  par (1.3.2) vérifie

$$\hat{m} \xrightarrow[n \rightarrow \infty]{a.s.} m^*. \tag{1.3.13}$$

Ce résultat nous donne une généralisation de la pénalité fortement consistante de Hannan et Quinn obtenue pour les processus linéaires.

### 1.3.3 Sélection de modèles efficiente et consistante pour les séries temporelles

Les deux premières contributions se sont focalisées sur la propriété de consistance pour la classe des séries affines causales. À présent, nous allons nous intéresser aux propriétés d'efficience notamment asymptotique.

Soit  $(X_1, X_2, \dots, X_n)$  une trajectoire observée du processus (1.3.1). L'objectif n'est plus de retrouver le vrai modèle  $m^*$ .

Nous avons défini en (1.2.3) le risque de prédiction par  $\theta$  de l'observation  $X_t$  au moyen d'un contraste basé sur la log-densité de  $X_t$  sachant le passé. Ainsi, le contraste empirique associé est

$$\gamma_n(\theta) := \mathbb{P}_n \gamma(\theta, \cdot) = \frac{1}{n} \sum_{t=1}^n \gamma(\theta, X_t),$$



avec  $\gamma(\theta, X_t) = q_t(\theta)$  et

$$\mathbb{P}_n = \frac{1}{n} \sum_{t=1}^n \delta_{X_t},$$

où  $\delta_{X_t}$  est la distribution de Dirac de l'observation  $X_t$ . Pour les mêmes raisons évoquées sur la nécessité d'approximer la fonction  $L_n$ , l'on définit

$$\hat{\gamma}_n(\theta) = \mathbb{P}_n \hat{\gamma}(\theta, \cdot) = \frac{1}{n} \sum_{t=1}^n \hat{\gamma}(\theta, X_t),$$

où  $\hat{\gamma}(\theta, X_t) = \hat{q}_t(\theta)$ .

Considérons  $\mathcal{M}$  une famille finie de modèles. L'on aimerait construire des critères de sélection de sorte à obtenir l'efficacité asymptotique (1.2.6). Pour cela, pour tout  $m \in \mathcal{M}$ , l'on définit le minimiseur du risque empirique

$$\hat{\theta}_m = \operatorname{argmin}_{\theta \in \Theta_m} \hat{\gamma}_n(\theta). \quad (1.3.14)$$

Définissons la pénalité  $\text{pen} : m \in \mathcal{M} \mapsto \text{pen}(m) \in \mathbb{R}^+$ . Considérons

$$\hat{m}_{\text{pen}} = \operatorname{argmin}_{m \in \mathcal{M}} \{ \hat{C}_{\text{pen}}(m) \} \quad \text{avec} \quad \hat{C}_{\text{pen}}(m) := \hat{\gamma}_n(\hat{\theta}_m) + \text{pen}(m). \quad (1.3.15)$$

Pour atteindre l'oracle (1.2.5), la pénalité idéale à considérer en (1.3.15) est

$$\text{pen}_{\text{id}}(m) = R(\hat{\theta}_m) - \hat{\gamma}_n(\hat{\theta}_m). \quad (1.3.16)$$

Cependant, cette pénalité comme son nom l'indique n'est pas du tout accessible. Dans un premier temps, nous chercherons à étudier sa distribution asymptotique afin d'y extraire une pénalité "proche" de celle idéale. Pour cela, il est indispensable d'établir un TCL général qui sera également valable pour des modèles ne contenant pas  $m^*$ . Outre les hypothèses **A1-A3** considérées précédemment, nous supposons:

### 1.3.3.a Hypothèses

**A4:** Pour tout  $m \in \mathcal{M}$ , si  $(\bar{\theta}_{m,n})$  est une suite de  $\Theta_m$  vérifiant  $\bar{\theta}_{m,n} \xrightarrow[n \rightarrow +\infty]{a.s.} \theta_m^*$ , alors

$$\limsup_{n \rightarrow \infty} \left\{ \mathbb{E} \left[ \left\| \frac{1}{n} (\partial_{\theta_i \theta_j}^2 L_n(\bar{\theta}_m))_{i,j \in m} \right\|^{-1} \right]^8 \right\} < \infty. \quad (1.3.17)$$

**Remarque 1.2.** Notons que dès que les fonctions  $\theta \rightarrow M_\theta$  et  $\theta \rightarrow f_\theta$  sont  $\mathcal{C}^2(\Theta)$  et si  $\bar{\theta}_{m,n} \xrightarrow[n \rightarrow +\infty]{a.s.} \theta_m^*$  alors

$$\left\| \left( \frac{1}{n} (\partial_{\theta_i \theta_j}^2 L_n(\bar{\theta}_{m,n}))_{i,j \in m} \right)^{-1} \right\|^8 \xrightarrow[n \rightarrow +\infty]{a.s.} \left\| \left( -\frac{1}{2} \partial_{\theta_i \theta_j}^2 \gamma(\theta_m^*)_{i,j \in m} \right)^{-1} \right\|^8.$$

Ainsi, à partir du théorème d'Egorov, on peut trouver un événement  $\tilde{\Omega}$  avec une probabilité suffisamment grande pour que la relation (1.3.17) dans l'hypothèse **A4** tienne si l'on prend l'espérance sur l'événement  $\tilde{\Omega}$ . Pour le cas particulier des processus linéaires, l'hypothèse **A4** est vraie sous une condition légère sur la distribution de  $X$ , voir par exemple Papangelou (1994) et Findley and Wei (2002).

Nous supposons enfin que la vitesse de décroissance de  $(\alpha_j(f_\theta, \Theta))_j$ ,  $(\alpha_j(M_\theta, \Theta))_j$ ,  $(\alpha_j(\partial_\theta f_\theta, \Theta))_j$  et  $(\alpha_j(\partial_\theta M_\theta, \Theta))_j$  doivent être suffisamment rapides pour garantir la forte consistance et la normalité asymptotique de  $\hat{\theta}_m$  :

**A5:** Supposons  $A_i(\Psi_\theta, \Theta)$ ,  $\Psi_\theta = f_\theta, M_\theta$  avec

$$\alpha_j^0(f_\theta, \Theta) + \alpha_j^0(M_\theta, \Theta) + \alpha_j^1(f_\theta, \Theta) + \alpha_j^1(M_\theta, \Theta) = O(j^{-\delta}) \quad \text{où } \delta > 7/2.$$

Notons que l'hypothèse A5 ne permet pas de considérer des processus à mémoire longue, mais les séries temporelles causales habituelles à mémoire courte satisfont à cette hypothèse.

### 1.3.3.b Résultats théoriques

Le résultat suivant montre que le TCL (1.1.12) s'étend même si le modèle candidat  $m$  est mal spécifié.

**Théorème 1.6.** Soit  $(X_1, \dots, X_n)$  une trajectoire observée d'un processus affine causal  $X$  appartenant à  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$  (ou  $\widetilde{\mathcal{AC}}(\tilde{H}_\theta)$ ) où  $\theta^*$  est un paramètre inconnu du compact  $\Theta$  contenu dans  $\Theta(r) \subset \mathbb{R}^d$  (ou  $\tilde{\Theta}(r) \subset \mathbb{R}^d$ ) avec  $r \geq 8$ . Supposons aussi **A1-A5** satisfaites, alors pour tout  $m \in \mathcal{M}$ ,

$$\sqrt{n}((\hat{\theta}_m)_i - (\theta_m^*)_i)_{i \in m} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}\left(0, (F(\theta_m^*))^{-1} G(\theta_m^*) (F(\theta_m^*))^{-1}\right), \quad (1.3.18)$$

Ainsi, nous pouvons asymptotiquement approximer (1.3.16).

**Proposition 1.5.** Sous les hypothèses du Théorème 1.6, il existe une suite bornée  $(v_n^*)_{n \in \mathbb{N}^*}$  telle que pour tout  $m \in \mathcal{M}$ , l'on ait

$$\mathbb{E}[\text{pen}_{id}(m)] \underset{n \rightarrow \infty}{\sim} -\frac{2}{n} \text{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right) + \frac{v_n^*}{n}. \quad (1.3.19)$$

L'on verra par la suite que le terme en trace dans (1.3.19) est connu pour la plupart des modèles classiques.

Nous énonçons un résultat de consistance qui montre que si la pénalité tends en probabilité vers 0, alors le critère  $\hat{C}_{\text{pen}}$  ne sélectionne pas asymptotiquement un modèle mal spécifié.

**Théorème 1.7.** Sous les hypothèses du Théorème 1.6, supposons aussi qu'il existe  $\varepsilon > 0$  tel que,

$$n \mathbb{P}(\text{pen}(m) \geq \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0 \quad \text{pour tout } m \in \mathcal{M}. \quad (1.3.20)$$

alors,

$$n \mathbb{P}(m^* \not\subset \hat{m}_{\text{pen}}) \xrightarrow[n \rightarrow \infty]{} 0. \quad (1.3.21)$$

À présent, nous pouvons spécifier la vitesse de convergence de  $\text{pen}$  à considérer afin d'obtenir un excès de risque proche de celui de l'oracle.

**Théorème 1.8.** Sous les hypothèses du Théorème 1.6 et si pour tout  $\varepsilon > 0$ , il existe  $K_\varepsilon > 0$  tel que

$$\mathbb{P}\left(\limsup_{n \rightarrow \infty} \max_{m \in \mathcal{M}} (n \text{pen}(m)) \geq K_\varepsilon\right) \leq \varepsilon. \quad (1.3.22)$$

Alors, pour tout  $\varepsilon > 0$ , il existe  $M_\varepsilon > 0$  et  $N_\varepsilon \in \mathbb{N}^*$  tels que pour tout  $n \geq N_\varepsilon$ ,

$$\mathbb{P}\left(\ell(\hat{\theta}_{\hat{m}_{\text{pen}}}, \theta^*) \leq \inf_{m \in \mathcal{M}} \{\ell(\hat{\theta}_m, \theta^*)\} + \frac{M_\varepsilon}{n}\right) \geq 1 - \varepsilon. \quad (1.3.23)$$

**Remarque 1.3.** Remarquons que cette optimalité asymptotique est assez différente de l'optimalité classique de l'efficacité asymptotique, où le cardinal de la collection  $\mathcal{M}$  et la dimension des modèles compétitifs dépendent de  $n$ . Cependant, cette dernière est généralement étudiée dans le cadre où le paramètre  $\theta^*$  est de dimension infinie (voir par exemple [Shibata \(1980\)](#), [Li \(1987\)](#), [Hsu et al. \(2019b\)](#)))

### 1.3.3.c Dérivation du BIC pour les modèles affines causaux

Le résultat suivant nous donne une approximation de  $\hat{S}_n(m, X) := \log(\pi_m \mathbb{P}(X | m))$  qui est le terme principal dans la probabilité a posteriori de choisir le modèle  $m$ . En gardant les notations de la sous-section 1.2.4, nous supposons qu'il existe une fonction borélienne positive  $\theta \rightarrow b_m(\theta)$  telle que  $d\mu_m(\theta) = b_m(\theta) d\theta$ .

**Théorème 1.9.** *Sous les hypothèses **A1**, **A2**, **A3**, **A5** et si pour tout  $x \in \mathbb{R}^\infty$ , les fonctions  $\theta \rightarrow M_\theta$  et  $\theta \rightarrow f_\theta$  sont  $\mathcal{C}^6(\Theta)$  et satisfont  $A_k(f_\theta, \Theta)$  et  $A_k(M_\theta, \Theta)$  pour tout  $0 \leq k \leq 2$ . Alors,*

$$\begin{aligned} \hat{S}_n(m, X) &= \hat{L}_n(\hat{\theta}_m) - \frac{\log(n)}{2} |m| + \log(b_m(\hat{\theta}_m)) \\ &\quad + \frac{\log(2\pi)}{2} |m| - \frac{1}{2} \log(\det(-\hat{F}_n(m))) - \log(|\mathcal{M}|) + O(n^{-1}) \quad \text{a.s.} \end{aligned} \quad (1.3.24)$$

avec  $\hat{F}_n(m) := \partial_{\theta^2}^2 \hat{L}_n(\hat{\theta}_m)$  et  $\pi_m = 1/|\mathcal{M}|$ .

Dans l'équation ci-dessus, il est clair que  $-2\hat{S}_n(m, X) \simeq -2\hat{L}_n(\hat{\theta}_m) + \log(n)|m|$  p.s.. Cela donne une légitimité au critère BIC habituel dans le cadre des processus affines causaux puisque :

$$\hat{m}_{BIC} = \operatorname{argmin}_{m \in \mathcal{M}} \left\{ -2\hat{L}_n(\hat{\theta}_m) + \log(n)|m| \right\},$$

et nous observons que  $\hat{m}_{BIC}$  maximise les principaux termes de  $\hat{S}_n(m, X)$ .

À partir de la relation (1.3.24), en considérant certains termes de second ordre dans l'approximation asymptotique de  $\hat{S}_n(m, X)$ , nous pouvons aussi obtenir le critère de Kashyap (voir [Kashyap \(1982\)](#), [Sclove \(1987\)](#), [Bozdogan \(1987\)](#)), définie pour tout  $m \in \mathcal{M}$  par

$$\begin{aligned} \widehat{KC}(m) &:= -2\hat{L}_n(\hat{\theta}_m) + \log(n)|m| + \log(\det(-\hat{F}_n(m))) \\ \text{et } \hat{m}_{KC} &= \operatorname{argmin}_{m \in \mathcal{M}} \{ \widehat{KC}(m) \}. \end{aligned} \quad (1.3.25)$$

Par conséquent, nous pourrions définir un nouveau critère consistant data driven, appelé  $\widehat{KC}'$ , tel que, pour tout  $m \in \mathcal{M}$

$$\begin{aligned} \widehat{KC}'(m) &:= -2\hat{L}_n(\hat{\theta}_m) + (\log(n) - \log(2\pi))|m| + \log(\det(-\hat{F}_n(m))) + 2\log(|m|) \\ \text{et } \hat{m}_{KC'} &= \operatorname{argmin}_{m \in \mathcal{M}} \{ \widehat{KC}'(m) \}. \end{aligned} \quad (1.3.26)$$

**Corollaire 1.1.** *Sous les hypothèses **A1**, **A2**, **A3** et **A5**, nous déduisons de [Bardet et al. \(2020b\)](#), que  $\hat{m}_{BIC}$ ,  $\hat{m}_{KC}$  et  $\hat{m}_{KC'}$  sont consistants, i.e.*

$$\hat{m}_{BIC} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} m^*, \quad \hat{m}_{KC} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} m^* \quad \text{et} \quad \hat{m}_{KC'} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} m^*.$$

### 1.3.3.d Quelques exemples de calcul du terme trace dans la pénalité idéale

Pour un modèle  $m$  vérifiant  $\theta_m^* \in \Theta_m$

1/ Si  $(\xi_t)$  est un bruit blanc gaussien, alors  $\mu_4 = 3$  et d'après [Bardet and Wintenberger \(2009\)](#), l'on a

$$G(\theta_m^*) = -F(\theta_m^*) \implies -2 \text{Trace}((F(\theta_m^*))^{-1} G(\theta_m^*)) = 2 \text{Trace}(I_{|m|}) = 2|m|;$$

2/ Si le paramètre  $\theta$  identifiant un modèle affine causal  $X_t = M_\theta^t \xi_t + f_\theta^t$  peut s'écrire  $\theta = (\theta_1, \theta_2)'$  avec  $f_\theta^t = \tilde{f}_{\theta_1}^t$  et  $M_\theta^t = \tilde{M}_{\theta_2}^t$ . Soit  $p_1, p_2$  tels que  $p_1 = |\theta_1|$ ,  $p_2 = |\theta_2|$  et  $|m| = p_1 + p_2$ . Alors nous obtenons

$$-2 \text{Trace}((F(\theta_m^*))^{-1} G(\theta_m^*)) = 2p_1 + (\mu_4 - 1)p_2. \quad (1.3.27)$$

Au Chapitre 4, l'on trouvera l'exemple de calcul dans un cas complexe (AR-ARCH).

### 1.3.4 Sélection de modèles data driven pour la prédiction d'un processus linéaire autorégressif

A la suite de notre troisième contribution sur l'efficacité asymptotique, une question naturelle qui se pose est celle de la construction de critères de sélection qui atteindront cette propriété d'efficacité mais à  $n$  fixé. Au vu de la complexité de notre modèle général, nous nous sommes restreints au processus  $\text{AR}(\infty)$ .

Supposons avoir observé  $(X_1, X_2, \dots, X_n)$  trajectoire du processus  $(X_t)$  vérifiant

$$X_t = \sum_{k=1}^{+\infty} \theta_k^* X_{t-k} + \sigma \xi_t \quad \text{pour tout } t \in \mathbb{Z}. \quad (1.3.28)$$

L'on aimerait prédire  $X_{n+1}$  par sélection de modèles. Ainsi, soit  $S_m$  (souvent noté  $m$ ) un modèle qui est un ensemble de fonctions linéaires  $f$  de  $\mathbb{R}^{D_m}$  dans  $\mathbb{R}$  telle que

$$f(x_1, x_2, \dots, x_{D_m}) = \sum_{i=1}^{D_m} \theta_i x_i, \quad (1.3.29)$$

avec  $\theta = (\theta_1, \dots, \theta_{D_m}) \in \Theta_m$  et  $\Theta_m$  un compact de  $\mathbb{R}^{D_m}$ .

Etant donné un prédicteur  $f_\theta \in S_m$ , sa qualité de prédiction est mesurée par le risque quadratique

$$R(\theta) = \mathbb{E}[(X_{n+1} - f_\theta^{n+1})^2]$$

où  $f_\theta^n = f_\theta(X_{n-1}, \dots, X_{n-D_m})$ . Le prédicteur de Bayes qui minimise  $R(\theta)$  sur l'ensemble de tous les prédicteurs est clairement la fonction inaccessible  $f_{\theta^*}$ . Introduisons alors l'excès de risque du prédicteur  $f_\theta$

$$\ell(\theta, \theta^*) := R(\theta) - R(\theta^*) = \mathbb{E}[(f_{\theta^*}^{n+1} - f_\theta^{n+1})^2] \geq 0.$$

we will consider that the excess loss is measured on the design points Pour un modèle  $m$ , le meilleur prédicteur est  $f_{\theta_m^*}$  défini par

$$\theta_m^* = \underset{\theta \in \Theta_m}{\operatorname{argmin}} R(\theta).$$

La version empirique de  $\theta_m^*$  est

$$\hat{\theta}_m = \operatorname{argmin}_{\theta \in \Theta_m} \gamma_n(\theta) \quad \text{avec} \quad \gamma_n(\theta) = \frac{1}{n} \sum_{t=1}^n (X_t - f_\theta^t)^2. \quad (1.3.30)$$

Considérons  $\mathcal{M}_n$  une famille dénombrable de modèles  $S_m$ . Le meilleur modèle de la collection  $\mathcal{M}_n$  est déterminé par pénalisation du contraste empirique

$$\hat{m} = \operatorname{argmin}_{m \in \mathcal{M}_n} \{C(m)\} \quad \text{avec} \quad C(m) := \gamma_n(\hat{\theta}_m) + \operatorname{pen}(S_m). \quad (1.3.31)$$

Le but est de trouver la fonction  $\operatorname{pen}$  de sorte que  $\hat{m}$  vérifie

$$\ell(\hat{\theta}_{\hat{m}}, \theta^*) \leq C_1 \inf_{m \in \mathcal{M}_n} \left\{ \ell(\theta_m^*, \theta^*) + \operatorname{pen}(S_m) \right\} + \frac{C_2}{n} \quad (1.3.32)$$

où  $C_1 = 1 + \delta$  avec  $\delta > 0$  (et proche de zéro) et  $C_2 > 0$  et  $\ell(\hat{\theta}, \theta^*) = \mathbb{E} \left[ \frac{1}{n} \sum_{t=1}^n (f_{\hat{\theta}}^t - f_{\theta^*}^t)^2 \right]$ . C'est-à-dire que le modèle sélectionné  $\hat{m}$  sera suffisamment grand pour réduire son biais, mais pas trop grand pour éviter d'avoir une variance élevée.

#### 1.3.4.a Hypothèses

Pour montrer (1.3.32), nous faisons les hypothèses suivantes.

La toute première n'est rien d'autre que la condition assurant l'existence et la stationnarité de la solution de (1.3.28)

$$\mathbf{A1} : \quad \sum_{i=1}^{\infty} |\theta_i^*| < 1.$$

La deuxième hypothèse concerne la sous-gaussianité des observations. Elle nous permettra d'appliquer les inégalités exponentielles pour de tels types de variables.

**A2**  $X_t$  est sous-gaussien avec proxy de variance  $\sigma_0^2 > 0$  i.e.

$$\mathbb{E}[e^{\lambda X_t}] \leq e^{\lambda^2 \sigma_0^2 / 2} \quad \text{pour tout } \lambda > 0.$$

Cette hypothèse est vérifiée si l'on considère un bruit gaussien ou borné.

Pour assurer l'inversibilité des matrices  $\Sigma_m$  et  $\hat{\Sigma}_m$  définies par  $\Sigma_m = \mathbb{E}[\hat{\Sigma}_m]$ ,  $\hat{\Sigma}_m = \mathbf{M}_m^\top \mathbf{M}_m$  avec  $\mathbf{M}_m = [X_{i-1}, \dots, X_{i-D_m}]_{i=1}^n$ , il est suffisant de supposer que

**A3:** Pour tout  $f_\theta \in S_m$ ,  $\langle \alpha, \partial_\theta f_\theta \rangle = 0$  p.s.  $\implies \alpha = 0$ .

Nous avons eu besoin de supposer une condition assez faible sur la densité spectrale du processus  $(X_t)_{t \in \mathbb{Z}}$ . Cette densité  $g : [-\pi, \pi] \rightarrow \mathbb{C}$  est définie ainsi qu'il suit

$$g(\lambda) = \frac{1}{2\pi} \sum_{h \in \mathbb{Z}} r(h) e^{-ih\lambda},$$

où  $r(h) := \mathbb{E}[X_t X_{t+h}]$ , est la fonction de covariance. La fonction  $g$  existe en vertu de **A1** avec  $|\theta_t^*| = O(t^{-\gamma})$  où  $\gamma \geq 1$ .

**A4:** Il existe une constante  $a > 0$  telle que  $\inf_{-\pi \leq \lambda < \pi} g(\lambda) \geq a$ .

Pour des raisons techniques, nous supposons que la dimension du plus grand modèle  $K_n$  de  $\mathcal{M}_n$  est de la forme

### 1.3.4.b Résultats théoriques

La proposition suivante, qui est une conséquence du Théorème 3.1 de [Doukhan and Wintenberger \(2008\)](#), établit un lien entre les coefficients de mélange  $\tau$  du processus  $(X_t)_{t \in \mathbb{Z}}$  et les coefficients  $\theta_i^*$  du modèle (1.3.28).

**Proposition 1.6.** *Supposons que **A1** tienne avec  $|\theta_t^*| = O(t^{-\gamma})$  et  $\gamma > 1$ , alors il existe une solution  $\tau$ -dependante stationnaire de (1.3.28) et une constante  $C_\tau > 0$  telle que pour tout  $r > 0$*

$$\tau_{X,\infty}^{(2)}(r) \leq C_\tau \left( \frac{\log r}{r} \right)^{\gamma-1} \quad (1.3.33)$$

Le résultat intermédiaire suivant nous donne une borne inférieure et supérieure uniformes pour la norme spectrale des matrices  $\Sigma_m$ .

**Proposition 1.7.** *Supposons que **A1** tienne avec  $|\theta_t^*| = O(t^{-\gamma})$  et  $\gamma \geq 2$ , alors pour tout  $m \in \mathcal{M}_n$*

$$\|\Sigma_m\|_{\text{op}} \leq \pi^{-1} \sum_{i=0}^{\infty} |\mathbb{E}[X_0 X_i]| < \infty. \quad (1.3.34)$$

Si **A3-A4** sont satisfaites par ailleurs, alors

$$\|\Sigma_m^{-1}\|_{\text{op}} \leq 1/a. \quad (1.3.35)$$

Pour prouver (1.3.32) en présence de la dépendance, il est courant de considérer l'ensemble

$$\Omega_n = \left\{ \omega : \left| \frac{\|F_\theta\|_n^2}{\|F_\theta\|_\mu^2} - 1 \right| \leq \frac{1}{2}, \quad \forall F_\theta \in \bigcup_{m,m' \in \mathcal{M}_n} (S_m + S_{m'}) \right\}$$

où  $\|F_\theta\|_\mu^2 := \frac{1}{n} \mathbb{E} \left[ \sum_{t=1}^n (f_\theta^t)^2 \right] = \int (f_\theta^1)^2 d\mu$  et  $\|F_\theta\|_n^2 = n^{-1} \sum_{t=1}^n (f_\theta^t)^2$ ,  $F_\theta = (f_\theta^1, \dots, f_\theta^n)^\top$ .

Sur  $\Omega_n$ , on a une relation entre la norme empirique  $\|\cdot\|_n$  et la norme  $\mathbb{L}_2$  (voir par exemple [Baraud et al. \(2001b\)](#), [Comte and Genon-Catalot \(2020\)](#)) et cela est très pratique pour prouver (1.3.32).

La proposition suivante montre que  $\Omega_n$  tient avec une grande probabilité.

**Proposition 1.8.** *Sous les hypothèses **A1** – **A4** et si  $|\theta_t^*| = O(t^{-\gamma})$  avec  $\gamma \geq 8$ , nous avons*

$$\mathbb{P}(\Omega_n^c) \leq \frac{c_0}{n^3}, \quad (1.3.36)$$

où  $c_0 > 0$ .

Nous sommes maintenant en mesure d'énoncer le résultat principal de notre quatrième contribution.

**Théorème 1.10.** *Considérons les observations  $(X_1, \dots, X_n)$  d'une solution du processus (1.3.28) satisfaisant **A1** avec  $|\theta_t^*| = O(t^{-\gamma})$  où  $\gamma \geq 8$  et vérifiant aussi **A2** et **A4**. Soit  $\mathcal{M}_n$  une famille dénombrable de modèles  $S_m$  satisfaisant **A3**. Pour  $x > 0$ , considérons la fonction pen:  $\mathcal{M}_n \rightarrow \mathbb{R}^+$  telle que*

$$\text{pen}(S_m) \geq 8x\sigma^2 \frac{D_m}{n}. \quad (1.3.37)$$

Alors, l'estimateur  $\hat{\theta}_{\hat{m}}$  avec  $\hat{m}$  donnée en (1.3.31), satisfait

$$\mathbb{E} \left[ \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta^*}\|_n^2 \mathbb{I}_{\Omega_n} \right] \leq C_1(x) \inf_{m \in \mathcal{M}_n} \left\{ \mathbb{E} \left[ \|F_{\theta^*} - F_{\theta_m^*}\|_n^2 \right] + 2 \text{pen}(S_m) \right\} + \frac{x(x+2)}{x-2} \frac{C_2}{n}$$

avec  $C_1(x) = \left( \frac{x+2}{x-2} \right)^2 > 1$  et  $C_2 > 0$ .

Ce théorème vient clore la liste des résultats principaux obtenus tout au long de cette thèse. Dans la suite, nous allons présenter de manière détaillée les quatre chapitres associés aux contributions résumées dans ce chapitre.

# 2

## Consistent model selection criteria and goodness-of-fit test for common time series models

---

2.1	Introduction	30
2.2	General Framework	33
2.2.1	The Gaussian quasi-maximum likelihood estimation and the model selection criterion	33
2.2.2	The affine causal framework	34
2.2.3	The special case of NLARCH( $\infty$ ) processes	36
2.3	Asymptotic results	36
2.3.1	Assumptions required for the asymptotic study	36
2.3.2	Asymptotic model selection	38
2.4	Examples	38
2.4.1	AR( $\infty$ ) models	39
2.4.2	ARCH( $\infty$ ) models	39
2.4.3	APARCH( $\delta, p, q$ ) models	40
2.4.4	ARMA( $p, q$ )-GARCH( $p', q'$ ) models	41
2.5	Portmanteau test	41
2.6	Numerical Results	43
2.6.1	Monte-Carlo experiments for common time series selection	43
2.6.2	Subset model selection	45
2.6.3	Application to real data	46
2.6.3.a	Air quality analysis	46
2.6.3.b	Financial index analysis	47
2.7	Proofs	49
2.7.1	Misspecified model	55
2.7.2	Proof of Theorem 2.1	55
2.7.3	Proof of Theorem 2.2	58
2.7.4	Proof of Theorem 2.3	58

---



The content of this chapter was taken from the article: J-M BARDET, K. KAMILA, W. KENGNE, "Consistent Model Selection Criteria and Goodness-of-fit Test for Affine Causal Processes", Electronic journal of statistics, Shaker Heights, OH : Institute of Mathematical Statistics, 2020, 14 (1), pp.2009-2052.

## Abstract

This chapter studies the model selection problem in a large class of causal time series models, which includes both the ARMA or  $AR(\infty)$  processes, as well as the GARCH or  $ARCH(\infty)$ , APARCH, ARMA-GARCH and many others processes. To tackle this issue, we consider a penalized contrast based on the quasi-likelihood of the model. We provide sufficient conditions for the penalty term to ensure the consistency of the proposed procedure as well as the consistency and the asymptotic normality of the quasi-maximum likelihood estimator of the chosen model. We also propose a tool for diagnosing the goodness-of-fit of the chosen model based on a Portmanteau test. Monte-Carlo experiments and numerical applications on illustrative examples are performed to highlight the obtained asymptotic results.

## 2.1 Introduction

Model selection is an important tool for statisticians and all those who process data. This issue has received considerable attention in the recent literature. There are several model selection procedures, the main ones are : cross validation and penalized contrast based.

The cross validation (Stone (1974), Allen (1974)) consists in splitting the data into learning sample, which will be used for computing estimators of the parameters and the test sample which allows to assess these estimators by evaluate their risks.

The procedures using penalized objective function search for a model, minimizing a trade-off between a sum of an empirical risk (for instance least squares,  $-2 \times \log$ -likelihood), which indicates how well the model fits the data, and a measure of model's complexity so-called a penalty.

The idea of penalizing dates back to the 1970s with the works of Mallows (1973) and Akaike (1973). Although it is likely that these ideas have already existed in other contexts such as subset selection by Beale et al. (1967), Hocking and Leslie (1967), and ridge regression by Hoerl and Kennard (1970).

By using the ordinary least squares in regression framework, Mallows obtained the  $C_p$  criterion. Meanwhile, Akaike derived AIC for density estimation using log-likelihood contrast. A few years later, following Akaike, Schwarz (1978) proposed an alternative approach to density estimation and derived the Bayesian Information Criteria (BIC). The penalty term of these criteria is proportional to the dimension of the model. In the recent decades, different approaches of penalization have emerged such as the  $\mathbb{L}^2$  norm for the Ridge penalisation Hoerl and Kennard (1970), the  $\mathbb{L}^1$  norm used by Tibshirani (1996) that provides the LASSO procedure and the elastic-net that mixes the  $\mathbb{L}^1$  and  $\mathbb{L}^2$  norms Zou and Hastie (2005).

Model selection procedures can have two different objectives: *consistency* and *efficiency*. A procedure is said to be consistent if given a family of models, including the "true model", the probability of choosing the correct model approaches one as the sample size tends to infinity. On the other hand, a procedure is efficient when its risk is asymptotically equivalent to the risk of the oracle. In this work, we are interested to construct a consistent procedure for the general class of times series known as *affine causal processes*, which includes the most common time series.

This class of affine causal time series can be defined as follows. Let  $\mathbb{R}^\infty$  be the space of sequences of real numbers with a finite number of non zero, if  $M, f : \mathbb{R}^\infty \rightarrow \mathbb{R}$  are two measurable functions, then an affine causal class is

**Class  $\mathcal{AC}(M, f)$  :** A process  $X = (X_t)_{t \in \mathbb{Z}}$  belongs to  $\mathcal{AC}(M, f)$  if it satisfies:

$$X_t = M((X_{t-i})_{i \in \mathbb{N}^*}) \xi_t + f((X_{t-i})_{i \in \mathbb{N}^*}) \quad \text{for any } t \in \mathbb{Z}; \quad (2.1.1)$$

where  $(\xi_t)_{t \in \mathbb{Z}}$  is a sequence of zero-mean independent identically distributed random vectors (i.i.d.r.v) satisfying  $\mathbb{E}(|\xi_0|^r) < \infty$  for some  $r \geq 2$  and  $\mathbb{E}[\xi_0^2] = 1$ .

For instance,

- if  $M((X_{t-i})_{i \in \mathbb{N}^*}) = \sigma$  and  $f((X_{t-i})_{i \in \mathbb{N}^*}) = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p}$ , then  $(X_t)_{t \in \mathbb{Z}}$  is an AR( $p$ ) process;
- if  $M((X_{t-i})_{i \in \mathbb{N}^*}) = \sqrt{a_0 + a_1 X_{t-1}^2 + \dots + a_p X_{t-p}^2}$  and  $f((X_{t-i})_{i \in \mathbb{N}^*}) = 0$ , then  $(X_t)_{t \in \mathbb{Z}}$  is an ARCH( $p$ ) process.

Numerous classical time series models such as ARMA( $p, q$ ), GARCH( $p, q$ ), ARMA( $p, q$ )-GARCH( $p, q$ ) (see [Ding et al. \(1993\)](#) and [Ling and McAleer \(2003\)](#)) or APARCH( $\delta, p, q$ ) processes (see [Ding et al. \(1993\)](#)) belongs to  $\mathcal{AC}(M, f)$ . The existence of stationary and ergodic solutions of this class has been studied in [Doukhan and Wintenberger \(2008\)](#) and [Bardet and Wintenberger \(2009\)](#).

We consider a trajectory  $(X_1, \dots, X_n)$  of a stationary affine causal process  $\mathcal{AC}(M^*, f^*)$ , where  $M^*$  and  $f^*$  are unknown. We also consider a finite set  $\mathcal{M}$  of parametric models  $m$ , which are affine causal time series. We assume that the "true" model  $m^*$  corresponds to  $M^*$  and  $f^*$ . The aim is to obtain an estimator  $\hat{m}$  of  $m^*$  and testing the goodness-of-fit of the chosen model.

There already exist several important contributions devoted to the model selection for time series ; we refer to the book of [McQuarrie and Tsai \(1998\)](#) and the references therein for an overview on this topic.

As we have pointed above, two properties are often used to evaluate a quality of a model selection procedure : consistency and efficiency. The consistency assumes that the true model exists and it is included in the collection of candidate models; while the efficiency does not necessarily require the existence of a true model. In many research in this framework, the main goal is to develop a procedure that fulfills one of these properties. So, in some classical linear time series models, the consistency of the BIC procedure has been established, see for instance [Hannan \(1980\)](#) or [Tsay \(1984\)](#) ; and the asymptotic efficiency of the AIC has been proved, see, among others, [Shibata \(1980\)](#), [Hurvich and Tsai \(1989\)](#) for a corrected version of AIC for small samples, [Ing and Wei \(2005\)](#), [Ing \(2007\)](#), [Ing et al. \(2012\)](#) for the case of infinite order autoregressive model. [Shi and Tsai \(2002\)](#) propose the (consistent) residual information criteria (RIC) for regression model (including regression models with ARMA errors) selection. In the framework of nonlinear threshold models, [Kapetanios \(2001\)](#) proved consistency results of a large class of information criteria, whereas [Gao and Tong \(2004\)](#) focused on cross-validation type procedure for model selection in a class of semiparametric time series regression model. Let us recall that, the time series model selection literature is very extensive and still growing ; we refer to the monograph of [Rao et al. \(2001\)](#), which provided an excellent summary of existing model selection procedure, including the case of time series models as well as the recent review paper of [Ding et al. \(2018\)](#).

The adaptive lasso, introduced by Zou (2006) for variable selection in linear regression models has been extended by Ren and Zhang (2010) to vector autoregressive models, Kock (2016) carried out this procedure in stationary and nonstationary autoregressive models; the oracle efficiency is established. Lerasle (2011) considers model selection for density estimation under mixing conditions and derived oracle inequalities of the slope heuristic procedure (Birgé and Massart (2007a) or Arlot and Massart (2009)) ; whereas Alquier and Wintenberger (2012) develop oracle inequalities for model selection for weakly dependent time series forecasting. Recently, Shao and Yang (2017) have considered the model selection for ARMA time series with trend, and proved the consistency of BIC for the detrended residual sequence, while Arkoun et al. (2019) developed oracle inequalities of sequential model selection method for nonparametric autoregression. Hsu et al. (2019a) pointed out that most existing model selection procedure cannot simultaneously enjoy consistency and (asymptotic) efficiency. They propose a misspecification-resistant information criterion that can achieve consistency and asymptotic efficiency for prediction using model selection.

In this paper, we focus on the class of models (5.1.1), and addressed the following questions :

1. What regularity conditions are sufficient to build a consistent model selection procedure? Does the classic criterion such as BIC, still have consistent property for choosing a model among the collection  $\mathcal{M}$ ?
2. How can we test the goodness-of-fit of the chosen model?

These questions have not yet been answered for the class of models and the framework considered here, in particular in case of infinite memory processes. This new contribution provides theoretical and numerical response of these issues. (i) The estimator  $\hat{m}$  of  $m^*$  is chosen by minimizing a penalized criterion  $\hat{C}(m) = -2\hat{L}_n(m) + |m|\kappa_n$ , where  $\hat{L}_n(m)$  is a Gaussian quasi-log-likelihood of the model  $m$ ,  $|m|$  is the number of estimated parameters of the model  $m$  and  $\kappa_n$  is a non-decreasing sequence of real numbers (see more details in Section 2.2). Note that, in the cases  $\kappa_n = 2$  or  $\kappa_n = \log n$  we respectively consider the usual AIC and BIC criteria. We provide sufficient conditions (essentially depending on the decreasing of the Lipschitz coefficients of the functions  $f$  and  $M$ ) for obtaining consistency of the model selection procedure.

(ii) We provide an asymptotic goodness-of-fit test for the selected model that is very simple to be used (with the usual Chi-square distribution limit), which successively completes the model selection procedure. Numerical applications show the accuracy of this test under the null hypothesis as well as an efficient test power under an alternative hypothesis. Note that, a similar test has been proposed by Li and Mak (1994) under the Gaussian assumption on the observations, whereas Ling and Li (1997) focused for multivariate time series with multivariate ARCH-type errors. These papers are also based on exact likelihood estimators that do not make feasible Portemanteau tests. Duchesne and Francq (2008) proposed an interesting Portmanteau test statistic directly based on the autocorrelations of residuals (and not squared residuals) computed from quasi-likelihood estimators for diagnostic checking in the class of model (5.1.1). Unlike these authors, we apply the test to a model obtained from a model selection procedure.

Monte-Carlo experiments and numerical applications on illustrative examples are also performed to highlight the obtained asymptotic results. We have also considered a data-driven choice of the penalty obtained from the slope heuristic procedure (see for instance Arlot and Massart (2009)) for avoiding an a priori choice of the penalty sequence. The

simulation study and real data applications show that the results of the proposed model selection procedure and the Portetmanteau test are overall satisfactory.

The paper is organized as follows. Some definitions, notations and assumptions are described in Section 2.2. The consistency of the criteria and the asymptotic normality of the post-model-selection estimator are studied in Section 2.3. In Section 2.4, the examples of  $AR(\infty)$ ,  $ARCH(\infty)$ ,  $APARCH(\delta, p, q)$  and  $ARMA(p, q)$ -GARCH( $p', q'$ ) processes are detailed. The goodness-of-fit test is studied in Section 2.5. Finally, numerical results are presented in Section 2.6 and Section 2.7 contains the proofs.

## 2.2 General Framework

In this section, we are going to present the model selection using Gaussian quasi-maximum likelihood estimators (QMLE) and give some notations in order to facilitate the presentation.

### 2.2.1 The Gaussian quasi-maximum likelihood estimation and the model selection criterion

In the sequel, for a model  $m \in \mathcal{M}$ , a family of models of  $\mathcal{AC}(M_\theta, f_\theta)$  with  $\theta \in \Theta \subset \mathbb{R}^d$ , where  $\theta \rightarrow M_\theta$  and  $\theta \rightarrow f_\theta$  are two fixed functions, we are going to consider QMLE of  $\theta$  for each specific model  $m$ .

This approach as semi-parametric estimation has been successively introduced for GARCH( $p, q$ ) processes in [Jeaneau \(1998\)](#) where its consistency is also proved, and the asymptotic normality of this estimator has been established in [Berkas et al. \(2003\)](#) and [Francq and Zakoian \(2004\)](#). In [Bardet and Wintenberger \(2009\)](#), those results have been extended to affine causal processes, and an extension to Laplacian QMLE has been also proposed in [Bardet et al. \(2017\)](#).

The Gaussian QMLE is derived from the conditional (with respect to the filtration  $\sigma\{(X_t)_{t \leq 0}\}$ ) log-likelihood of  $(X_1, \dots, X_n)$  when  $(\xi_t)$  is supposed to be a Gaussian standard white noise. Due to the linearity of a causal affine process, we deduce that this conditional log-likelihood (up to an additional constant)  $L_n$  is defined for all  $\theta \in \Theta$  by:

$$L_n(\theta) := -\frac{1}{2} \sum_{t=1}^n q_t(\theta), \text{ with } q_t(\theta) := \frac{(X_t - f_\theta^t)^2}{H_\theta^t} + \log(H_\theta^t) \quad (2.2.1)$$

where  $f_\theta^t := f_\theta(X_{t-1}, X_{t-2}, \dots)$ ,  $M_\theta^t := M_\theta(X_{t-1}, X_{t-2}, \dots)$  and  $H_\theta^t = (M_\theta^t)^2$ . Since  $L_n(\theta)$  depends on  $(X_t)_{t \leq 0}$  that are unobserved, the idea of the quasi log-likelihood is to replace  $q_t(\theta)$  by an approximation  $\hat{q}_t(\theta)$  and to compute  $\hat{\theta}$  as in equation (2.2.3) even if the white noise is not Gaussian. Hence, the conditional Gaussian quasi log-likelihood (up to an additional constant) is given for all  $\theta \in \Theta$  by

$$\begin{aligned} \hat{L}_n(\theta) &:= -\frac{1}{2} \sum_{t=1}^n \hat{q}_t(\theta), \text{ with } \hat{q}_t(\theta) := \frac{(X_t - \hat{f}_\theta^t)^2}{\hat{H}_\theta^t} + \log(\hat{H}_\theta^t) \\ \text{where } \begin{cases} \hat{f}_\theta^t &:= f_\theta(X_{t-1}, X_{t-2}, \dots, X_1, u) \\ \hat{M}_\theta^t &:= M_\theta(X_{t-1}, X_{t-2}, \dots, X_1, u) \\ \hat{H}_\theta^t &:= (\hat{M}_\theta^t)^2 \end{cases} \quad (2.2.2) \end{aligned}$$

for any deterministic sequence  $u = (u_n)$  with finitely many non-zero values ( $u = 0$  is very often chosen without loss of generality).

Finally, for each specific model  $m \in \mathcal{M}$ , we define the Gaussian QMLE  $\hat{\theta}(m)$  as

$$\hat{\theta}(m) = \operatorname{argmax}_{\theta \in \Theta(m)} \hat{L}_n(\theta). \quad (2.2.3)$$

To select the "best" model  $m \in \mathcal{M}$ , we chose a penalized contrast  $\hat{C}(m)$  ensuring a trade-off between  $-2$  times the maximized quasi log-likelihood, which decreases with the size of the model, and a penalty increasing with the size of the model. Therefore, the choice of the "best" model  $\hat{m}$  among the estimated can be performed by minimizing the following criteria

$$\hat{m} = \operatorname{argmin}_{m \in \mathcal{M}} \hat{C}(m) \quad \text{with} \quad \hat{C}(m) = -2\hat{L}_n(\hat{\theta}(m)) + |m| \kappa_n, \quad (2.2.4)$$

where

- $(\kappa_n)_n$  an increasing sequence depending on the number of observations  $n$ .
- $|m|$  denotes the dimension of the model  $m$ , *i.e.* the cardinal of  $m$ , subset of  $\{1, \dots, d\}$ , which is also the number of estimated components of  $\theta$  (the others are fixed to zero).

The consistency of the criterion  $\hat{C}$ , *i.e.*

$$\mathbb{P}(\hat{m} = m^*) \xrightarrow[n \rightarrow \infty]{} 1; \quad (2.2.5)$$

will be established after showing that both of following probabilities are zero:

- the asymptotic probability of selecting a larger model containing the true model (overfitting case);
- the asymptotic probability of selecting a false model that is a model not containing  $m^*$ .

## 2.2.2 The affine causal framework

In the introduction, to be more concise, we have presented the problem of time series model selection in a very general form. In reality, we will limit our field of study a little bit by considering a semi-parametric framework. Hence, let  $(f_\theta)_{\theta \in \Theta}$  and  $(M_\theta)_{\theta \in \Theta}$  be two families of known functions such as for any  $\theta \in \Theta$ , both  $f_\theta, M_\theta$  with real values defined on  $\mathbb{R}^\infty$ .

Before diving in details, let's give some notations that will be useful throughout the paper. We will consider a subset  $\Theta$  of  $\mathbb{R}^d$  ( $d \in \mathbb{N}$ ). We will use the following norms:

- $\|\cdot\|$  denotes the usual Euclidean norm on  $\mathbb{R}^\nu$ , with  $\nu \geq 1$ ;
- if  $X$  is  $\mathbb{R}^\nu$ -random variable with  $r \geq 1$  order moment, we set  $\|X\|_r = (\mathbb{E}(\|X\|^r))^{1/r}$ ;
- for any set  $\Theta \subseteq \mathbb{R}^d$  and for any  $g : \Theta \rightarrow \mathbb{R}^{d'}$ ,  $d' \geq 1$ , denote  $\|g\|_\Theta = \sup_{\theta \in \Theta} \{\|g(\theta)\|\}$ .

Let us start with an example to better understand the framework and the approach of model selection we will follow.

**Example:** Assume that the observed trajectory  $(X_1, \dots, X_n)$  is generated from a model

belonging to a collection  $\mathcal{M}$ , for instance a set of  $\text{ARMA}(p, q)$  and  $\text{GARCH}(p', q')$  processes for  $0 \leq p \leq p_{\max}$ ,  $0 \leq q \leq q_{\max}$ ,  $0 \leq p' \leq p'_{\max}$ ,  $0 \leq q' \leq q'_{\max}$  (where  $p_{\max}, q_{\max}, p'_{\max}, q'_{\max}$  are the upper bounds of orders). Then, we would like to chose in this family a "best" model for fitting the data  $(X_1, \dots, X_n)$ . For instance, if  $p_{\max} = q_{\max} = p'_{\max} = q'_{\max} = 9$ , in the collection above, there is 200 possible models and we expect to recognize the true process (which is unknown to the analyst) as the selected model, at least when  $n$  is large enough.

We begin with the following property that allow to enlarge the family of models by extending the dimension  $d$  of the parameter  $\theta$ :

**Proposition 2.1.** *Let  $d_1, d_2 \in \mathbb{N}$ ,  $\Theta_1 \subset \mathbb{R}^{d_1}$  and  $\Theta_2 \subset \mathbb{R}^{d_2}$ , and for  $i = 1, 2$ , define  $f_{\theta_i}^{(i)}, M_{\theta_i}^{(i)} : \mathbb{R}^\infty \rightarrow \mathbb{R}$  and for  $\theta_i \in \Theta_i$ . Then there exist  $\max(d_1, d_2) \leq d \leq d_1 + d_2$ ,  $\Theta \subset \mathbb{R}^d$ , and a family of functions  $f_\theta : \mathbb{R}^\infty \rightarrow \mathbb{R}$  and  $M_\theta : \mathbb{R}^\infty \rightarrow [0, \infty)$  with  $\theta \in \Theta$ , such that for any  $\theta_1 \in \Theta_1$  and  $\theta_2 \in \Theta_2$ , there exists  $\theta \in \Theta$  satisfying*

$$\mathcal{AC}(M_{\theta_1}^{(1)}, f_{\theta_1}^{(1)}) \cup \mathcal{AC}(M_{\theta_2}^{(2)}, f_{\theta_2}^{(2)}) \subset \mathcal{AC}(M_\theta, f_\theta).$$

The proof of this proposition, as well as the other proofs, can be found in Section 2.7. This proposition says that it is always possible to embed two parametric causal affine models in a larger one. Hence, for instance, we can consider as well AR processes and ARCH processes in a unique representation, *i.e.*

$$\begin{cases} \text{AR} & \begin{cases} M_{\theta_1}^{(1)}((X_{t-i})_{i \in \mathbb{N}^*}) = \sigma \\ f_{\theta_1}^{(1)}((X_{t-i})_{i \in \mathbb{N}^*}) = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p} \end{cases} \\ \text{ARCH} & \begin{cases} M_{\theta_2}^{(2)}((X_{t-i})_{i \in \mathbb{N}^*}) = \sqrt{a_0 + a_1 X_{t-1}^2 + \dots + a_q X_{t-q}^2} \\ f_{\theta_2}^{(2)}((X_{t-i})_{i \in \mathbb{N}^*}) = 0 \end{cases} \end{cases} \implies \begin{cases} M_\theta((X_{t-i})_{i \in \mathbb{N}^*}) = \sqrt{\theta_0 + \theta_1 X_{t-1}^2 + \dots + \theta_q X_{t-q}^2} \\ f_\theta((X_{t-i})_{i \in \mathbb{N}^*}) = \theta_{q+1} X_{t-1} + \dots + \theta_{q+p} X_{t-p} \end{cases}.$$

From now and in all the sequel, we fix  $d \in \mathbb{N}^*$ , and the family of functions  $f_\theta, M_\theta : \mathbb{R}^\infty \rightarrow \mathbb{R}$  for  $\theta \in \Theta \subset \Theta(r) \subset \mathbb{R}^d$ .

Let  $(X_1, \dots, X_n)$  be an observed trajectory of an affine causal process  $X$  belonging to  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$ , where  $\theta^*$  is an unknown vector of  $\Theta$ , and therefore:

$$X_t = M_{\theta^*}((X_{t-i})_{i \in \mathbb{N}^*}) \xi_t + f_{\theta^*}((X_{t-i})_{i \in \mathbb{N}^*}) \quad \text{for any } t \in \mathbb{Z}. \quad (2.2.6)$$

In the sequel, we will consider several models, which all are particular cases of  $\mathcal{AC}(M_\theta, f_\theta)$  with  $\theta \in \Theta \subset \mathbb{R}^d$ . More precisely define:

- a model  $m$  as a subset of  $\{1, \dots, d\}$  and denote  $|m| = \#(m)$ ;
- $\Theta(m) = \{(\theta_i)_{1 \leq i \leq d} \in \mathbb{R}^d, \theta_i = 0 \text{ if } i \notin m\} \cap \Theta$ ;
- $\mathcal{M}$  as a finite family of models, *i.e.*  $\mathcal{M} \subset \mathcal{P}(\{1, \dots, d\})$ .

Finally, for all  $m \in \mathcal{M}$ ,  $m \in \mathcal{AC}(M_\theta, f_\theta)$  when  $\theta \in \Theta(m)$  and denote  $m^*$  the "true" model. We could as well consider hierarchical or exhaustive families of models.

**Example:** From the previous example, we can consider:



- a family  $\mathcal{M}_1$  of models  $m_1$  such as  $\mathcal{M}_1 = \{\{1\}, \{1, 2\}, \dots, \{1, \dots, q+1\}\}$ : this family is the hierarchical one of ARCH processes with orders varying from 0 to  $q$ .
- a family  $\mathcal{M}_2$  of models  $m_2$  such as  $\mathcal{M}_2 = \mathcal{P}(\{1, \dots, p+q+1\})$ : this family is the exhaustive one and contains as well the AR(2) process  $X_t = \phi_2 X_{t-2} + \theta_0 \xi_t$  as the process  $X_t = \phi_1 X_{t-1} + \phi_3 X_{t-3} + \xi_t \sqrt{\theta_0 + a_2 X_{t-2}^2}$ .

To establish the consistency of the selected model, we will need to assume that the "true" model  $m^*$  with the parameter  $\theta^*$ , is included in the model family  $\mathcal{M}$ .

### 2.2.3 The special case of NLARCH( $\infty$ ) processes

As in [Bardet and Wintenberger \(2009\)](#), in the special case of NLARCH( $\infty$ ) processes, including for instance GARCH( $p, q$ ) or ARCH( $\infty$ ) processes, a particular treatment can be realized for obtaining sharper results than using the previous framework. In such case, define the class:

**Class  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$ :** A process  $X = (X_t)_{t \in \mathbb{Z}}$  belongs to  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$  if it satisfies:

$$X_t = \xi_t \sqrt{\widetilde{H}_\theta((X_{t-i}^2)_{i \in \mathbb{N}^*})} \text{ for any } t \in \mathbb{Z}. \quad (2.2.7)$$

Therefore, if  $M_\theta^2((X_{t-i})_{i \in \mathbb{N}^*}) = H_\theta((X_{t-i})_{i \in \mathbb{N}^*}) = \widetilde{H}_\theta((X_{t-i}^2)_{i \in \mathbb{N}^*})$  then,  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta) = \mathcal{AC}(M_\theta, 0)$ . In case of the class  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$ , we will use the assumption  $A(\widetilde{H}_\theta, \Theta)$ . By this way, we will obtain a new set of stationary solutions. For  $r \geq 2$  define:

$$\widetilde{\Theta}(r) = \left\{ \theta \in \mathbb{R}^d, A(\widetilde{H}_\theta, \{\theta\}) \text{ holds with } (\|\xi_0\|_r)^2 \sum_{k=1}^{\infty} \alpha_k(\widetilde{H}_\theta, \{\theta\}) < 1 \right\}. \quad (2.2.8)$$

Then, for  $\theta \in \Theta(r)$ , a process  $(X_t)_{t \in \mathbb{Z}}$  belonging to the class  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$  is stationary ergodic and satisfies  $\|X_0\|_r < \infty$ .

## 2.3 Asymptotic results

### 2.3.1 Assumptions required for the asymptotic study

We begin by giving a condition on  $f_\theta$  and  $M_\theta$  which ensure the existence of a  $r$ -order moment, stationary and ergodic time series belonging to  $\mathcal{AC}(M_\theta, f_\theta)$ . This condition, initially obtained in [Doukhan and Wintenberger \(2008\)](#), is written in terms of Lipschitz coefficients of both these functions. Hence, for  $\Psi_\theta = f_\theta$  or  $M_\theta$ , define:

**Assumption A( $\Psi_\theta, \Theta$ ):** Assume that  $\|\Psi_\theta(0)\|_\Theta < \infty$  and there exists a sequence of non-negative real numbers  $(\alpha_k(\Psi_\theta, \Theta))_{k \geq 1}$  such that  $\sum_{k=1}^{\infty} \alpha_k(\Psi_\theta, \Theta) < \infty$  satisfying:

$$\|\Psi_\theta(x) - \Psi_\theta(y)\|_\Theta \leq \sum_{k=1}^{\infty} \alpha_k(\Psi_\theta, \Theta) |x_k - y_k| \text{ for all } x, y \in \mathbb{R}^\infty.$$

Now for  $r \geq 1$ , where  $\|\xi_0\|_r < \infty$ , define:

$$\Theta(r) = \left\{ \theta \in \mathbb{R}^d, A(f_\theta, \{\theta\}) \text{ and } A(M_\theta, \{\theta\}) \text{ hold with } \sum_{k=1}^{\infty} \alpha_k(f_\theta, \{\theta\}) + \|\xi_0\|_r \sum_{k=1}^{\infty} \alpha_k(M_\theta, \{\theta\}) < 1 \right\}. \quad (2.3.1)$$

Then, for any  $\theta \in \Theta(r)$ , there exists a stationary and ergodic solution with  $r$ -order moment belonging to  $\mathcal{AC}(M_\theta, f_\theta)$ . (see Doukhan and Wintenberger (2008) and Bardet and Wintenberger (2009)).

Secondly, note that the definitions of the conditional log-likelihood (2.2.1) and quasi log-likelihood (2.2.2) require that their denominators do not vanish. Hence, we will suppose in the sequel that the lower bound of  $H_\theta(\cdot) = (M_\theta(\cdot))^2$  (which is reached since  $\Theta$  is compact) is strictly positive:

**Assumption D( $\Theta$ ):**  $\exists \underline{h} > 0$  such that  $\inf_{\theta \in \Theta} (H_\theta(x)) \geq \underline{h}$  for all  $x \in \mathbb{R}^\infty$ .

The following classical assumption ensures the identifiability of the considered model.

**Assumption Id( $\Theta$ ):** For all  $\theta, \theta' \in \Theta$ ,  $(f_\theta^0 = f_{\theta'}^0 \text{ and } M_\theta^0 = M_{\theta'}^0) \text{ a.s.} \implies \theta = \theta'$ .

Another required assumption concerns the differentiability of  $\Psi_\theta = f_\theta$  or  $M_\theta$  on  $\Theta$ . This type of assumption has already been considered in order to apply the QMLE procedure (see Bardet and Wintenberger (2009), Straumann and Mikosch (2006), White (1982)). First, the following Assumption Var( $\Theta$ ) provides the invertibility of the "Fisher's information matrix" of  $X$  and is important to prove the asymptotic normality of the QMLE.

**Assumption Var( $\Theta$ ):** For any  $\theta \in \Theta$ ,  $(\sum_{i=1}^d \beta_i \frac{\partial f_\theta^0}{\partial \theta^{(i)}} = 0 \implies \forall i = 1, \dots, d, \beta_i = 0 \text{ a.s.})$  or  $(\sum_{i=1}^d \beta_i \frac{\partial H_\theta^0}{\partial \theta^{(i)}} = 0 \implies \forall i = 1, \dots, d, \beta_i = 0 \text{ a.s.})$ .

Moreover, one of the following technical assumption is required to establish the consistency of the model selection procedure.

**Assumption K( $\Theta$ ):** Assumptions  $A(f_\theta, \Theta)$ ,  $A(M_\theta, \Theta)$ ,  $A(\partial_\theta f_\theta, \Theta)$ ,  $A(\partial_\theta M_\theta, \Theta)$  and  $B(\Theta)$  hold and there exists  $r \geq 2$  such that  $\theta^* \in \Theta(r)$ . Moreover, with  $s = \min(1, r/3)$ , assume that the sequence  $(\kappa_n)_{n \in \mathbb{N}}$  satisfies

$$\sum_{k \geq 1} \left(\frac{1}{\kappa_k}\right)^s \left(\sum_{j \geq k} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) + \alpha_j(\partial_\theta f_\theta, \Theta) + \alpha_j(\partial_\theta M_\theta, \Theta)\right)^s < \infty.$$

**Assumption  $\tilde{K}(\Theta)$ :** Assumptions  $A(\tilde{H}_\theta, \Theta)$ ,  $A(\partial_\theta \tilde{H}_\theta, \Theta)$  and  $B(\Theta)$  hold and there exists  $r \geq 2$  such that  $\theta^* \in \Theta(r)$ . Moreover, with  $s = \min(1, r/4)$ , assume that the sequence  $(\kappa_n)_{n \in \mathbb{N}}$  satisfies

$$\sum_{k \geq 1} \left(\frac{1}{\kappa_k}\right)^s \left(\sum_{j \geq k} \alpha_j(\tilde{H}_\theta, \Theta) + \alpha_j(\partial_\theta \tilde{H}_\theta, \Theta)\right)^s < \infty.$$

**Remark 2.1.** These conditions on  $(\kappa_n)_{n \in \mathbb{N}}$  have been deduced from conditions for strong law of large numbers obtained in Kounias and Weng (1969) and are not too restrictive: for instance, if the Lipschitz coefficients of  $f_\theta$ ,  $M_\theta$  (the case using  $\tilde{H}_\theta$  can be treated similarly) and their derivatives are bounded by a geometric or Riemanian decrease:

1. Geometric case:  $\alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) + \alpha_j(\partial_\theta f_\theta, \Theta) + \alpha_j(\partial_\theta M_\theta, \Theta) = O(a^j)$  with  $0 \leq a < 1$ , then any  $(\kappa_n)$  such as  $1/\kappa_n = o(1)$  can be chosen; for instance  $\kappa_n = \log n$  or  $\log(\log n)$ ; this is the case for instance of ARMA, GARCH, APARCH or ARMA-GARCH processes.
2. Riemanian case:  $\alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) + \alpha_j(\partial_\theta f_\theta, \Theta) + \alpha_j(\partial_\theta M_\theta, \Theta) = O(j^{-\gamma})$  with  $\gamma > 1$ :



- if  $r \geq 3$  then
  - if  $\gamma > 2$  then any sequence such as  $1/\kappa_n = o(1)$  can be chosen;
  - if  $1 < \gamma < 2$ , any  $(\kappa_n)$  such as  $\kappa_n = O(n^\delta)$  with  $\delta > 2 - \gamma$  can be chosen.
- if  $1 \leq r < 3$ 
  - if  $\gamma > (r+3)/r$  then any sequence such as  $1/\kappa_n = o(1)$  can be chosen;
  - if  $1 < \gamma < (r+3)/r$  then any  $(\kappa_n)$  such as  $\kappa_n = n^\delta$  with  $\delta > (r+3)/r - \gamma$  can be chosen.

In the last case of these two conditions on  $r$ , we can see the usual BIC choice,  $\kappa_n = \log n$  does not fulfill the assumption in general.

### 2.3.2 Asymptotic model selection

Using the above assumptions, we can establish the limit theorem below, which provides sufficient conditions for the consistency of the model selection procedure.

**Theorem 2.1.** *Let  $(X_1, \dots, X_n)$  be an observed trajectory of an affine causal process  $X$  belonging to  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$  (or  $\widetilde{\mathcal{AC}}(\tilde{H}_\theta)$ ) where  $\theta^*$  is an unknown vector of  $\Theta$  a compact set included in  $\Theta(r) \subset \mathbb{R}^d$  (or  $\tilde{\Theta}(r) \subset \mathbb{R}^d$ ) with  $r \geq 4$ . If assumptions  $D(\Theta)$ ,  $Id(\Theta)$ ,  $K(\Theta)$  (or  $\tilde{K}(\Theta)$ ),  $A(\partial_{\theta^2}^2 f_\theta, \Theta)$  and  $A(\partial_{\theta^2}^2 M_\theta, \Theta)$  (or  $A(\partial_{\theta^2}^2 \tilde{H}_\theta, \Theta)$ ) also hold, then*

$$\mathbb{P}(\hat{m} = m^*) \xrightarrow[n \rightarrow \infty]{} 1 \quad \text{and} \quad \hat{\theta}(\hat{m}) \xrightarrow[n \rightarrow \infty]{\mathcal{P}} \theta^*. \quad (2.3.2)$$

The following theorem shows the asymptotic normality of the QMLE of the chosen model.

**Theorem 2.2.** *Under the assumptions of Theorem 2.1 and if  $\theta^* \in \overset{\circ}{\Theta}$  and  $Var(\Theta)$  holds, then*

$$\sqrt{n} \left( (\hat{\theta}(\hat{m}))_i - (\theta^*)_i \right)_{i \in m^*} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}_{|m^*|} \left( 0, F(\theta^*, m^*)^{-1} G(\theta^*, m^*) F(\theta^*, m^*)^{-1} \right) \quad (2.3.3)$$

where  $(F(\theta^*, m^*))_{i,j} = \mathbb{E} \left[ \frac{\partial^2 q_0(\theta^*)}{\partial \theta_i \partial \theta_j} \right]$  and  $(G(\theta^*, m^*))_{i,j} = \mathbb{E} \left[ \frac{\partial q_0(\theta^*)}{\partial \theta_i} \frac{\partial q_0(\theta^*)}{\partial \theta_j} \right]$  for  $i, j \in m^*$ .

**Remark 2.2.** *In Remark 2.1, we detailed some situations where the assumption  $K(\Theta)$  (or  $\tilde{K}(\Theta)$ ) holds, which leads to the results of Theorem 2.1 and 2.2. In particular, the  $\log n$  penalty usually linked to BIC is consistent in the case of a geometric decrease of the Lipschitz coefficients of the functions  $f_\theta$  and  $M_\theta$  (and their first order derivative). In the case of a Riemannian rate, the consistency of BIC is not ensured; see also the next section.*

## 2.4 Examples

In this section, some examples of time series satisfying the conditions of previous results are considered. These examples include  $AR(\infty)$ ,  $ARCH(\infty)$ ,  $APARCH(\delta, p, q)$  and  $ARMA(p, q)$ -GARCH( $p', q'$ ).

### 2.4.1 AR( $\infty$ ) models

For  $(\psi_k(\theta))_{k \in \mathbb{N}}$  a sequence of real numbers depending on  $\theta \in \mathbb{R}^d$ , let us consider an AR( $\infty$ ) process defined by:

$$X_t = \sum_{k \geq 1} \psi_k(\theta^*) X_{t-k} + \sigma \xi_t \quad \text{for any } t \in \mathbb{Z}, \quad (2.4.1)$$

where  $(\xi_t)_t$  admits 4-order moments, and  $\theta^* \in \Theta \subset \Theta(4)$ , the set of  $\theta \in \mathbb{R}^d$  such that  $\sum_{k \geq 1} \|\psi_k(\theta)\|_{\Theta} < 1$  and  $\sigma > 0$ . This process corresponds to (2.2.6) with  $f_{\theta}((x_i)_{i \geq 1}) = \sum_{k \geq 1} \psi_k(\theta) x_k$  and  $M_{\theta} \equiv \sigma$  for any  $\theta \in \Theta$ . The Lipschitz coefficients of  $f_{\theta}$  are  $\alpha_k(f_{\theta}) = \|\psi_k(\theta)\|_{\Theta}$ . Moreover, Assumption  $D(\Theta)$  holds with  $\underline{h} = \sigma^2 > 0$ .

Let us consider  $\mathcal{M}$  a finite family of models. Of course, the main example of such family of models is given by the one of ARMA( $p, q$ ) processes with  $0 \leq p \leq p_{\max}$  and  $0 \leq q \leq q_{\max}$ , providing  $(p_{\max} + 1)(q_{\max} + 1)$  models and  $\theta \in \mathbb{R}^{p_{\max} + q_{\max} + 1}$ .

Besides, assume that  $Id(\Theta)$ ,  $\text{Var}(\Theta)$  hold and that the sequence  $(\psi_k)$  is twice differentiable (with respect to  $\theta$ ) on  $\Theta$ , with  $\sum_k \|\partial_{\theta}^2 \psi_k(\theta)\|_{\Theta} < \infty$  and  $\|\psi_k(\theta)\|_{\Theta} + \|\partial_{\theta} \psi_k(\theta)\|_{\Theta} = O(k^{-\gamma})$  with  $\gamma > 1$ . From Remark 2.1,

- if  $\gamma > 2$ , the condition  $\kappa_n \xrightarrow{n \rightarrow \infty} \infty$  (for instance, the BIC penalization with  $\kappa_n = \log(n)$ , or  $\kappa_n = \sqrt{n}$ ) ensures the consistency of  $\hat{m}$  and the Theorem (2.2) holds if in addition  $\theta^* \in \overset{\circ}{\Theta}$ ;
- if  $1 < \gamma < 2$ ,  $\kappa_n = O(n^{\delta})$  with  $\delta > 2 - \gamma$  has to be chosen (and we cannot insure the consistency of  $\hat{m}$  in case of classical BIC penalization).

Finally, in the particular case of the family of ARMA processes, the stationarity condition implies that any  $\kappa_n \xrightarrow{n \rightarrow \infty} \infty$  can be chosen (BIC penalization with  $\kappa_n = \log(n)$ , or  $\kappa_n = \sqrt{n}$ ), since the decreases of  $\psi_k$  and its derivative are exponential.

### 2.4.2 ARCH( $\infty$ ) models

For  $(\psi_k(\theta))_{k \in \mathbb{N}}$  a sequence of nonnegative real numbers depending on  $\theta \in \mathbb{R}^d$ , with  $\psi_0 > 0$ , let us consider an ARCH( $\infty$ ) process defined by :

$$X_t = \left( \psi_0(\theta^*) + \sum_{k=1}^{\infty} \psi_k(\theta^*) X_{t-k}^2 \right)^{1/2} \xi_t \quad \text{for any } t \in \mathbb{Z}, \quad (2.4.2)$$

where  $\mathbb{E}[\xi_0^4] < \infty$ , and  $\theta^* \in \Theta \subset \tilde{\Theta}(4)$ , the set of  $\theta \in \mathbb{R}^d$  such that  $\sum_{k \geq 1} \|\psi_k(\theta)\|_{\Theta} < 1$ . This process corresponds to (2.2.6) with  $f_{\theta}((x_i)_{i \geq 1}) \equiv 0$  and  $H_{\theta}((x_i)_{i \geq 1}) = \psi_0(\theta) + \sum_{k=1}^{\infty} \psi_k(\theta) x_k^2$ , i.e.  $\tilde{H}_{\theta}((y_i)_{i \geq 1}) = \psi_0(\theta) + \sum_{k=1}^{\infty} \psi_k(\theta) y_k$ , for any  $\theta \in \Theta$ . The Lipschitz coefficients of  $\tilde{H}_{\theta}$  are  $\alpha_k(\tilde{H}_{\theta}) = \|\psi_k(\theta)\|_{\Theta}$ . Moreover, Assumption  $D(\Theta)$  holds if  $\underline{h} = \inf_{\theta \in \Theta} \psi_0(\theta) > 0$ .

Let us consider  $\mathcal{M}$  a finite family of models. The main example of such family of models is given by the GARCH( $p, q$ ) processes with  $0 \leq p \leq p_{\max}$  and  $0 \leq q \leq q_{\max}$ , providing  $(p_{\max} + 1)(q_{\max} + 1)$  models and  $\theta \in \mathbb{R}^{p_{\max} + q_{\max} + 1}$ .

Moreover, assume that  $Id(\Theta)$ ,  $\text{Var}(\Theta)$  hold and that the sequence  $(\psi_k)$  is twice differentiable (with respect to  $\theta$ ) on  $\Theta$ , with  $\sum_k \|\partial_{\theta}^2 \psi_k(\theta)\|_{\Theta} < \infty$  and  $\|\psi_k(\theta)\|_{\Theta} + \|\partial_{\theta} \psi_k(\theta)\|_{\Theta} = O(k^{-\gamma})$  with  $\gamma > 1$ . From Remark 2.1,

- if  $\gamma > 2$ , the condition  $\kappa_n \xrightarrow[n \rightarrow \infty]{} \infty$  (for instance, the BIC penalization with  $\kappa_n = \log(n)$ , or  $\kappa_n = \sqrt{n}$ ) ensures the consistency of  $\hat{m}$  and the Theorem (2.2) holds if in addition,  $\theta^* \in \overset{\circ}{\Theta}$ ;
- if  $1 < \gamma < 2$ ,  $\kappa_n = O(n^\delta)$  with  $\delta > 2 - \gamma$  has to be chosen (and we cannot insure the consistency of  $\hat{m}$  in the case of the classical BIC penalization).

Finally, in the particular case of the family of GARCH processes, the stationarity condition implies that any  $\kappa_n \xrightarrow[n \rightarrow \infty]{} \infty$  can be chosen (BIC penalization with  $\kappa_n = \log(n)$ , or  $\kappa_n = \sqrt{n}$ ), since the decreases of  $\psi_k$  and its derivative are exponential.

### 2.4.3 APARCH( $\delta, p, q$ ) models

For  $\delta \geq 1$  and from Ding et al. (1993),  $(X_t)_{t \in \mathbb{Z}}$  is an APARCH( $\delta, p, q$ ) process with  $p, q \geq 0$  if:

$$\begin{cases} X_t = \sigma_t \xi_t \\ (\sigma_t)^\delta = \omega + \sum_{i=1}^p \alpha_i (|X_{t-i}| - \gamma_i X_{t-i})^\delta + \sum_{j=1}^q \beta_j (\sigma_{t-j})^\delta \end{cases} \quad \text{for any } t \in \mathbb{Z}, \quad (2.4.3)$$

where  $\omega > 0$ ,  $-1 < \gamma_i < 1$ ,  $\alpha_i \geq 0$ ,  $\beta_j \geq 0$  for  $1 \leq i \leq p$  and  $1 \leq j \leq q$ ,  $\alpha_p > 0$ ,  $\beta_q > 0$  and  $\sum_{j=1}^q \beta_j < 1$ . From Bardet et al. (2017), with  $\theta = (\omega, \alpha_1, \dots, \alpha_p, \gamma_1, \dots, \gamma_p, \beta_1, \dots, \beta_p)'$ , the conditional variance  $\sigma_t$  can be rewritten as follows

$$\sigma_t^\delta = b_0(\theta) + \sum_{k \geq 1} \left( b_k^+(\theta) (\max(X_{t-k}, 0))^\delta - b_k^-(\theta) (\min(X_{t-k}, 0))^\delta \right);$$

with  $f_\theta \equiv 0$  and  $M_\theta^t = \sigma_t$ , we deduce that  $\alpha_k(M_\theta, \Theta) = \max(\|b_k^+(\theta)\|_\Theta^{1/\delta}, \|b_k^-(\theta)\|_\Theta^{1/\delta})$ , and from the assumption  $\sum_{j=1}^q \beta_j < 1$ , the Lipschitz coefficients  $\alpha_k(M_\theta, \Theta)$  decrease exponentially fast. Then, the stationarity set for  $r \geq 1$  is

$$\Theta(r) = \left\{ \theta \in \mathbb{R}^{2p+q+1} \mid \|\xi_0\|_r \sum_{j=1}^{\infty} \max(|b_j^+(\theta)|^{1/\delta}, |b_j^-(\theta)|^{1/\delta}) < 1 \right\}.$$

Now, assume that  $(X_t)_{t \in \mathbb{Z}}$  is an APARCH( $\delta, p^*, q^*$ ) where  $0 \leq p^* \leq p_{\max}$  and  $0 \leq q^* \leq q_{\max}$  are unknown orders as well as the other parameters:  $\omega^* > 0$ ,  $-1 < \gamma_i^* < 1$ ,  $\alpha_i^* \geq 0$ ,  $\beta_j^* \geq 0$  for  $1 \leq i \leq p_{\max}$  and  $1 \leq j \leq q_{\max}$ ,  $\alpha_{p^*} > 0$ ,  $\beta_{q^*} > 0$ .

Let  $\mathcal{M}$  be the family of APARCH( $\delta, p, q$ ) processes, with  $0 \leq p \leq p_{\max}$  and  $0 \leq q \leq q_{\max}$ . As a consequence, we consider here  $d = 2p_{\max} + q_{\max} + 1$ , and

$$\theta^* = {}^t(\omega^*, \alpha_1^*, \dots, \alpha_{p^*}^*, 0, \dots, 0, \gamma_1^*, \dots, \gamma_{p^*}^*, 0, \dots, 0, \beta_1^*, \dots, \beta_{q^*}^*, 0, \dots, 0) \in \mathbb{R}^d.$$

With all the previous conditions, assumptions  $D(\Theta)$ ,  $\text{Id}(\Theta)$ ,  $\text{Var}(\Theta)$  are satisfied. Moreover, since the Lipschitz coefficients decrease exponentially fast,  $K(\Theta)$  is satisfied when  $\kappa_n \rightarrow \infty$ . Therefore, the consistency Theorem (2.1) and the Theorem (2.2) of the estimator of the chosen model are satisfied when  $r = 4$  and  $\kappa_n \rightarrow \infty$  (for instance with the typical BIC penalty  $\kappa_n = \log n$ ).

#### 2.4.4 ARMA( $p, q$ )-GARCH( $p', q'$ ) models

From Ding et al. (1993) and Ling and McAleer (2003), we define  $(X_t)_{t \in \mathbb{Z}}$  as an (invertible) ARMA( $p, q$ )-GARCH( $p', q'$ ) process with  $p, q, p', q' \geq 0$  if:

$$\begin{cases} X_t = \sum_{i=1}^p a_i X_{t-i} + \varepsilon_t - \sum_{i=1}^q b_i \varepsilon_{t-i} \\ \varepsilon_t = \sigma_t \xi_t, \text{ with } \sigma_t^2 = c_0 + \sum_{i=1}^{p'} c_i \varepsilon_{t-i}^2 + \sum_{i=1}^{q'} d_i \sigma_{t-i}^2 \end{cases} \quad \text{for all } t \in \mathbb{Z},$$

where

- $c_0 > 0$ ,  $c_{p'} > 0$ ,  $c_i \geq 0$  for  $i = 1, \dots, p' - 1$  and  $d_{q'} > 0$ ,  $d_i \geq 0$  for  $i = 1, \dots, q' - 1$ ;
- $P(x) = 1 - \sum_{i=1}^p a_i x^i$  and  $Q(x) = 1 - \sum_{i=1}^{q'} b_i x^i$  are coprime polynomials.

Here we will consider the case of a stationary invertible ARMA( $p, q$ )-GARCH( $p', q'$ ) process such as  $\|X_0\|_4 < \infty$  and therefore we will consider:

$$\begin{aligned} \Theta_{p,q,p',q'} = \left\{ (a_1, \dots, d_{q'}) \in \mathbb{R}^{p+q+p'+1+q'}, \sum_{j=1}^{q'} d_j + \|\xi_0\|_4 \sum_{j=1}^{p'} c_j < 1 \right. \\ \left. \text{and } \left(1 - \sum_{j=1}^p a_j z^j\right) \left(1 - \sum_{j=1}^{q'} b_j z^j\right) \neq 0 \text{ for all } |z| \leq 1 \right\}. \end{aligned}$$

Therefore, if  $(a_1, \dots, d_{q'}) \in \Theta_{p,q,p',q'}$ ,  $(\varepsilon_t)_t$  is a stationary GARCH( $p', q'$ ) process and  $(X_t)_t$  is a stationary weak invertible ARMA( $p, q$ ) process.

Moreover, following Lemma 2.1. of Bardet et al. (2017), we know that a stationary ARMA( $p, q$ )-GARCH( $p', q'$ ) process is a stationary affine causal process with functions  $f_\theta$  and  $M_\theta$  satisfying the Assumption A( $f_\theta, \Theta$ ) and A( $M_\theta, \Theta$ ) with Lipschitzian coefficients decreasing exponentially fast, as well as their derivatives. Finally, if  $\Theta$  is a bounded subset of  $\Theta_{p,q,p',q'}$ , then assumptions D( $\Theta$ ), Id( $\Theta$ ) and Var( $\Theta$ ) are automatically satisfied.

Assume now that  $(X_t)_{t \in \mathbb{Z}}$  is an ARMA( $p^*, q^*$ )-GARCH( $p'^*, q'^*$ ) process where  $0 \leq p^* \leq p_{\max}$ ,  $0 \leq q^* \leq q_{\max}$ ,  $0 \leq p'^* \leq p'_{\max}$  and  $0 \leq q'^* \leq q'_{\max}$  are unknown orders with also unknown parameters:  $c_0^*, \dots, c_{p'^*}^*, d_1^*, \dots, d_{q'^*}^*, a_1^*, \dots, a_{p^*}^*, b_1^*, \dots, b_{q^*}^*$ .

Let  $\mathcal{M}$  be the family of ARMA( $p, q$ )-GARCH( $p', q'$ ) processes, with  $0 \leq p \leq p_{\max}$ ,  $0 \leq q \leq q_{\max}$ ,  $0 \leq p' \leq p'_{\max}$  and  $0 \leq q' \leq q'_{\max}$ . Hence, we consider here  $d = p_{\max} + q_{\max} + p'_{\max} + q'_{\max} + 1$ , and

$$\begin{aligned} \theta^* = (c_0^*, \dots, c_{p'^*}^*, 0, \dots, 0, d_1^*, \dots, d_{q'^*}^*, 0, \dots, 0, \\ a_1^*, \dots, a_{p^*}^*, 0, \dots, 0, b_1^*, \dots, b_{q^*}^*, 0, \dots, 0) \in \mathbb{R}^d. \end{aligned}$$

With  $\Theta$  a bounded subset of  $\Theta_{p_{\max}, q_{\max}, p'_{\max}, q'_{\max}}$ , all the previous assumptions D( $\Theta$ ), Id( $\Theta$ ), Var( $\Theta$ ) are satisfied and K( $\Theta$ ) is also satisfied as soon as  $\kappa_n \rightarrow \infty$ . As a consequence, in this framework the consistency Theorem (2.1) and the Theorem (2.2) of the estimator of the chosen model are satisfied when  $r = 4$  and  $\kappa_n \rightarrow \infty$  (for instance with the typical BIC penalty  $\kappa_n = \log n$ ).

## 2.5 Portmanteau test

From the above section, we are now able to asymptotically pick up a best model in a family of models. We can also obtain asymptotic confident regions of the estimated parameter

of the chosen model. However, it is also important to check whether the chosen model is appropriate. This section attempts to answer this question by constructing a portmanteau test as a diagnostic tool based on the squares of the residuals sequence of the chosen model. This test has been widely considered in the time series literature, with procedures based on the squared residual correlogram (see for instance Li and Mak (1994), Ling and Li (1997) ) and the absolute residual (or usual residuals) correlogram (see for instance Li (1992), Duchesne and Francq (2008), Li and Li (2008)), among others.

Since our goal is to provide an efficient test for the entire affine class that contains weak white noise processes. We consider in this setting the autocorrelation of the squared residuals and follow the same scheme of procedure used in (Li and Mak (1994), Ling and Li (1997)) while relying on some of their results. But three main differences need to be pointed out:

- the results of Li and Mak (1994) are based on the exact likelihood of the data, which is then assumed to be known. But it is not at all the case even for simple ARMA(1,1) or GARCH(1,1) processes. By working directly on the quasi-likelihood, we really proposes a feasible Portemanteau test;
- we provide more detailed sufficient conditions to get the asymptotic results of the Portmanteau test;
- our procedure is also applied to the selected model, which is not necessarily the true model.

For  $m \in \mathcal{M}$ , for  $K$  a positive integer, denote the vector of adjusted correlogram of squared residuals by:

$$\hat{\rho}(m) := (\hat{\rho}_1(m), \dots, \hat{\rho}_K(m))',$$

where for  $k = 1, \dots, K$ ,  $\hat{\rho}_k(m) := \frac{\hat{\gamma}_k(m)}{\hat{\gamma}_0(m)}$  with

$$\hat{\gamma}_k(m) := \frac{1}{n} \sum_{t=k+1}^n (\hat{e}_t^2(m) - 1)(\hat{e}_{t-k}^2(m) - 1) \text{ and } \hat{e}_t(m) := (\hat{M}_{\hat{\theta}(m)}^t)^{-1}(X_t - \hat{f}_{\hat{\theta}(m)}^t).$$

Finally, the following theorem provides central limit theorems for  $\hat{\rho}(m^*)$  and  $\hat{\rho}(\hat{m})$  as well as for a portmanteau test statistic.

**Theorem 2.3.** *Under the assumptions of Theorem 2.2, with also*

- $\mathbb{E}[\xi_0^3] = 0$ ;
- 

$$\sum_{t=1}^{\infty} t^{-1/4} \left( \sum_{j \geq t} \alpha_j(f_{\theta}, \Theta) + \alpha_j(M_{\theta}, \Theta) \right)^{1/2} < \infty$$

$$\text{or } \sum_{t=1}^{\infty} t^{-1/4} \left( \sum_{j \geq t} \alpha_j(\tilde{H}_{\theta}, \Theta) \right)^{1/2} < \infty,$$

then,

1. With  $V(\theta^*, m^*)$  defined in (2.7.39), it holds that

$$\sqrt{n} \hat{\rho}(m^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}_K(0, V(\theta^*, m^*)). \quad (2.5.1)$$

2. With  $\hat{Q}_K(m^*) := n \hat{\rho}(m^*)' (V(\hat{\theta}(m^*), m^*))^{-1} \hat{\rho}(m^*)$ , we have

$$\hat{Q}_K(m^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \chi^2(K). \quad (2.5.2)$$

3. The previous points 1. and 2. also hold when  $m^*$  is replaced by  $\hat{m}$ .

Using the Theorem 2.3, we can asymptotically test:

$$\begin{cases} H_0 : \exists m^* \in \mathcal{M}, \text{ such as } (X_1, \dots, X_n) \text{ is a trajectory of } X \in \mathcal{AC}(M_\theta, f_{\theta^*}) \\ H_1 : \nexists m^* \in \mathcal{M}, \text{ such as } (X_1, \dots, X_n) \text{ is a trajectory of } X \in \mathcal{AC}(M_\theta, f_{\theta^*}) \end{cases}.$$

with  $\theta^* \in \Theta(m^*)$  in both cases.

Therefore,  $\hat{Q}_K(\hat{m})$  can be used as a portmanteau test statistic to decide between  $H_0$  and  $H_1$  and diagnose the goodness-of-fit of the selected model.

**Remark 2.3.** 1. In practice the constant  $\mu_4$  and the columns of the matrix  $J_K(m^*)$  (see (2.7.35)) involved in  $V(\theta^*, m^*)$  are estimated by the correspondent sample average; they are respectively  $\hat{\mu}_4 = \frac{1}{n} \sum_{t=1}^n (\hat{e}_t(\hat{m}))^4$  and  $(\hat{J}_K(\hat{\theta}(\hat{m})))_{.,k} = \frac{1}{n} \sum_{t=1}^{n-k} [(\hat{e}_t(\hat{m}))^2 - 1] \partial_\theta \log(M_{\hat{\theta}(\hat{m})}^{t+k})$ .

2. For  $AR(\infty)$  models (and then for causal invertible  $ARMA(p, q)$ ), since  $M_\theta = \sigma$  as we have seen in Sub-section 2.4.1, we deduce from (2.7.39) that  $V(\theta^*, m^*) = I_K$  as  $J_K(m^*) = 0$ . Hence, in such a case, we simply obtained:

$$\hat{Q}_K(\hat{m}) = n \|\hat{\rho}(\hat{m})\|^2 \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \chi^2(K). \quad (2.5.3)$$

Note that working with autocorrelations of squared residuals rather than those of residuals, avoids the need to subtract the number of estimated parameters in the asymptotic chi-square distribution. Hence our result is valid for any  $K \in \mathbb{N}^*$ .

## 2.6 Numerical Results

This section features some simulation experiments that are performed to assess the usefulness of the asymptotic results obtained in Section 2.3. Each model is generated independently 1000 times over a trajectory of length  $n$ . Different sample sizes are considered to identify possible discrepancies between asymptotically expected properties and those obtained at finite distance. We will consider  $n$  belongs to  $\{100, 500, 1000, 2000\}$ . The process used to generate the trajectory is indicated each time. Throughout this section,  $(\xi_t)$  represents a Gaussian white noise with variance unity.

### 2.6.1 Monte-Carlo experiments for common time series selection

We first generate some classical models as "true" models  $m^*$ :

1. Model 1,  $AR(2)$  process:  $X_t = 0.4X_{t-1} + 0.4X_{t-2} + \xi_t$ ;
2. Model 2,  $ARMA(1, 1)$  process:  $X_t = 0.3X_{t-1} + \xi_t + 0.5\xi_{t-1}$ ;
3. Model 3,  $ARCH(2)$  process:  $X_t = \xi_t \sqrt{0.2 + 0.4X_{t-1}^2 + 0.2X_{t-2}^2}$ ;

4. Model 4, GARCH(1, 1) process:  $X_t = \sigma_t \xi_t$ , with  $\sigma_t^2 = 0.2 + 0.3X_{t-1}^2 + 0.5\sigma_{t-1}^2$ .

We considered as competitive models all the models in the family  $\mathcal{M}$  defined by:

$$\mathcal{M} = \{\text{ARMA}(p, q) \text{ or GARCH}(p', q') \text{ processes with } 0 \leq p, q, p' \leq 5, 1 \leq q' \leq 5\}.$$

As a consequence, there are 66 candidate models. Note also that in our simulations, since we have more than one model per dimension, slope estimation is done after considering the "best model" (which maximizes quasi-log likelihood) within each dimension.

The results of the model selection procedure are displayed in Table 2.1. More precisely, for each penalty ( $\log n$ ,  $\sqrt{n}$ ) the frequency that the associated criterion selects respectively a wrong model, the true model and an overfitted model (here a model that contains the true model).

**Table 2.1:** Percentage of selected order based on 1000 replications depending on sample's length for Model 1, 2, 3 and 4 respectively.

Sample length $n$		100		500		1000		2000	
	Penalty	$\log n$	$\sqrt{n}$	$\log n$	$\sqrt{n}$	$\log n$	$\sqrt{n}$	$\log n$	$\sqrt{n}$
Model 1	Wrong	21.4	32.3	1.7	0.8	0.8	0.1	0.2	0
	True	74.2	67.6	97.2	99.2	98.2	99.9	99.2	100
	Overfitted	4.4	0.1	1.1	0	1.0	0	0.6	0
Model 2	Wrong	30.4	57.7	4.8	4.2	0.7	0.3	0.4	0
	True	64.1	42.1	93.6	95.8	98.2	99.7	99.2	100
	Overfitted	5.5	0.2	1.6	0	1.1	0	0.4	0
Model 3	Wrong	76.1	90.8	27.3	67.1	14.0	41.5	4.6	12.0
	True	23.8	9.2	72.7	32.9	85.9	58.5	95.4	88.0
	Overfitted	0.1	0	0	0	0.1	0	0	0
Model 4	Wrong	83.8	94.3	22.1	61.5	5.8	31.3	1.8	6.2
	True	15.9	5.7	77.5	38.5	93.2	68.7	98.0	93.8
	Overfitted	0.3	0	0.4	0	1.0	0	0.2	0

From these results, it is clear that the consistency of our model selection procedure is numerically convincing, which is in accordance with Theorem 2.1, where both penalties ( $\log n$ ,  $\sqrt{n}$ ) lead to consistent criteria for the four models under consideration. Note also that the typical BIC  $\log n$  penalty is more interesting for retrieving the true model than the  $\sqrt{n}$ -penalized likelihood for a small sample size. But the larger the sample size, the more accurate the  $\sqrt{n}$  penalty is, compared to the  $\log n$  penalty.

In addition, for each of the three models, we also applied the portmanteau test statistic  $\hat{Q}_K(\hat{m})$ , using the  $\sqrt{n}$  penalty. Table 2.2 shows the empirical size and empirical power of this test. We call by empirical size, the percentage of falsely rejecting the null hypothesis  $H_0$ . On the other hand, the empirical power represents the percentage of rejection of  $H_0$  when we arbitrary chose a false model, which is a AR(3) process  $X_t = 0.2X_{t-1} + 0.2X_{t-2} + 0.4X_{t-3} + \xi_t$  for Model 1 and 2, and a ARCH(3) process  $X_t = \xi_t \sqrt{0.4 + 0.2X_{t-1}^2 + 0.2X_{t-2}^2 + 0.2X_{t-3}^2}$  for Model 3 and 4.

It is important to note that choosing the maximum number of lags  $K$  is sometimes tricky.

To our knowledge, there is no real theoretical study to justify the choice of one value or another. However, some Monte Carlo simulations have suggested some ways to make a good choice. For instance Li and Mak (1994) suggested that the autocorrelations  $\hat{\rho}_k(\hat{m})$  with  $1 \leq k \leq K$  have a better asymptotic behaviour for small values of  $k$ . Therefore, the finite sample performance of the size and power of the test may also vary with the choice of  $K$  and could be better for small values of  $K$ . On the other hand, Tse and Zuo (1997) suggested that  $K = p + q + 1$  may be an appropriate choice for the GARCH( $p, q$ ) family. Thus, in our tests, we consider  $K = 3$  and  $K = 6$  so that the rejection is based on the upper 5th percentile of the  $\chi^2(3)$  distribution on the one hand and  $\chi^2(6)$  on the other hand.

**Table 2.2:** The empirical size and empirical power of the portmanteau test statistic  $\hat{Q}_K(\hat{m})$  based on 1000 independent replications (in %) with  $K = 3$  and  $K = 6$ .

n		100		500		1000		2000	
		size	power	size	power	size	power	size	power
$K = 3$	Model 1	3.3	10.9	6.2	52.2	3.5	84.8	5.0	98.2
	Model 2	3.3	7.0	4.8	23.3	6.2	42.4	4.9	70.4
	Model 3	4.6	6.4	8.4	44.1	14.3	81.0	36.9	99.4
	Model 4	9.5	23.2	21.3	38.5	33.6	57.2	39.4	88.3
$K = 6$	Model 1	2.9	9.1	4.9	42.0	4.4	76.3	4.5	97.6
	Model 2	3.0	6.3	5.2	18.0	5.1	35.1	4.6	60.2
	Model 3	4.5	12.6	11.1	64.4	14.7	92.5	27.9	99.9
	Model 4	4.3	52.7	4.2	98.6	3.2	99.6	3.6	99.9

Once again, the results of Table 2.2 numerically confirms the asymptotic results of Theorem 2.3. Remark that the test is more powerful by using values of  $K$  not too large as mentioned above especially for small samples.

## 2.6.2 Subset model selection

Now, we exhibit the performance of the previously considered criteria on a particular case of dimension selection. The process generated data is considered as follows:

$$\text{Model 5 : } X_t = 0.4X_{t-3} + 0.4X_{t-4} + \xi_t.$$

Here, we will consider the case of a nonhierarchical but exhaustive family  $\mathcal{M}$  of AR(4) models, *i.e.*

$$\mathcal{M} = \mathcal{P}(\{1, 2, \dots, 10\})$$

$$\implies X_t = \theta_1 X_{t-1} + \theta_2 X_{t-2} + \dots + \theta_{10} X_{t-10} + \xi_t$$

$$\text{and } \theta = (\theta_1, \theta_2, \dots, \theta_{10})' \in \Theta(m).$$

As a consequence,  $1024 = 2^{10}$  candidate models are considered and Table 2.3 presents the results of the selection procedure.

**Table 2.3:** Percentage of selected model based on 1000 replications depending on sample's length for Model 4



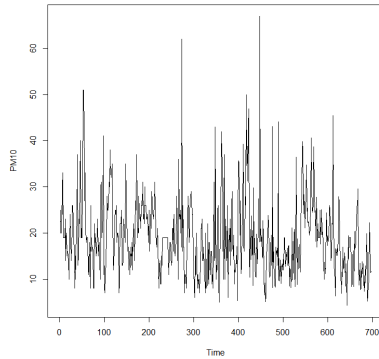
Sample length	100		500		1000		2000	
	$\log n$	$\sqrt{n}$	$\log n$	$\sqrt{n}$	$\log n$	$\sqrt{n}$	$\log n$	$\sqrt{n}$
True model	70.4	67.3	90.0	100	93.2	100	95.3	100
Overfitted	25.0	1.6	10	0	6.8	0	4.7	0
False model	4.6	31.1	0	0	0	0	0	0

We deduce that the consistency of our model selection procedure is also numerically convincing in this case of exhaustive model selection, which is in accordance with Theorem 2.1.

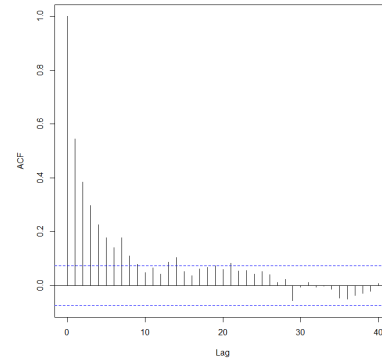
## 2.6.3 Application to real data

### 2.6.3.a Air quality analysis

Air quality, which can be defined as the level of cleanliness of the air, is probably one of the first health and environmental concerns of this new century. With the increasing number of human activities, the air is being degraded by a wide variety of pollutants, including PM. PM stands for particulate matter [government \(2017\)](#): the term for a mixture of solid particles and liquid droplets found in the air. Some particles, such as dust, dirt, soot, or smoke, are large or dark enough to be seen with the naked eye. Let consider daily observations of PM10 (downloaded from [Air PACA](#)) at Marseille Kaddouz station (France) from January 1, 2018 to November 30, 2019. This is a time series trajectory of length  $n = 698$  (see Figure 2.1a). We are going to use our model selection criteria to identify the "best" model for this time series.



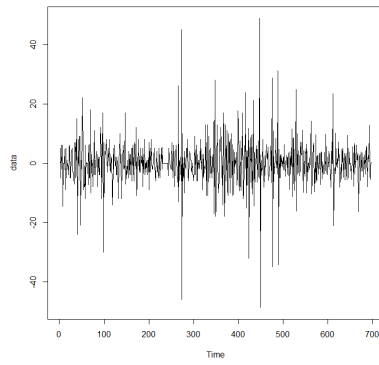
(a) Time plot of PM10 levels.



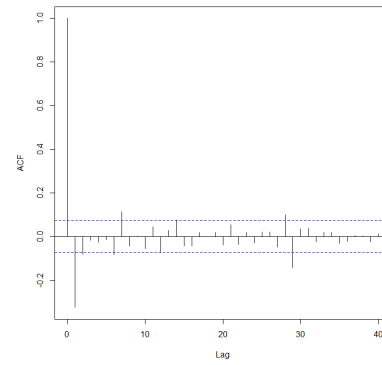
(b) The sample autocorrelation function for the Marseille PM10 showing the bounds  $\pm 1.96/\sqrt{n}$ .

**Figure 2.1:** The Marseille PM10 levels (January 1st, 2018 to November 30, 2019).

An inspection of the Figure 2.1 may suggest us a family of candidate models. First, the slow decrease of the sample autocorrelation (up to lag 6), suggests that there is a component trend in the variability of the PM10. Also, a close inspection of the data shows that pollution is on average much lower on weekends than on working days. So before identifying a plausible family of models, let consider the detrended time series by differencing (see Figure 2.2). Therefore, we use the same family  $\mathcal{M}$  already considered in Subsection 2.6.1 that provides us 66 candidate models. For each model, we compute the criterion (2.2.4) with  $\kappa_n = \log(n)$  and  $\kappa_n = \sqrt{n}$ . The selection results and also the goodness-of-fit of the selected model are featured in the Table 2.4.



(a) PM10 levels.



(b) The sample autocorrelation function for the Marseille PM10 showing the bounds  $\pm 1.96/\sqrt{n}$ .

**Figure 2.2:** Elimination of trend and seasonality in Marseille PM10 levels (January 1st, 2018 to November 30, 2019).

**Table 2.4:** Summary of the results of the model selection and goodness-of-fit analysis on PM10.

	$\kappa_n = \log(n)$	$\kappa_n = \sqrt{n}$
$\hat{m}$	ARMA(1, 2)	ARMA(1, 1)
$\hat{Q}_{10}(\hat{m})$	11.09	18.02
$p - value$	0.35	0.055

This table shows that all p-values are greater than 0.05, and then none of the test statistics leads us to reject the null hypothesis at this level even though the case of the ARMA(1, 1) is somehow limit. The chosen ARMA(1, 2) seems to be the more suitable model for PM10 time series.

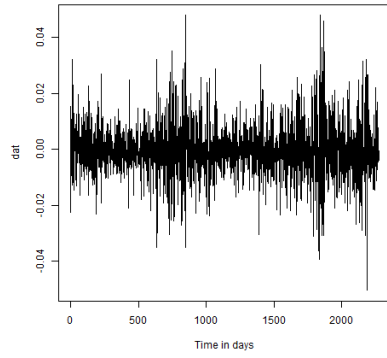
### 2.6.3.b Financial index analysis

We consider the returns of the daily closing prices of the FTSE 100 index and also the SP 500. They are respectively 2273 and 2264 observations from January 4th, 2010 to December 31st, 2018 for FTSE 100 and SP500. The time plot and the correlograms for the log-returns and squared log-returns are plotted in Figure 2.3. Figures 2.3a and 2.3c exhibit the conditional heteroskedasticity in the log-return time series. Moreover, Figure 2.3b shows that more than 5 per cent of the autocorrelations are out of the confidence interval  $\pm 1.96/\sqrt{2273}$  and specially the Figure 2.3d suggests that the strong white noise assumption cannot be sustained for this log-returns sequence of FTSE index. We also have the same conclusion for SP 500 (see Figure 2.4)

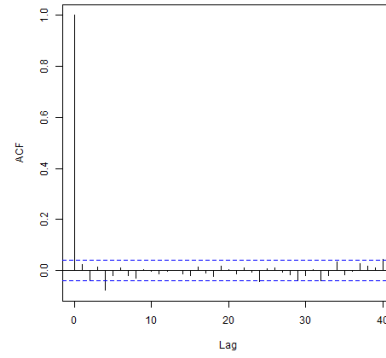
As in the previous illustrative example, the ARMA-GARCH is a plausible family for modeling of the FTSE 100 and SP 500 index. The  $\log n$  and  $\sqrt{n}$  penalizations have been applied to identify the best order and the goodness-of-fit of the selected model has been tested by the Portmanteau test.

**Table 2.5:** Summary of the results of the model selection and goodness-of-fit analysis on FTSE index.

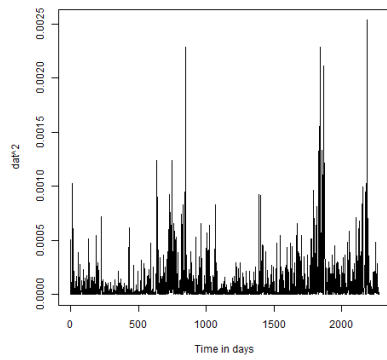
	$\kappa_n = \log(n)$	$\kappa_n = \sqrt{n}$
$\hat{m}$	GARCH(1, 1)	GARCH(1, 1)
$\hat{Q}_{10}(\hat{m})$	9.30	9.30
$p - value$	0.50	0.50



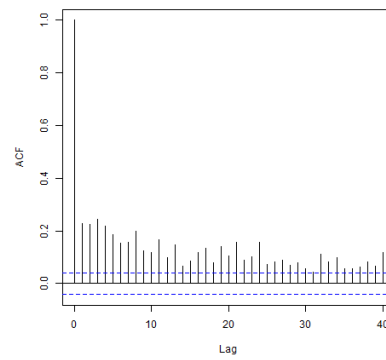
(a) Time plot of log-returns.



(b) Correlograms of log-returns.



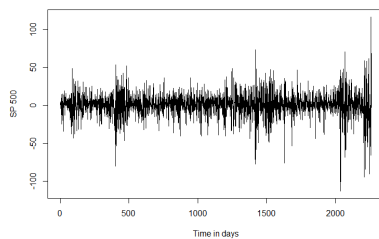
(c) Time plot of squared log-returns.



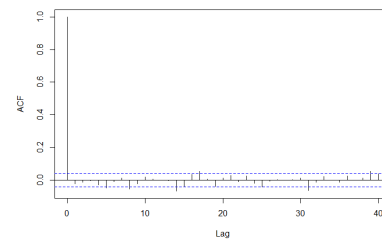
(d) Correlograms of squared log-returns.

**Figure 2.3:** Daily closing FTSE 100 index (January 4th, 2010 to December 31 st, 2018).

The GARCH(1,1) is the "best" model based on the three criteria considered and it is adequate (at level 0.95) to model the FTSE 100 index. Regarding the SP 500 index, the GARCH(1,1) is still the best model based on all three criteria and  $\hat{Q}_{10}(\hat{m}) = 15.2$  associated with a p-value of 0.12. These results are not surprising since the GARCH(1,1) is the reference model and the most commonly used in empirical studies. In addition, [Francq and Zakoian \(2010\)](#) found the GARCH(1,1) to be adequate using a FTSE 100 trajectory from April 3, 1984 to April 3, 2007 and January 3, 1950 to July 24, 2009 for SP 500.



(a) Time plot of log returns .



(b) The sample autocorrelation function showing the bounds  $\pm 1.96/\sqrt{n}$ .

**Figure 2.4:** Daily closing price of SP500 (January 4th, 2010 to December 31 st, 2018).

## 2.7 Proofs

We start with the proof of the Proposition 2.1.

*Proof.* For ease of writing, consider only the general case where  $f_{\theta_i}^{(i)} = g_{\alpha_i}^{(i)}$  and  $M_{\theta_i}^{(i)} = N_{\beta_i}^{(i)}$  where  $\theta_i = {}^t(\alpha_i, \beta_i)$  for  $i = 1, 2$ . Now, assume that there exist  $\alpha \in \mathbb{R}^\delta$ , where  $0 \leq \delta \leq \min(d_1, d_2)$  and a function  $h_\alpha$  such as  $g_{\alpha_1}^{(1)} = h_\alpha + \ell_{\alpha'_1}^{(1)}$ ,  $f_{\alpha_2}^{(2)} = h_\alpha + \ell_{\alpha'_2}^{(2)}$  with  $\alpha_1 = {}^t(\alpha, \alpha'_1)$  and  $\alpha_2 = {}^t(\alpha, \alpha'_2)$  and  $\ell_0^{(i)} = 0$ .

Similarly, assume that there exist  $\beta \in \mathbb{R}^{\delta'}$ , where  $0 \leq \delta' \leq \min(d_1, d_2)$  and a function  $R_\beta$  such as  $N_{\beta_1}^{(1)} = R_\beta + m_{\beta'_1}^{(1)}$ ,  $N_{\beta_2}^{(2)} = R_\beta + m_{\beta'_2}^{(2)}$  with  $\beta_1 = {}^t(\beta, \beta'_1)$  and  $\beta_2 = {}^t(\beta, \beta'_2)$  and  $m_0^{(i)} = 0$ .

Consider now  $\theta = {}^t(\alpha, \alpha'_1, \alpha'_2, \beta, \beta'_1, \beta'_2) \in \mathbb{R}^d$  (and therefore  $\max(d_1, d_2) \leq d \leq d_1 + d_2$ ),  $f_\theta = h_\alpha + \ell_{\alpha'_1}^{(1)} + \ell_{\alpha'_2}^{(2)}$  and  $M_\theta = R_\beta + m_{\beta'_1}^{(1)} + m_{\beta'_2}^{(2)}$ . Then if  $X \in \mathcal{AC}(M_\theta, f_\theta)$ , for any  $t \in \mathbb{Z}$ ,

$$X_t = (R_\beta((X_{t-k})_{k \geq 1}) + m_{\beta'_1}^{(1)}((X_{t-k})_{k \geq 1}) + m_{\beta'_2}^{(2)}((X_{t-k})_{k \geq 1})) \xi_t \\ + (h_\alpha((X_{t-k})_{k \geq 1}) + \ell_{\alpha'_1}^{(1)}((X_{t-k})_{k \geq 1}) + \ell_{\alpha'_2}^{(2)}((X_{t-k})_{k \geq 1})).$$

Then, for  $\alpha'_2 = \beta'_2 = 0$ ,  $X \in \mathcal{AC}(M_{\theta_1}^{(1)}, f_{\theta_1}^{(1)})$  and for  $\alpha'_1 = \beta'_1 = 0$ ,  $X \in \mathcal{AC}(M_{\theta_2}^{(2)}, f_{\theta_2}^{(2)})$ .  $\square$

In the sequel, some lemmas are stated and their proofs are given.

**Lemma 2.1.** *Let  $X \in \mathcal{AC}(M_\theta, f_\theta)$  (or  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$ ) and  $\Theta \subseteq \Theta(r)$  (or  $\Theta \subseteq \widetilde{\Theta}(r)$ ) with  $r \geq 2$ . Assume that the assumptions  $D(\Theta)$  and  $K(\Theta)$  (or  $\widetilde{K}(\Theta)$ ) hold. Then:*

$$\frac{1}{\kappa_n} \|\widehat{L}_n(\theta) - L_n(\theta)\|_\Theta \xrightarrow[n \rightarrow +\infty]{a.s.} 0. \quad (2.7.1)$$

*Proof.* We have  $|\widehat{L}_n(\theta) - L_n(\theta)| \leq \sum_{t=1}^n |\widehat{q}_t(\theta) - q_t(\theta)|$ . Then,

$$\frac{1}{\kappa_n} \|\widehat{L}_n(\theta) - L_n(\theta)\|_\Theta \leq \frac{1}{\kappa_n} \sum_{t=1}^n \|\widehat{q}_t(\theta) - q_t(\theta)\|_\Theta.$$

By Corollary 1 of Kounias and Weng (1969), with  $r \leq 3$ , (2.7.1) is established when:

$$\sum_{k \geq 1} \left(\frac{1}{\kappa_k}\right)^{r/3} \mathbb{E}(\|\widehat{q}_k(\theta) - q_k(\theta)\|_\Theta^{r/3}) < \infty. \quad (2.7.2)$$

With  $r \geq 3$ , and under the assumptions, we first recall some results already obtained in Bardet and Wintenberger (2009): for any  $t \in \mathbb{Z}$ ,

$$\begin{aligned} & \bullet \quad \mathbb{E}[|X_t|^r + \|f_\theta^t\|_\Theta^r + \|\widehat{f}_\theta^t\|_\Theta^r + \|M_\theta^t\|_\Theta^r + \|\widehat{M}_\theta^t\|_\Theta^r + \|H_\theta^t\|_\Theta^{r/2} + \|\widehat{H}_\theta^t\|_\Theta^{r/2}] < \infty \quad (2.7.3) \\ & \bullet \quad \begin{cases} \mathbb{E}[\|f_\theta^t - \widehat{f}_\theta^t\|_\Theta^r] \leq C \left( \sum_{j \geq t} \alpha_j(f_\theta, \Theta) \right)^r \\ \mathbb{E}[\|M_\theta^t - \widehat{M}_\theta^t\|_\Theta^r] \leq C \left( \sum_{j \geq t} \alpha_j(M_\theta, \Theta) \right)^r \\ \mathbb{E}[\|H_\theta^t - \widehat{H}_\theta^t\|_\Theta^{r/2}] \leq C \left( \min \left\{ \sum_{j \geq t} \alpha_j(M_\theta, \Theta), \sum_{j \geq t} \alpha_j(H_\theta, \Theta) \right\} \right)^{r/2}. \end{cases} \quad (2.7.4) \end{aligned}$$

For any  $\theta \in \Theta$ , we have:

$$\begin{aligned}
|\hat{q}_t(\theta) - q_t(\theta)| &= \left| \frac{(X_t - \hat{f}_\theta^t)^2}{\hat{H}_\theta^t} + \log(\hat{H}_\theta^t) - \frac{(X_t - f_\theta^t)^2}{H_\theta^t} - \log(H_\theta^t) \right| \\
&\leq (H_\theta^t \hat{H}_\theta^t)^{-1} |H_\theta^t (X_t - \hat{f}_\theta^t)^2 - \hat{H}_\theta^t (X_t - f_\theta^t)^2| + |\log(\hat{H}_\theta^t) - \log(H_\theta^t)| \\
&\leq (H_\theta^t \hat{H}_\theta^t)^{-1} |(H_\theta^t - \hat{H}_\theta^t)(X_t - f_\theta^t)^2 - H_\theta^t (X_t - f_\theta^t)^2 + H_\theta^t (X_t - \hat{f}_\theta^t)^2| \\
&\quad + |\log(\hat{H}_\theta^t) - \log(H_\theta^t)| \\
&\leq \underline{h}^{-3/2} (|X_t|^2 + 2|X_t|f_\theta^t + |f_\theta^t|^2) |M_\theta^t - \widehat{M}_\theta^t| + \underline{h}^{-1} (2|X_t| + |f_\theta^t| + |\hat{f}_\theta^t|) |f_\theta^t - \hat{f}_\theta^t| \\
&\quad + 2 |\log(\widehat{M}_\theta^t) - \log(M_\theta^t)| \\
&\leq \underline{h}^{-3/2} (|X_t|^2 + 2|X_t| \times \|f_\theta^t\|_\Theta + \|f_\theta^t\|_\Theta^2) \|M_\theta^t - \widehat{M}_\theta^t\|_\Theta \\
&\quad + \underline{h}^{-1} (2|X_t| + \|f_\theta^t\|_\Theta + \|\hat{f}_\theta^t\|_\Theta) \|f_\theta^t - \hat{f}_\theta^t\|_\Theta + 2 \underline{h}^{-1/2} \|\widehat{M}_\theta^t - M_\theta^t\|_\Theta.
\end{aligned}$$

1/ If  $X \subset \mathcal{AC}(M_\theta, f_\theta)$ , we deduce

$$\begin{aligned}
\mathbb{E}[\|\hat{q}_t(\theta) - q_t(\theta)\|_\Theta^{r/3}] &\leq C \left( \mathbb{E}[(\|X_t + f_\theta^t\|_\Theta^2 + 1)^{r/3} \|M_\theta^t - \widehat{M}_\theta^t\|_\Theta^{r/3}] \right. \\
&\quad \left. + \mathbb{E}[(2|X_t| + \|f_\theta^t\|_\Theta + \|\hat{f}_\theta^t\|_\Theta)^{r/3} \|f_\theta^t - \hat{f}_\theta^t\|_\Theta^{r/3}] \right). \quad (2.7.5)
\end{aligned}$$

Then, by Hölder's inequality and (2.7.3) we have:

$$\begin{aligned}
&\mathbb{E}[(\|X_t + f_\theta^t\|_\Theta^2 + 1)^{r/3} \|M_\theta^t - \widehat{M}_\theta^t\|_\Theta^{r/3}] \\
&\leq \left( \mathbb{E}[\|X_t + f_\theta^t + 1\|_\Theta^r] \right)^{2/3} \left( \mathbb{E}[\|M_\theta^t - \widehat{M}_\theta^t\|_\Theta^r] \right)^{1/3} \leq C \left( \mathbb{E}[\|M_\theta^t - \widehat{M}_\theta^t\|_\Theta^r] \right)^{1/3}. \quad (2.7.6)
\end{aligned}$$

Again with Hölder's inequality and (2.7.3),

$$\mathbb{E}[(2|X_t| + \|f_\theta^t\|_\Theta + \|\hat{f}_\theta^t\|_\Theta) \|f_\theta^t - \hat{f}_\theta^t\|_\Theta^{r/3}] \leq C \left( \mathbb{E}[\|f_\theta^t - \hat{f}_\theta^t\|_\Theta^r] \right)^{1/3}. \quad (2.7.7)$$

Therefore, from (2.7.6), (2.7.7) and (2.7.4), there exists a constant  $C$  such that

$$\mathbb{E}[\|\hat{q}_t(\theta) - q_t(\theta)\|_\Theta^{r/3}] \leq C \left( \sum_{j \geq t} \alpha_j(f_\theta, \Theta) + \sum_{j \geq t} \alpha_j(M_\theta, \Theta) \right)^{r/3}. \quad (2.7.8)$$

Hence,

$$\sum_{k \geq 1} \left( \frac{1}{\kappa_k} \right)^{r/3} \mathbb{E}[\|\hat{q}_k(\theta) - q_k(\theta)\|_\Theta^{r/3}] \leq C \sum_{k \geq 1} \left( \frac{1}{\kappa_k} \right)^{r/3} \left( \sum_{j \geq k} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) \right)^{r/3},$$

which is finite by assumption  $K(\Theta)$ , and this achieves the proof.

2/ If  $X \subset \widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$  and using Corollary 1 of [Kounias and Weng \(1969\)](#), with  $r \leq 4$ , (2.7.1) is established when:

$$\sum_{k \geq 1} \left( \frac{1}{\kappa_k} \right)^{r/4} \mathbb{E}[\|\hat{q}_k(\theta) - q_k(\theta)\|_\Theta^{r/4}] < \infty. \quad (2.7.9)$$

By proceeding as in the previous case, we deduce

$$|\hat{q}_t(\theta) - q_t(\theta)| \leq \underline{h}^{-2} |X_t|^2 \|H_\theta^t - \widehat{H}_\theta^t\|_\Theta + \underline{h}^{-1} \|\widehat{H}_\theta^t - H_\theta^t\|_\Theta.$$

In addition, we deduce that there exists a constant  $C$  such that

$$\mathbb{E}[\|(\hat{q}_t(\theta) - q_t(\theta))\|_{\Theta}^{r/4}] \leq C \left( \sum_{j \geq t} \alpha_j(H_\theta, \Theta) \right)^{r/4}. \quad (2.7.10)$$

□

**Lemma 2.2.** *Let  $X \in \mathcal{AC}(M_\theta, f_\theta)$  (or  $\tilde{\mathcal{AC}}(\tilde{H}_\theta)$ ) and  $\Theta \subseteq \Theta(r)$  (or  $\Theta \subseteq \tilde{\Theta}(r)$ ) with  $r \geq 2$ . Assume that the assumptions  $D(\Theta)$  and  $K(\Theta)$  (or  $\tilde{K}(\Theta)$ ) hold. Then:*

$$\frac{1}{\kappa_n} \left\| \frac{\partial \hat{L}_n(\theta)}{\partial \theta} - \frac{\partial L_n(\theta)}{\partial \theta} \right\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0. \quad (2.7.11)$$

*Proof.* We will go along similar lines as in the proof of Lemma 2.1. We have:

$$\frac{1}{\kappa_n} \left\| \frac{\partial \hat{L}_n(\theta)}{\partial \theta} - \frac{\partial L_n(\theta)}{\partial \theta} \right\|_{\Theta} \leq \frac{1}{\kappa_n} \sum_{t=1}^n \left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta}.$$

Using again Corollary 1 of Kounias and Weng (1969), it is sufficient to prove for  $r \leq 3$  that

$$\sum_{k \geq 1} \left( \frac{1}{\kappa_k} \right)^{r/3} \mathbb{E} \left[ \left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta}^{r/3} \right] < \infty. \quad (2.7.12)$$

For any  $\theta \in \Theta$ , with  $H_\theta = M_\theta^2$ , the first partial derivatives of  $q_t(\theta)$  are

$$\begin{aligned} \frac{\partial q_t(\theta)}{\partial \theta_i} &= \frac{-2(X_t - f_\theta^t) \frac{\partial f_\theta^t}{\partial \theta_i}}{H_\theta^t} - \frac{(X_t - f_\theta^t)^2 \frac{\partial H_\theta^t}{\partial \theta_i}}{(H_\theta^t)^2} + \frac{1}{H_\theta^t} \frac{\partial H_\theta^t}{\partial \theta_i} \\ &= -2(H_\theta^t)^{-1}(X_t - f_\theta^t) \frac{\partial f_\theta^t}{\partial \theta_i} + (X_t - f_\theta^t)^2 \frac{\partial (H_\theta^t)^{-1}}{\partial \theta_i} + (H_\theta^t)^{-1} \frac{\partial H_\theta^t}{\partial \theta_i}, \end{aligned}$$

for  $i = 1, \dots, d$ . Hence,

$$\begin{aligned} \left| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right| &\leq 2 \left| (h_\theta^t)^{-1}(X_t - f_\theta^t) \frac{\partial f_\theta^t}{\partial \theta_i} - (\hat{h}_\theta^t)^{-1}(X_t - \hat{f}_\theta^t) \frac{\partial \hat{f}_\theta^t}{\partial \theta_i} \right| \\ &\quad + \left| (X_t - \hat{f}_\theta^t)^2 \frac{\partial (\hat{H}_\theta^t)^{-1}}{\partial \theta_i} - (X_t - f_\theta^t)^2 \frac{\partial (H_\theta^t)^{-1}}{\partial \theta_i} \right| + \left| (\hat{H}_\theta^t)^{-1} \frac{\partial \hat{H}_\theta^t}{\partial \theta_i} - (H_\theta^t)^{-1} \frac{\partial H_\theta^t}{\partial \theta_i} \right|. \end{aligned}$$

Then, using  $|a_1 b_1 c_1 - a_2 b_2 c_2| \leq |a_1 - a_2| |b_2| |c_2| + |a_1| |b_1 - b_2| |c_2| + |a_1| |b_1| |c_1 - c_2|$  for any  $a_1, a_2, b_1, b_2, c_1, c_2$  in  $\mathbb{R}$ , we obtain

$$\begin{aligned} \left| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right| &\leq 2 \left( |(H_\theta^t)^{-1} - (\hat{H}_\theta^t)^{-1}| \times |X_t - \hat{f}_\theta^t| \left| \frac{\partial f_\theta^t}{\partial \theta_i} \right| + |(H_\theta^t)^{-1}| \times |\hat{f}_\theta^t - f_\theta^t| \left| \frac{\partial f_\theta^t}{\partial \theta_i} \right| \right. \\ &\quad \left. + |(H_\theta^t)^{-1}| \times |X_t - f_\theta^t| \left| \frac{\partial f_\theta^t}{\partial \theta_i} - \frac{\partial \hat{f}_\theta^t}{\partial \theta_i} \right| + |X_t - \hat{f}_\theta^t|^2 \left| \frac{\partial (\hat{H}_\theta^t)^{-1}}{\partial \theta_i} - \frac{\partial (H_\theta^t)^{-1}}{\partial \theta_i} \right| \right. \\ &\quad \left. + 2 \left| \frac{\partial (H_\theta^t)^{-1}}{\partial \theta_i} \right| |X_t| |f_\theta^t - \hat{f}_\theta^t| + |(\hat{H}_\theta^t)^{-1}| \left| \frac{\partial \hat{H}_\theta^t}{\partial \theta_i} - \frac{\partial H_\theta^t}{\partial \theta_i} \right| + \left| \frac{\partial H_\theta^t}{\partial \theta_i} \right| |(\hat{H}_\theta^t)^{-1} - (H_\theta^t)^{-1}| \right). \end{aligned}$$

Thus,

$$\begin{aligned}
\left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta} &\leq 2h^{-1} \left( \|\hat{f}_{\theta}^t - f_{\theta}^t\|_{\Theta} \left\| \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta} + \|X_t - f_{\theta}^t\|_{\Theta} \left\| \frac{\partial f_{\theta}^t}{\partial \theta_i} - \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta} \right) \\
&\quad + 2 \left\| (H_{\theta}^t)^{-1} - (\hat{H}_{\theta}^t)^{-1} \right\|_{\Theta} \|X_t - \hat{f}_{\theta}^t\|_{\Theta} \left\| \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta} + \|X_t - \hat{f}_{\theta}^t\|_{\Theta}^2 \left\| \frac{\partial (\hat{H}_{\theta}^t)^{-1}}{\partial \theta_i} - \frac{\partial (H_{\theta}^t)^{-1}}{\partial \theta_i} \right\|_{\Theta} \\
&\quad + 2 \|X_t\|_{\Theta} \|f_{\theta}^t - \hat{f}_{\theta}^t\|_{\Theta} \left\| \frac{\partial (H_{\theta}^t)^{-1}}{\partial \theta_i} \right\|_{\Theta} + \|(\hat{H}_{\theta}^t)^{-1}\|_{\Theta} \left\| \frac{\partial \hat{H}_{\theta}^t}{\partial \theta_i} - \frac{\partial H_{\theta}^t}{\partial \theta_i} \right\|_{\Theta} \\
&\quad + \|(\hat{H}_{\theta}^t)^{-1} - (H_{\theta}^t)^{-1}\|_{\Theta} \left\| \frac{\partial H_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}.
\end{aligned}$$

Using again the results of [Bardet and Wintenberger \(2009\)](#), we know that:

$$\begin{aligned}
&\bullet \quad \mathbb{E} \left[ \left\| \frac{\partial f_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^r + \left\| \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^r + \left\| \frac{\partial M_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^r + \left\| \frac{\partial \hat{M}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^r + \left\| \frac{\partial H_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^{r/2} + \left\| \frac{\partial (H_{\theta}^t)^{-1}}{\partial \theta_i} \right\|_{\Theta}^r \right] < \infty \quad (2.7.13) \\
&\bullet \quad \begin{cases} \mathbb{E} \left[ \left\| (H_{\theta}^t)^{-1} - (\hat{H}_{\theta}^t)^{-1} \right\|_{\Theta}^r \right] \leq C \left( \sum_{j \geq t} \alpha_j(M_{\theta}, \Theta) \right)^r \\ \mathbb{E} \left[ \left\| \frac{\partial f_{\theta}^t}{\partial \theta_i} - \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^r \right] \leq C \left( \sum_{j \geq t} \alpha_j(\partial f_{\theta}, \Theta) \right)^r \\ \mathbb{E} \left[ \left\| \frac{\partial H_{\theta}^t}{\partial \theta_i} - \frac{\partial \hat{H}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^{r/2} \right] \leq C \left( \sum_{j \geq t} (\alpha_j(M_{\theta}, \Theta) + \alpha_j(\partial M_{\theta}, \Theta)) \right)^{r/2} \\ \mathbb{E} \left[ \left\| \frac{\partial (H_{\theta}^t)^{-1}}{\partial \theta_i} - \frac{\partial (\hat{H}_{\theta}^t)^{-1}}{\partial \theta_i} \right\|_{\Theta}^{r/2} \right] \leq C \left( \sum_{j \geq t} (\alpha_j(M_{\theta}, \Theta) + \alpha_j(\partial M_{\theta}, \Theta)) \right)^{r/2} \end{cases} \quad (2.7.14)
\end{aligned}$$

1. If  $X \subset \mathcal{AC}(M_{\theta}, f_{\theta})$ , we deduce from the Hölder's Inequality that,

$$\begin{aligned}
\mathbb{E} \left[ \left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta}^{r/3} \right] &\leq C \left[ (\mathbb{E} [\|\hat{f}_{\theta}^t - f_{\theta}^t\|_{\Theta}^r])^{1/3} (\mathbb{E} [\left\| \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^{r/2}])^{2/3} \right. \\
&\quad + (\mathbb{E} [\|X_t - f_{\theta}^t\|_{\Theta}^{2r/3}])^{1/2} (\mathbb{E} [\left\| \frac{\partial f_{\theta}^t}{\partial \theta_i} - \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^r])^{1/3} \\
&\quad + (\mathbb{E} [\|(H_{\theta}^t)^{-1} - (\hat{H}_{\theta}^t)^{-1}\|_{\Theta}^r])^{1/3} (\mathbb{E} [\|X_t - \hat{f}_{\theta}^t\|_{\Theta}^r] \mathbb{E} [\left\| \frac{\partial \hat{f}_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^r])^{1/3} \\
&\quad + (\mathbb{E} [\|X_t - \hat{f}_{\theta}^t\|_{\Theta}^r])^{1/3} (\mathbb{E} [\left\| \frac{\partial (\hat{H}_{\theta}^t)^{-1}}{\partial \theta_i} - \frac{\partial (H_{\theta}^t)^{-1}}{\partial \theta_i} \right\|_{\Theta}^{r/2}])^{2/3} \\
&\quad + (\mathbb{E} [\left\| \frac{\partial (H_{\theta}^t)^{-1}}{\partial \theta_i} \right\|_{\Theta}^r])^{1/3} (\mathbb{E} [|X_t|^r] \mathbb{E} [\|f_{\theta}^t - \hat{f}_{\theta}^t\|_{\Theta}^r])^{1/3} \\
&\quad \left. + (\mathbb{E} [\left\| \frac{\partial \hat{H}_{\theta}^t}{\partial \theta_i} - \frac{\partial H_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^{r/3}]) + (\mathbb{E} [\left\| \frac{\partial H_{\theta}^t}{\partial \theta_i} \right\|_{\Theta}^{r/2}])^{2/3} (\mathbb{E} [\|(H_{\theta}^t)^{-1} - (\hat{H}_{\theta}^t)^{-1}\|_{\Theta}^r])^{1/3} \right].
\end{aligned}$$

Using (2.7.13) and (2.7.14), we deduce

$$\begin{aligned}
\mathbb{E} \left[ \left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta}^{r/3} \right] &\leq C \left( \sum_{j \geq t} \alpha_j(f_{\theta}, \Theta) + \alpha_j(M_{\theta}, \Theta) \right. \\
&\quad \left. + \alpha_j(\partial f_{\theta}, \Theta) + \alpha_j(\partial M_{\theta}, \Theta) \right)^{r/3}.
\end{aligned}$$

Therefore,

$$\begin{aligned} \sum_{k \geq 1} \frac{1}{\kappa_k^{r/3}} \mathbb{E} \left[ \left\| \frac{\partial \hat{q}_k(\theta)}{\partial \theta_i} - \frac{\partial q_k(\theta)}{\partial \theta_i} \right\|_{\Theta}^{r/3} \right] \\ \leq C \sum_{k \geq 1} \frac{1}{\kappa_k^{r/3}} \left( \sum_{j \geq t} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) + \alpha_j(\partial f_\theta, \Theta) + \alpha_j(\partial M_\theta, \Theta) \right)^{r/3}. \end{aligned}$$

We conclude the proof of (2.7.12) from assumption  $K(\Theta)$ .

2. If  $X \subset \widetilde{\mathcal{AC}}(\tilde{H}_\theta)$ , we deduce

$$\begin{aligned} \left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta} \leq |X_t|^2 \left\| \frac{\partial(\hat{H}_\theta^t)^{-1}}{\partial \theta_i} - \frac{\partial(H_\theta^t)^{-1}}{\partial \theta_i} \right\|_{\Theta} \\ + \underline{h}^{-1} \left\| \frac{\partial \hat{H}_\theta^t}{\partial \theta_i} - \frac{\partial H_\theta^t}{\partial \theta_i} \right\|_{\Theta} + \|(\hat{H}_\theta^t)^{-1} - (H_\theta^t)^{-1}\|_{\Theta} \left\| \frac{\partial H_\theta^t}{\partial \theta_i} \right\|_{\Theta}. \end{aligned}$$

As a consequence,

$$\begin{aligned} \mathbb{E} \left[ \left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta}^{r/4} \right] \leq \left( \mathbb{E}[|X_t|^r] \mathbb{E} \left[ \left\| \frac{\partial(\hat{H}_\theta^t)^{-1}}{\partial \theta_i} - \frac{\partial(H_\theta^t)^{-1}}{\partial \theta_i} \right\|_{\Theta}^{r/2} \right] \right)^{1/2} \\ + \underline{h}^{-r/4} \mathbb{E} \left[ \left\| \frac{\partial \hat{H}_\theta^t}{\partial \theta_i} - \frac{\partial H_\theta^t}{\partial \theta_i} \right\|_{\Theta}^{r/4} \right] + \left( \mathbb{E}[\|(\hat{H}_\theta^t)^{-1} - (H_\theta^t)^{-1}\|_{\Theta}^{r/2}] \mathbb{E} \left[ \left\| \frac{\partial H_\theta^t}{\partial \theta_i} \right\|_{\Theta}^{r/2} \right] \right)^{1/2}, \end{aligned}$$

implying

$$\mathbb{E} \left[ \left\| \frac{\partial \hat{q}_t(\theta)}{\partial \theta_i} - \frac{\partial q_t(\theta)}{\partial \theta_i} \right\|_{\Theta}^{r/4} \right] \leq C \left( \sum_{j \geq t} \alpha_j(H_\theta, \Theta) + \alpha_j(\partial H_\theta, \Theta) \right)^{r/4},$$

which achieves the proof, according to Corollary 1 of [Kounias and Weng \(1969\)](#).  $\square$

**Lemma 2.3.** *Under the assumptions of Theorem 2.1 and if a model  $m \in \mathcal{M}$  is such that  $\theta^* \in \Theta(m)$ , then:*

$$\frac{1}{\kappa_n} |\hat{L}_n(\hat{\theta}(m)) - \hat{L}_n(\hat{\theta}(m^*))| = o_P(1). \quad (2.7.15)$$

*Proof.* We have:

$$\begin{aligned} \frac{1}{\kappa_n} |\hat{L}_n(\hat{\theta}(m)) - \hat{L}_n(\hat{\theta}(m^*))| &= \frac{1}{\kappa_n} |\hat{L}_n(\hat{\theta}(m)) - L_n(\hat{\theta}(m)) + L_n(\hat{\theta}(m)) - L_n(\hat{\theta}(m^*)) \\ &\quad + L_n(\hat{\theta}(m^*)) - \hat{L}_n(\hat{\theta}(m^*))| \\ &\leq \frac{2}{\kappa_n} \|\hat{L}_n(\theta) - L_n(\theta)\|_{\Theta} + \frac{1}{\kappa_n} |L_n(\hat{\theta}(m)) - L_n(\hat{\theta}(m^*))|. \end{aligned}$$

According to Lemma 2.1,  $\frac{1}{\kappa_n} \|\hat{L}_n(\theta) - L_n(\theta)\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$ . The proof will be achieved if we can show that

$$\frac{1}{\kappa_n} |L_n(\hat{\theta}(m)) - L_n(\theta^*)| = o_P(1). \quad (2.7.16)$$



Since

$$\frac{1}{\kappa_n} \left| L_n(\hat{\theta}(m)) - L_n(\hat{\theta}(m^*)) \right| \leq \frac{1}{\kappa_n} \left| L_n(\hat{\theta}(m)) - L_n(\theta^*) \right| + \frac{1}{\kappa_n} \left| L_n(\hat{\theta}(m^*)) - L_n(\theta^*) \right|.$$

Applying a second order Taylor expansion of  $L_n$  around  $\hat{\theta}(m)$  for  $n$  sufficiently large such that  $\bar{\theta}(m) \in \Theta(m)$  which are between  $\hat{\theta}(m)$  and  $\theta^*$ , yields:

$$\begin{aligned} \frac{1}{\kappa_n} (L_n(\hat{\theta}(m)) - L_n(\theta^*)) &= \\ \frac{1}{\kappa_n} (\hat{\theta}(m) - \theta^*) \frac{\partial L_n(\hat{\theta}(m))}{\partial \theta} &+ \frac{1}{2\kappa_n} (\hat{\theta}(m) - \theta^*)' \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta^2} (\hat{\theta}(m) - \theta^*). \end{aligned} \quad (2.7.17)$$

Let us deal first with the first term on the right hand side of last equality:

$$\frac{1}{\kappa_n} (\hat{\theta}(m) - \theta^*) \frac{\partial L_n(\hat{\theta}(m))}{\partial \theta} = \frac{1}{\kappa_n} \sqrt{n} (\hat{\theta}(m) - \theta^*) \frac{1}{\sqrt{n}} \frac{\partial L_n(\hat{\theta}(m))}{\partial \theta}.$$

Since  $\frac{1}{\kappa_n} = o(1)$  and from [Bardet and Wintenberger \(2009\)](#) we have  $\sqrt{n}(\hat{\theta}(m) - \theta^*) = O_P(1)$  and  $\frac{1}{\sqrt{n}} \frac{\partial L_n(\hat{\theta}(m))}{\partial \theta} = o_P(1)$ , it follows that:

$$\frac{1}{\kappa_n} (\hat{\theta}(m) - \theta^*) \frac{\partial L_n(\hat{\theta}(m))}{\partial \theta} = o_P(1). \quad (2.7.18)$$

On the other hand, for the second term of the right hand side of equality (2.7.17), let us note that, we have from [Bardet and Wintenberger \(2009\)](#):

- $\sqrt{n} (\hat{\theta}(m) - \theta^*) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{A}_{\theta^*, m}$  a Gaussian random variable from Theorem 2 of [Bardet and Wintenberger \(2009\)](#).
- $-\frac{2}{n} \left( \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta_i \partial \theta_j} \right)_{i,j \in m} \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta^*, m)$  since  $\hat{\theta}(m) \xrightarrow[n \rightarrow +\infty]{a.s.} \theta^*$  and using the assumption  $\text{Var}(\Theta)$  insuring that the matrix  $F(\theta^*, m)$  exists and is definite positive (see [Bardet and Wintenberger \(2009\)](#)).

Hence,

$$\begin{aligned} (\hat{\theta}(m) - \theta^*)' \left( \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta_i \partial \theta_j} \right)_{i,j \in m} (\hat{\theta}(m) - \theta^*) &= \\ = \frac{-1}{2} \sqrt{n} (\hat{\theta}(m) - \theta^*)' (F(\theta^*, m) + o_P(1)) \sqrt{n} (\hat{\theta}(m) - \theta^*) & \\ \xrightarrow[n \rightarrow \infty]{\mathcal{P}} \frac{-1}{2} \mathcal{A}'_{\theta^*, m} F(\theta^*, m) \mathcal{A}_{\theta^*, m}. \end{aligned}$$

We deduce that

$$\begin{aligned} (\hat{\theta}(m) - \theta^*)' \left( \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta_i \partial \theta_j} \right)_{i,j \in m} (\hat{\theta}(m) - \theta^*) &= O_P(1) \\ \implies \frac{1}{\kappa_n} (\hat{\theta}(m) - \theta^*)' \left( \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta_i \partial \theta_j} \right)_{i,j \in m} (\hat{\theta}(m) - \theta^*) &= o_P(1). \end{aligned} \quad (2.7.19)$$

Thus, (2.7.16) follows from (2.7.17), (2.7.18) and (2.7.19); which completes the proof of Lemma 2.3.  $\square$

### 2.7.1 Misspecified model

When a model  $m$  is misspecified ( $\theta^* \notin \Theta(m)$ ), we will show that  $\mathbb{P}(m^* \not\subseteq \hat{m}) \xrightarrow[n \rightarrow \infty]{} 0$  by following the key idea of similar proof in [Sin and White \(1996\)](#) by defining the "best" parameter  $\theta^*(m) \in \Theta(m)$  which will play the role of  $\theta^*$  in cases of "true" or overfitted model. For model  $m \in \mathcal{M}$ , let define

$$\theta^*(m) := \operatorname{argmax}_{\theta \in \Theta(m)} L(\theta) \quad \text{with} \quad L(\theta) := -\frac{1}{2} \mathbb{E}[q_0(\theta)]. \quad (2.7.20)$$

Given that the function  $L$  is continuous and the parameter set  $\Theta(m)$  is assumed compact, the parameter  $\theta^*(m)$  exists and is unique by virtue of the identifiability condition  $Id(\Theta)$ .

It is worth noting, since  $L(\theta)$  has a unique maximum reached at  $\theta^*$  (see [Bardet and Wintenberger \(2009\)](#)), and along with the fact that  $\theta^* \in \Theta(m)$ , it follows that  $\theta^*(m) = \theta^*$  when  $m$  is the true model or an overfitted one.

Let us show that even in the presence of misspecification, the QMLE still remains consistent but for  $\theta^*(m)$ . This important result will allow us to show that our model selection procedure can not choose a misspecified model.

**Proposition 2.2.** *Let  $X \in \mathcal{AC}(M_{\theta^*}, f_{\theta^*})$  (or  $\widetilde{\mathcal{AC}}(\tilde{H}_{\theta^*})$ ) and  $\Theta \subseteq \Theta(r)$  (or  $\Theta \subseteq \tilde{\Theta}(r)$ ) with  $r \geq 2$ . Under the assumptions  $Id(\Theta)$ ,  $D(\Theta)$  and  $K(\Theta)$ , it holds*

$$\left\| \frac{1}{n} L_n(\theta) - L(\theta) \right\|_{\Theta(m)} \xrightarrow[n \rightarrow +\infty]{a.s.} 0 \quad \text{and} \quad (2.7.21)$$

$$\hat{\theta}(m) \xrightarrow[n \rightarrow +\infty]{a.s.} \theta^*(m). \quad (2.7.22)$$

*Proof.* The proof of (2.7.21) follows from a consequence of uniform strong law of large numbers for stationary ergodic sequence (see the proof of Theorem 1 in [Bardet and Wintenberger \(2009\)](#)). The second result holds by applying (2.7.21) and Lemma 2.1.  $\square$

### 2.7.2 Proof of Theorem 2.1

Before diving into the proof, remark first that:

$$\mathbb{P}(\hat{m} = m^*) = 1 - \mathbb{P}(m^* \subset \hat{m}) - \mathbb{P}(m^* \not\subseteq \hat{m}). \quad (2.7.23)$$

As we point out in Subsection 2.2.1, the proof is divided into two parts; the first part shows that our selection criterion chooses an overfitted model with probability decreasing to zero while the second part shows a similar behavior for the probability of selecting a misspecified model.

*Proof.* 1. Since  $\mathcal{M}$  is finite, let  $m_0 \in \mathcal{M}$  such as  $\hat{m} = m_0$  and  $m^* \subset m_0$ , (i.e an overfitted model was selected, but let show that this cannot happen). Let compute  $\mathbb{P}(\hat{C}(m_0) \subseteq \hat{C}(m^*))$  for large  $n$ .

We have:

$$\begin{aligned}
\mathbb{P}(\widehat{C}(m_0) \leq \widehat{C}(m^*)) &= \mathbb{P}\left(-2\widehat{L}_n(\widehat{\theta}(m_0)) + |m_0|\kappa_n \leq -2\widehat{L}_n(\widehat{\theta}(m^*)) + |m^*|\kappa_n\right) \\
&= \mathbb{P}\left(-2\widehat{L}_n(\widehat{\theta}(m_0)) + 2\widehat{L}_n(\widehat{\theta}(m^*)) \leq \kappa_n(|m^*| - |m_0|)\right) \\
&= \mathbb{P}\left(\frac{1}{\kappa_n}(\widehat{L}_n(\widehat{\theta}(m^*)) - \widehat{L}_n(\widehat{\theta}(m_0))) \leq \frac{(|m^*| - |m_0|)}{2}\right) \\
&\xrightarrow{n \rightarrow \infty} 0
\end{aligned}$$

by virtue of Lemma 2.3 and because  $|m_0| - |m^*| \geq 1$ . This shows,  $\widehat{C}(m_0) > \widehat{C}(m^*)$  with probability going to 1, i.e.  $\widehat{C}(\widehat{m}) > \widehat{C}(m^*)$ . We get a contradiction along with definition of  $\widehat{m}$  (2.2.4), and then the selection criteria can not choose  $\widehat{m}$  which stricly contains the true model, thus

$$\mathbb{P}(m^* \subset \widehat{m}) \xrightarrow{n \rightarrow \infty} 0.$$

2. Since  $\mathcal{M}$  is finite, let  $m_0 \in \mathcal{M}$  such as  $\widehat{m} = m_0$  and  $m^* \not\subseteq \widehat{m}$ . Let compute  $n^{-1}[\widehat{C}(m_0) - \widehat{C}(m^*)]$  for large  $n$ . First,

$$\begin{aligned}
\frac{1}{n}[\widehat{L}_n(\widehat{\theta}(m^*)) - \widehat{L}_n(\widehat{\theta}(m_0))] &= \frac{1}{n}[L_n(\widehat{\theta}(m^*)) - L_n(\widehat{\theta}(m_0))] + o_{a.s.}(1) \text{ using Lemma 2.1} \\
&= L(\widehat{\theta}(m^*)) - L(\widehat{\theta}(m_0)) + o_{a.s.}(1) \text{ using Proposition 2.2} \\
&= [L(\widehat{\theta}(m^*)) - L(\theta^*)] - [L(\widehat{\theta}(m_0)) - L(\theta^*(m_0))] \\
&\quad + [L(\theta^*) - L(\theta^*(m_0))] + o_{a.s.}(1).
\end{aligned}$$

Since  $L$  is continuous over  $\Theta$ , using continuous mapping theorem and the relation (2.7.22) of Proposition 2.2, it holds for  $n$  large enough

$$L(\widehat{\theta}(m^*)) - L(\theta^*) = o_{a.s.}(1) \quad \text{and} \quad L(\widehat{\theta}(m_0)) - L(\theta^*(m_0)) = o_{a.s.}(1).$$

Hence,

$$\frac{1}{n}[\widehat{L}_n(\widehat{\theta}(m^*)) - \widehat{L}_n(\widehat{\theta}(m_0))] = L(\theta^*) - L(\theta^*(m_0)) + o_{a.s.}(1). \quad (2.7.24)$$

Now, let us show that that  $A(m) := L(\theta^*) - L(\theta^*(m)) > 0$  for  $m = m_0$ .

Let us denote by  $\mathcal{F}_t := \sigma(X_{t-1}, X_{t-2}, \dots)$ . Using conditional expectation, we obtain

$$L(\theta^*) - L(\theta) = \frac{1}{2} \mathbb{E}[\mathbb{E}[q_0(\theta) - q_0(\theta^*) \mid \mathcal{F}_0]]. \quad (2.7.25)$$

But,

$$\begin{aligned}
\mathbb{E}[q_0(\theta) - q_0(\theta^*) \mid \mathcal{F}_0] &= \mathbb{E}\left[\frac{(X_0 - f_\theta^0)^2}{H_\theta^0} + \log(H_\theta^0) - \frac{(X_0 - f_{\theta^*}^0)^2}{H_{\theta^*}^0} - \log(H_{\theta^*}^0) \mid \mathcal{F}_0\right] \\
&= \log\left(\frac{H_\theta^0}{H_{\theta^*}^0}\right) + \frac{\mathbb{E}[(X_0 - f_\theta^0)^2 \mid \mathcal{F}_0]}{H_\theta^0} - \frac{\mathbb{E}[(X_0 - f_{\theta^*}^0)^2 \mid \mathcal{F}_0]}{H_{\theta^*}^0} \\
&= \log\left(\frac{H_\theta^0}{H_{\theta^*}^0}\right) - 1 + \frac{\mathbb{E}[(X_0 - f_{\theta^*}^0 + f_{\theta^*}^0 - f_\theta^0)^2 \mid \mathcal{F}_0]}{H_\theta^0} \\
&= \frac{H_{\theta^*}^0}{H_\theta^0} - \log\left(\frac{H_{\theta^*}^0}{H_\theta^0}\right) - 1 + \frac{(f_{\theta^*}^0 - f_\theta^0)^2}{H_\theta^0}
\end{aligned}$$

Thus from (2.7.25),

$$\begin{aligned} 2 A(m_0) &= \mathbb{E} \left[ \frac{H_{\theta^*}^0}{H_{\theta^*(m_0)}^0} - \log \left( \frac{H_{\theta^*}^0}{H_{\theta^*(m_0)}^0} \right) - 1 + \frac{(f_{\theta^*}^0 - f_{\theta^*(m_0)}^0)^2}{H_{\theta^*(m_0)}^0} \right] \\ &\geq \mathbb{E} \left[ \frac{H_{\theta^*}^0}{H_{\theta^*(m_0)}^0} \right] - \log \left( \mathbb{E} \left[ \frac{H_{\theta^*}^0}{H_{\theta^*(m_0)}^0} \right] \right) - 1 + \mathbb{E} \left[ \frac{(f_{\theta^*}^0 - f_{\theta^*(m_0)}^0)^2}{H_{\theta^*(m_0)}^0} \right] \quad \text{by Jensen Inequality.} \end{aligned}$$

Since  $x - \log(x) - 1 > 0$  for any  $x > 0$ ,  $x \neq 1$  and  $x - \log(x) - 1 = 0$  for  $x = 1$ , we deduce that

- If  $f_{\theta^*}^0 \neq f_{\theta^*(m_0)}^0$ , then  $\mathbb{E} \left[ \left( \frac{f_{\theta^*}^0 - f_{\theta^*(m_0)}^0}{H_{\theta^*(m_0)}^0} \right)^2 \right] > 0$  and  $A(m_0) > 0$ .
- Otherwise, then

$$2 A(m_0) = \mathbb{E} \left[ \frac{H_{\theta^*}^0}{H_{\theta^*(m_0)}^0} - \log \left( \frac{H_{\theta^*}^0}{H_{\theta^*(m_0)}^0} \right) - 1 \right],$$

and from assumption  $Id(\Theta)$ , since  $\theta^* \notin \Theta(m)$  and  $f_{\theta^*}^0 = f_{\theta^*(m_0)}^0$ , we necessarily have  $H_{\theta^*}^0 \neq H_{\theta^*(m_0)}^0$  so that  $\frac{H_{\theta^*}^0}{H_{\theta^*(m_0)}^0} \neq 1$ . Therefore,  $A(m_0) > 0$ .

As a consequence,

$$\frac{\widehat{C}(m_0) - \widehat{C}(m^*)}{n} = A(m_0) + \frac{\kappa_n}{n} (|m_0| - |m^*|) + o_{a.s.}(1). \quad (2.7.26)$$

Moreover, since  $\kappa_n = o(n)$  and all the considered models are finite dimensional, the equality (2.7.26) implies for large  $n$  that  $\widehat{C}(m_0) > \widehat{C}(m^*)$  almost surely. This means that it was possible to select a model  $\widehat{m}$  with  $\widehat{C}(\widehat{m}) > \widehat{C}(m^*)$ , which is impossible according to the definition (2.2.4). Therefore the event  $m^* \not\leq \widehat{m}$  can not happen and then

$$\mathbb{P}(m^* \not\leq \widehat{m}) \xrightarrow{n \rightarrow \infty} 0.$$

Thus we have proved the first and most difficult part of Theorem (2.1). The next lines show the second part which is about the consistency of  $\widehat{\theta}(\widehat{m})$ .

Given  $\epsilon > 0$ , we have :

$$\begin{aligned} \mathbb{P} \left( \|\widehat{\theta}(\widehat{m}) - \theta^*\|_{i \in m^*} > \epsilon \right) &= \mathbb{P} \left( \|\widehat{\theta}(\widehat{m}) - \theta^*\|_{i \in m^*} > \epsilon | \widehat{m} = m^* \right) \mathbb{P}(\widehat{m} = m^*) \\ &\quad + \mathbb{P} \left( \|\widehat{\theta}(\widehat{m}) - \theta^*\|_{i \in m^*} > \epsilon | \widehat{m} \neq m^* \right) \mathbb{P}(\widehat{m} \neq m^*). \end{aligned}$$

From the strong consistency of the QMLE (see New version of Theorem 1 of [Bardet and Wintenberger \(2009\)](#)), the first term of the right hand side of the above equation is asymptotically zero and also the second one under the assumptions of the first part of Theorem 2.1 which gives  $\mathbb{P}(\widehat{m} \neq m^*) \xrightarrow{n \rightarrow \infty} 0$ .

□

### 2.7.3 Proof of Theorem 2.2

*Proof.* For  $x = (x_i)_{1 \leq i \leq d} \in \mathbb{R}^d$ , denote  $F_n(x) = \mathbb{P}\left(\bigcap_{1 \leq i \leq d} \sqrt{n}(\hat{\theta}(\hat{m}) - \theta^*)_i \leq x_i\right)$ .

First, we have:

$$\begin{aligned} F_n(x) &= \mathbb{P}\left(\bigcap_{1 \leq i \leq d} \sqrt{n}(\hat{\theta}(\hat{m}) - \theta^*)_i \leq x_i \mid \hat{m} = m^*\right) \mathbb{P}(\hat{m} = m^*) \\ &\quad + \mathbb{P}\left(\bigcap_{1 \leq i \leq d} \sqrt{n}(\hat{\theta}(\hat{m}) - \theta^*)_i \leq x_i \mid \hat{m} \neq m^*\right) \mathbb{P}(\hat{m} \neq m^*). \end{aligned}$$

Under the assumptions of Theorem 2.1,  $\mathbb{P}(\hat{m} = m^*) \xrightarrow{n \rightarrow \infty} 1$  and  $\mathbb{P}(\hat{m} \neq m^*) \xrightarrow{n \rightarrow \infty} 0$ . Therefore the second term in the right side of the previous equality asymptotically vanishes. For the first term, we can write,

$$\begin{aligned} &\mathbb{P}\left(\bigcap_{1 \leq i \leq d} \sqrt{n}(\hat{\theta}(\hat{m}) - \theta^*)_i \leq x_i \mid \hat{m} = m^*\right) \\ &= \mathbb{P}\left(\left\{\bigcap_{i \in m^*} \sqrt{n}(\hat{\theta}(m^*) - \theta^*)_i \leq x_i\right\} \cap \left\{\bigcap_{i \notin m^*} \sqrt{n}(\hat{\theta}(m^*) - \theta^*)_i \leq x_i\right\}\right). \end{aligned}$$

Since  $\theta(m^*) \in \Theta(m^*)$ ,  $((\hat{\theta}(m^*))_i)_{i \notin m^*} = (\theta_i^*)_{i \notin m^*} = 0$ , for  $(x_i)_{i \notin m^*}$  a family of non negative real numbers we have:

$$\begin{aligned} &\mathbb{P}\left(\left\{\bigcap_{i \in m^*} \sqrt{n}(\hat{\theta}(m^*) - \theta^*)_i \leq x_i\right\} \cap \left\{\bigcap_{i \notin m^*} \sqrt{n}(\hat{\theta}(m^*) - \theta^*)_i \leq x_i\right\}\right) \\ &= \mathbb{P}\left(\bigcap_{i \in m^*} \sqrt{n}(\hat{\theta}(m^*) - \theta^*)_i \leq x_i\right) \\ &\xrightarrow{n \rightarrow \infty} \mathbb{P}\left((F(\theta^*, m^*)^{-1} G(\theta^*, m^*) F(\theta^*, m^*)^{-1})^{-1/2} Z \leq (x_i)_{i \in m^*}\right), \end{aligned}$$

with  $Z$  a standard Gaussian random vector in  $\mathbb{R}^{|m^*|}$  from the central limit theorem in Theorem 2 of [Bardet and Wintenberger \(2009\)](#), and this achieves the proof of 2.3.3 of Theorem 2.2.  $\square$

### 2.7.4 Proof of Theorem 2.3

Consider the following notation: for  $\theta \in \Theta$  and  $m \in \mathcal{M}$ , denote the residuals and quasi-residuals by:

$$\begin{cases} e_t(\theta) &:= (M_\theta^t)^{-1}(X_t - f_\theta^t) & \text{and} & \hat{e}_t(\theta) &:= (\widehat{M}_\theta^t)^{-1}(X_t - \hat{f}_\theta^t) \\ e_t(m) &:= (M_{\hat{\theta}(m)}^t)^{-1}(X_t - f_{\hat{\theta}(m)}^t) & \text{and} & \hat{e}_t(m) &:= (M_{\hat{\theta}(m)}^t)^{-1}(X_t - \hat{f}_{\hat{\theta}(m)}^t) \end{cases}.$$

For  $k \in \{0, 1, \dots, n-1\}$ ,  $\theta \in \Theta$  and  $m \in \mathcal{M}$ , define also the adjusted lag- $k$  covariograms and correlograms of the squared (standardized) residual by:

$$\begin{cases} \gamma_k(\theta) &:= \frac{1}{n} \sum_{t=1}^{n-k} (e_t^2(\theta) - 1)(e_{t+k}^2(\theta) - 1); \hat{\gamma}_k(\theta) &:= \frac{1}{n} \sum_{t=1}^{n-k} (\hat{e}_t^2(\theta) - 1)(\hat{e}_{t+k}^2(\theta) - 1) \\ \gamma_k(m) &:= \frac{1}{n} \sum_{t=1}^{n-k} (e_t^2(m) - 1)(e_{t+k}^2(m) - 1); \hat{\gamma}_k(m) &:= \frac{1}{n} \sum_{t=1}^{n-k} (\hat{e}_t^2(m) - 1)(\hat{e}_{t+k}^2(m) - 1) \end{cases}$$

and  $\rho_k(\theta) := \frac{\gamma_k(\theta)}{\gamma_0(\theta)}$ ,  $\hat{\rho}_k(\theta) := \frac{\hat{\gamma}_k(\theta)}{\hat{\gamma}_0(\theta)}$ ,  $\rho_k(m) := \frac{\gamma_k(m)}{\gamma_0(m)}$  and  $\hat{\rho}_k(m) := \frac{\hat{\gamma}_k(m)}{\hat{\gamma}_0(m)}$ .

Finally, for  $K$  a positive integer, denote the vector of adjusted correlogram:

$$\hat{\rho}(\theta) := (\hat{\rho}_1(\theta), \dots, \hat{\rho}_K(\theta))' \quad \text{and} \quad \hat{\rho}(m) := (\hat{\rho}_1(m), \dots, \hat{\rho}_K(m))'.$$

*Proof.* (1) This proof is divided into two parts. In (i) we prove a result that ensures that the asymptotic distributions of the vectors  $\hat{\rho}(\theta)$  and  $\rho(\theta)$  are the same. In (ii) we show that the large sample distribution of  $\sqrt{n}\rho(m^*)$  is normal with a covariance matrix  $V(\theta^*, m^*)$ . Those two conditions do lead well to the asymptotic normality (2.5.1).

(i) In this part, we first show that for any  $k \in \mathbb{N}$ ,

$$\sqrt{n} \|\hat{\gamma}_k(\theta) - \gamma_k(\theta)\|_{\Theta} \xrightarrow[n \rightarrow \infty]{a.s.} 0. \quad (2.7.27)$$

We have:

$$\begin{aligned} \sqrt{n}(\hat{\gamma}_k(\theta) - \gamma_k(\theta)) &= \frac{1}{\sqrt{n}} \sum_{t=k+1}^n (\hat{e}_t^2(\theta) - 1)(\hat{e}_{t-k}^2(\theta) - 1) \\ &\quad - \frac{1}{\sqrt{n}} \sum_{t=k+1}^n (e_t^2(\theta) - 1)(e_{t-k}^2(\theta) - 1) \\ &= \frac{1}{\sqrt{n}} \sum_{t=k+1}^n (\hat{e}_t^2(\theta)\hat{e}_{t-k}^2(\theta) - e_t^2(\theta)e_{t-k}^2(\theta)) + \frac{1}{\sqrt{n}} \sum_{t=k+1}^n (\hat{e}_t^2(\theta) - e_t^2(\theta)) \\ &\quad + \frac{1}{\sqrt{n}} \sum_{t=k+1}^n (e_{t-k}^2(\theta) - \hat{e}_{t-k}^2(\theta)) \\ &=: I_1 + I_2 + I_3. \end{aligned}$$

Now, we show that  $\|I_1\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$ . We can rewrite  $I_1$  as follows

$$\begin{aligned} I_1 &= \frac{1}{\sqrt{n}} \sum_{t=k+1}^n \hat{e}_{t-k}^2(\theta)(\hat{e}_t^2(\theta) - e_t^2(\theta)) + \frac{1}{\sqrt{n}} \sum_{t=k+1}^n e_t^2(\theta)(\hat{e}_{t-k}^2(\theta) - e_{t-k}^2(\theta)) \\ &= \frac{1}{\sqrt{n}} \sum_{t=k+1}^n (\hat{e}_{t-k}^2(\theta) - e_{t-k}^2(\theta))(\hat{e}_t^2(\theta) - e_t^2(\theta)) + \frac{1}{\sqrt{n}} \sum_{t=k+1}^n e_{t-k}^2(\theta)(\hat{e}_t^2(\theta) - e_t^2(\theta)) \\ &\quad + \frac{1}{\sqrt{n}} \sum_{t=k+1}^n e_t^2(\theta)(\hat{e}_{t-k}^2(\theta) - e_{t-k}^2(\theta)) \\ &:= I_1^1 + I_1^2 + I_1^3. \end{aligned}$$

Let us show that  $\|I_1^1\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$  in our two frameworks.

a/ If  $X \subset AC(M_{\theta}, f_{\theta})$ , by Hölder's inequality, it follows from (2.7.8) that,

$$\begin{aligned} \mathbb{E} \left[ \left\| (\hat{e}_{t-k}^2(\theta) - e_{t-k}^2(\theta))(\hat{e}_t^2(\theta) - e_t^2(\theta)) \right\|_{\Theta}^{1/2} \right] &\leq \left( \mathbb{E} [\|\hat{e}_t^2(\theta) - e_t^2(\theta)\|_{\Theta}] \right. \\ &\quad \left. \times \mathbb{E} [\|\hat{e}_{t-k}^2(\theta) - e_{t-k}^2(\theta)\|_{\Theta}] \right)^{1/2}. \end{aligned}$$

But we have

$$\|\widehat{e}_t^2(\theta) - e_t^2(\theta)\|_{\Theta} \leq \frac{1}{\underline{h}} (2|X_t| + \|\widehat{f}_\theta^t\|_{\Theta} + \|f_\theta^t\|_{\Theta}) \|\widehat{f}_\theta^t - f_\theta^t\|_{\Theta} + \frac{4}{\underline{h}^{3/2}} (|X_t|^2 + \|f_\theta^t\|_{\Theta}^2) \|\widehat{M}_\theta^t - M_\theta^t\|_{\Theta}.$$

Therefore,

$$\begin{aligned} \mathbb{E}[\|\widehat{e}_t^2(\theta) - e_t^2(\theta)\|_{\Theta}] &\leq C \left( \mathbb{E}[|X_t|^2 + \|\widehat{f}_\theta^t\|_{\Theta}^2 + \|f_\theta^t\|_{\Theta}^2] \times \mathbb{E}[\|\widehat{f}_\theta^t - f_\theta^t\|_{\Theta}^2] \right)^{1/2} \\ &\quad + C \left( \mathbb{E}[|X_t|^4 + \|f_\theta^t\|_{\Theta}^2] \times \mathbb{E}[\|\widehat{M}_\theta^t - M_\theta^t\|_{\Theta}^2] \right)^{1/2} \\ &\leq C \left( \mathbb{E}\left[\left|\sum_{j \geq t} \alpha_j(f_\theta, \Theta) X_{t-j}\right|^2\right] \right)^{1/2} \\ &\quad + C \left( \mathbb{E}\left[\left|\sum_{j \geq t} \alpha_j(M_\theta, \Theta) X_{t-j}\right|^2\right] \right)^{1/2} \\ &\leq C \sum_{j \geq t} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta), \end{aligned}$$

using  $\mathbb{E}[|X_t|^4 + \|f_\theta^t\|_{\Theta}^2 + \|\widehat{f}_\theta^t\|_{\Theta}^2] < \infty$  and Cauchy-Schwarz Inequality. Hence,

$$\mathbb{E}\left[\left\|(\widehat{e}_{t-k}^2(\theta) - e_{t-k}^2(\theta))(\widehat{e}_t^2(\theta) - e_t^2(\theta))\right\|_{\Theta}^{1/2}\right] \leq C \sum_{j \geq t-k} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta).$$

Therefore, from Kounias and Weng (1969),  $\|I_1^1\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$  when

$$\sum_{t=1}^{\infty} t^{-1/4} \sum_{j \geq t} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) < \infty. \quad (2.7.28)$$

b/ if  $X \subset \widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$ , same computations imply  $\|I_1^1\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$  when

$$\sum_{t=1}^{\infty} t^{-1/4} \sum_{j \geq t} \alpha_j(\widetilde{H}_\theta, \Theta) < \infty. \quad (2.7.29)$$

Since  $\mathbb{E}[\|\widehat{e}_t^2(\theta)\|_{\Theta}] \leq 2\underline{h}^{-1} \mathbb{E}[X_t^2 + \|f_\theta^t\|_{\Theta}^2] < \infty$  and similarly  $\mathbb{E}[\|\widehat{e}_t^2(\theta)\|_{\Theta}] < \infty$ , we deduce from the same inequalities as in the first case of  $I_1^1$  that  $\|I_1^2\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$  and  $\|I_1^3\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$  when

$$\sum_{t=1}^{\infty} t^{-1/4} \left( \sum_{j \geq t} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) + \alpha_j(\widetilde{H}_\theta, \Theta) \right)^{1/2} < \infty, \quad (2.7.30)$$

which is also the condition for insuring that  $\|I_2\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$  and  $\|I_3\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0$ . This ends the proof of (2.7.22).

Finally, since  $\widehat{\rho}_k(\theta) = \widehat{\gamma}_k(\theta)/\widehat{\gamma}_0(\theta)$  and  $\rho_k(\theta) = \gamma_k(\theta)/\gamma_0(\theta)$ , with  $\gamma_0(\theta) > 0$ , we deduce under condition (2.7.30) that

$$\sqrt{n} \|\widehat{\rho}_k(\theta) - \rho_k(\theta)\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0 \quad \text{for any } k \geq 1. \quad (2.7.31)$$

This also implies

$$\sqrt{n} |\widehat{\rho}_k(m^*) - \rho_k(m^*)| \xrightarrow[n \rightarrow +\infty]{a.s.} 0 \quad \text{for any } k \geq 1. \quad (2.7.32)$$

(ii) The proof of this result has already been done in [Li and Mak \(1994\)](#) but in a Gaussian framework. We recall here the main lines while avoiding the Gaussian assumption. The first step is to use a Taylor expansion of the function  $\gamma$ . Hence, we have for each  $k = 1, \dots, K$ ,

$$\sqrt{n} \gamma_k(m^*) = \sqrt{n} \gamma_k(\hat{\theta}(m^*)) = \sqrt{n} \gamma_k(\theta^*) + \partial_{\theta} \gamma_k(\bar{\theta}^{(k)}) \sqrt{n} ((\hat{\theta}(m^*))_i - \theta_i^*)_{i \in m^*}, \quad (2.7.33)$$

where  $\partial_{\theta} \gamma_k = {}^t(\partial \gamma_k / \partial \theta_i)_{i \in m^*}$ , and  $\bar{\theta}^{(k)}$  is in the ball of centre  $\theta^*$  and radius  $\|(\hat{\theta}(m^*) - \theta^*)_{i \in m^*}\|$ . We also have

$$\begin{aligned} \partial_{\theta} \gamma_k(\theta) = & -\frac{2}{n} \left( \sum_{t=k+1}^n e_t^2(\theta) (e_{t-k}^2(\theta) - 1) \frac{\partial_{\theta} M_{\theta}^t}{M_{\theta}^t} + e_t(\theta) (e_{t-k}^2(\theta) - 1) \frac{\partial_{\theta} f_{\theta}^t}{M_{\theta}^t} \right. \\ & \left. + e_{t-k}(\theta) (e_t^2(\theta) - 1) \frac{\partial_{\theta} f_{\theta}^{t-k}}{M_{\theta}^{t-k}} + e_{t-k}^2(\theta) (e_t^2(\theta) - 1) \frac{\partial_{\theta} M_{\theta}^{t-k}}{M_{\theta}^{t-k}} \right). \end{aligned} \quad (2.7.34)$$

We have  $\mathbb{E}[e_{t-k}(\theta^*) (e_t^2(\theta^*) - 1) \frac{\partial f_{\theta^*}^{t-k}}{M_{\theta^*}^{t-k}} \mid \sigma((\xi_s)_{s \leq t-k})] = e_{t-k}(\theta^*) \frac{\partial f_{\theta^*}^{t-k}}{M_{\theta^*}^{t-k}} \mathbb{E}[e_t^2(\theta^*) - 1] = 0$  since we have assumed  $\mathbb{E}[\xi_0^2] = 1$ . Moreover,  $\mathbb{E}[e_t(\theta^*) \frac{\partial f_{\theta^*}^t}{M_{\theta^*}^t}] = \mathbb{E}[\xi_t \frac{\partial f_{\theta^*}^t}{M_{\theta^*}^t}] = 0$  and this implies  $\mathbb{E}[e_t(\theta^*) (e_{t-k}^2(\theta^*) - 1) \frac{\partial f_{\theta^*}^t}{M_{\theta^*}^t}] = 0$ . As a consequence, the expectation of the three last terms of (2.7.34) vanishes for  $\theta = \theta^*$ . By using the Ergodic Theorem, we finally obtained:

$$\partial_{\theta} \gamma_k(\theta^*) \xrightarrow[n \rightarrow +\infty]{a.s.} -2 \mathbb{E} \left[ e_k^2(\theta^*) (e_0^2(\theta^*) - 1) \frac{\partial_{\theta} M_{\theta^*}^k}{M_{\theta^*}^k} \right] = -2 \mathbb{E} \left[ (\xi_0^2 - 1) \partial_{\theta} \log(M_{\theta^*}^k) \right].$$

Moreover, since  $\partial_{\theta^2}^2 f_{\theta}$  and  $\partial_{\theta^2}^2 M_{\theta}$  exist, and since  $\hat{\theta}(m^*) \xrightarrow[n \rightarrow +\infty]{a.s.} \theta^*$ , we deduce that the same almost sure convergence occurs for  $\partial_{\theta} \gamma_k(\bar{\theta}^{(k)})$ . Then, we finally obtain

$$(\partial_{\theta} \gamma_k(\bar{\theta}^{(k)}))_{1 \leq k \leq K} \xrightarrow[n \rightarrow +\infty]{a.s.} J_K(m^*) = -2 \left( \mathbb{E} \left[ (\xi_0^2 - 1) \frac{\partial}{\partial \theta_j} \log(M_{\theta^*}^j) \right] \right)_{1 \leq j \leq K, j \in m^*}. \quad (2.7.35)$$

Under the assumptions, a central limit theorem for  $\hat{\theta}(m^*)$  has been established in [Bardet and Wintenberger \(2009\)](#), and this implies

$$\begin{aligned} (\partial_{\theta} \gamma_k(\bar{\theta}^{(k)}))_{1 \leq k \leq K} \sqrt{n} ((\hat{\theta}(m^*))_i - \theta_i^*)_{i \in m^*} \\ \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}_K \left( 0, J_K(m^*) F(\theta^*, m^*)^{-1} G(\theta^*, m^*) F(\theta^*, m^*)^{-1} J_K'(m^*) \right). \end{aligned} \quad (2.7.36)$$

On the other hand, when  $\theta = \theta^*$ ,  $e_t^2(\theta^*) = \xi_t^2$  for any  $t \in \mathbb{Z}$  and since  $\mathbb{E}[\xi_0^2] = 1$ , we deduce that  $(e_t^2(\theta^*) - 1)_t$  is a sequence of centred iid random variables with variance  $\mu_4 - 1$  with  $\mu_4 = \mathbb{E}[\xi_0^4]$ . In such as case, the asymptotic behavior of the covariograms is well known and we deduce:

$$\sqrt{n} (\gamma_k(\theta^*))_{1 \leq k \leq K} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}_K(0, (\mu_4 - 1)^2 I_K), \quad (2.7.37)$$

with  $I_k$  the  $(K \times K)$  identity matrix.

We would like to use (2.7.33) for obtaining the asymptotic behavior of  $\gamma(m^*)$ . In (2.7.36) and (2.7.37), we obtained the asymptotic normality of each of the two terms composing



$\gamma(m^*)$ . Now we need to study the joint asymptotic behavior of  $\sqrt{n} \gamma(\theta^*)$  and  $\sqrt{n} ((\hat{\theta}(m^*))_i - \theta_i^*)_{i \in m^*}$ .

Using the proof of the asymptotic normality of the QMLE (see for instance [Bardet and Wintenberger \(2009\)](#)), a Taylor expansion of log-likelihood for large  $n$  leads to

$$((\hat{\theta}(m^*))_i - \theta_i^*)_{i \in m^*} \approx 2 (F(\theta^*, m^*))^{-1} \frac{1}{n} \frac{\partial}{\partial \theta} L_n(\theta^*).$$

Therefore, the asymptotic cross expectation between  $(\partial_\theta \gamma_k(\bar{\theta}^{(k)}))_k \sqrt{n} ((\hat{\theta}(m^*))_i - \theta_i^*)_{i \in m^*}$  and  $\sqrt{n} \gamma(\theta^*)$  is equal to:

$$- J_K(m^*) F(\theta^*, m^*)^{-1} \mathbb{E} \left[ \frac{\partial}{\partial \theta} L_n(\theta^*) \gamma(\theta^*)' \right]. \quad (2.7.38)$$

From (2.2.1), a direct differentiation of  $L_n$  provides

$$\frac{\partial}{\partial \theta} L_n(\theta^*) = \sum_{t=1}^n (e_t^2(\theta^*) - 1) \frac{\partial}{\partial \theta} \log(M_{\theta^*}^t) + \sum_{t=1}^n e_t(\theta^*) \frac{\partial}{\partial \theta} f_{\theta^*}^t$$

so that,

$$\begin{aligned} \mathbb{E} \left[ \frac{\partial}{\partial \theta} L_n(\theta^*) \gamma_k(\theta^*) \right] &= \frac{1}{n} \mathbb{E} \left[ \sum_{i=1}^n (e_i^2(\theta^*) - 1) \frac{\partial}{\partial \theta} \log(M_{\theta^*}^i) \right. \\ &\quad \times \left. \sum_{j=k+1}^n (e_j^2(\theta^*) - 1) (e_{j-k}^2(\theta^*) - 1) \right] \\ &+ \frac{1}{n} \mathbb{E} \left[ \sum_{i=1}^n e_i(\theta^*) \frac{\partial}{\partial \theta} f_{\theta^*}^i \sum_{j=k+1}^n (e_j^2(\theta^*) - 1) (e_{j-k}^2(\theta^*) - 1) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=k+1}^n \mathbb{E} \left[ (\xi_i^2 - 1) (\xi_j^2 - 1) (\xi_{j-k}^2 - 1) \frac{\partial}{\partial \theta} \log(M_{\theta^*}^i) \right] \\ &+ \frac{1}{n} \sum_{i=1}^n \sum_{j=k+1}^n \mathbb{E} \left[ \xi_i (\xi_j^2 - 1) (\xi_{j-k}^2 - 1) \frac{\partial}{\partial \theta} f_{\theta^*}^i \right]. \end{aligned}$$

Using conditional expectations, we have  $\mathbb{E} \left[ (\xi_i^2 - 1) (\xi_j^2 - 1) (\xi_{j-k}^2 - 1) \frac{\partial}{\partial \theta} \log(M_{\theta^*}^i) \right] = 0$  for  $i \neq j$  since  $k \geq 1$ . Moreover, for  $i = j$ , we obtain:

$$\mathbb{E} \left[ (\xi_i^2 - 1) (\xi_j^2 - 1) (\xi_{j-k}^2 - 1) \frac{\partial}{\partial \theta} \log(M_{\theta^*}^i) \right] = (\mu_4 - 1) \mathbb{E} \left[ (\xi_{i-k}^2 - 1) \frac{\partial}{\partial \theta} \log(M_{\theta^*}^i) \right],$$

which is the row  $k$  of matrix  $-\frac{(\mu_4-1)}{2} J_K(m^*)$ . Similarly, and using the assumption  $\mathbb{E}[\xi_0^3] = 0$ , we obtain  $\mathbb{E} \left[ \xi_i (\xi_j^2 - 1) (\xi_{j-k}^2 - 1) \frac{\partial}{\partial \theta} f_{\theta^*}^i \right] = 0$  for any  $i, j$  and  $k$ . As a consequence,

$$\begin{aligned} \text{Cov} \left( \sqrt{n} \gamma(\theta^*), (\partial_\theta \gamma_k(\bar{\theta}^{(k)}))_k \sqrt{n} ((\hat{\theta}(m^*))_i - \theta_i^*)_{i \in m^*} \right) \\ \xrightarrow{n \rightarrow \infty} \frac{1}{2} (\mu_4 - 1) J_K(m^*) F(\theta^*, m^*)^{-1} J'_K(m^*). \end{aligned}$$

Finally, we deduce the asymptotic covariance matrix of  $\sqrt{n} \gamma(m^*)$ , which is

$$\begin{aligned} (\mu_4 - 1)^2 I_K + J_K(m^*) F(\theta^*, m^*)^{-1} G(\theta^*, m^*) F(\theta^*, m^*)^{-1} J'_K(m^*) \\ + (\mu_4 - 1) J_K(m^*) F(\theta^*, m^*)^{-1} J'_K(m^*). \end{aligned}$$

Moreover the vector  $\gamma(m^*)$  is normal distributed from Lemma 3.3 of [Ling and Li \(1997\)](#). Thus, using Slutsky Lemma and with  $\gamma_0(m^*) \xrightarrow[n \rightarrow +\infty]{a.s.} \mu_4 - 1$ , and with  $\rho_k(m^*) = \gamma_k(m^*)/\gamma_0(m^*)$ , the limit theorem (2.5.1) holds with

$$V(\theta^*, m^*) := I_K + (\mu_4 - 1)^{-2} J_K(m^*) F(\theta^*, m^*)^{-1} G(\theta^*, m^*) F(\theta^*, m^*)^{-1} J'_K(m^*) \\ - 2(\mu_4 - 1)^{-1} J_K(m^*) F(\theta^*, m^*)^{-1} J'_K(m^*). \quad (2.7.39)$$

The proof is achieved after using the limit theorem (2.7.32).

(2) (2.5.2) follows directly from (2.5.1).

(3) We follow a same reasoning like in the proof of Theorem 2.2. For  $x = (x_k)_{1 \leq k \leq K} \in \mathbb{R}^K$ , denote by  $F_n(x) = \mathbb{P}\left(\bigcap_{1 \leq k \leq K} \sqrt{n}(\hat{\rho}(\hat{m}))_k \leq x_k\right)$  the distribution function of  $\sqrt{n}\hat{\rho}(\hat{m})$ .

Applying the Total Probability Rule and by virtue of Theorem 2.1, we obtain:

$$F_n(x) = \mathbb{P}\left(\bigcap_{1 \leq k \leq K} \sqrt{n}(\hat{\rho}(m^*))_k \leq x_k\right).$$

Therefore, the vectors  $\sqrt{n}\hat{\rho}(\hat{m})$  and  $\sqrt{n}\hat{\rho}(m^*)$  have exactly the same distribution.  $\square$



# 3

## General Hannan and Quinn Criterion for Common Time Series

---

3.1	INTRODUCTION .....	66
3.2	MODEL SELECTION CONSISTENCY .....	68
3.2.1	Model Selection Procedure .....	68
3.2.2	Assumptions .....	69
3.2.3	Consistency Result .....	70
3.2.4	Another Quasi-Likelihood Function .....	72
3.2.5	Algorithm of Calibration of the minimal constant .....	73
3.3	NUMERICAL EXPERIMENTS .....	74
3.3.1	Monte Carlo: Consistency .....	74
3.3.2	Real Data Analysis: financial time series .....	76
3.4	Proofs .....	76
3.4.1	Proof of Theorem 3.1 .....	76
3.4.1.a	Overfitting Case .....	77
3.4.1.b	Misspecification/Underfitting Case .....	77
3.4.1.c	Proof of Theorem 3.3 .....	79
3.4.1.d	Proof of Proposition 3.1 .....	79
3.4.1.e	Proof of Proposition 3.2 .....	81
3.4.2	Technical Lemmas .....	82

---

The content of this chapter was taken from the article submitted for publication:  
K. KAMILA "General Hannan and Quinn Criterion for Common Time Series "  
<https://arxiv.org/pdf/2101.04210.pdf> .

### Abstract

This chapter aims to study data driven model selection criteria for a large class of time

series, which includes ARMA or  $\text{AR}(\infty)$  processes, as well as GARCH or  $\text{ARCH}(\infty)$ , APARCH and many others processes. We tackled the challenging issue of designing adaptive criteria which enjoys the strong consistency property. When the observations are generated from one of the aforementioned models, the new criteria, select the true model almost surely asymptotically. The proposed criteria are based on the minimization of a penalized contrast akin to the Hannan and Quinn's criterion and then involved a term which is known for most classical time series models and for more complex models, this term can be data driven calibrated. Monte-Carlo experiments and an illustrative example on the CAC 40 index are performed to highlight the obtained results.

### 3.1 INTRODUCTION

A common solution in model selection is to choose the model, minimizing a penalized based criterion which is the sum of two terms: the first one is the empirical risk (least squares, likelihood) that measures the goodness of fit and the second one is an increasing function of the complexity which aims to penalize large models and control the bias. Therefore a challenging task when designing a penalized criterion is the specification of the penalty term. Considering leading model selection criteria (BIC, AIC,  $C_p$ , HQ to name a few), one can see that the penalty term is a product of the model dimension with a sequence which is specific to the criteria. Indeed, a criterion is designed according to the goal one would like to achieve. The classical properties for model selection criteria include *consistency*, *efficiency* (oracle inequality, asymptotic optimality), *adaptive in the minimax sense*.

In this chapter, we focus on consistency property which aims at identifying the data generating process with high probability or almost surely. Hence, it requires the assumption whereby there exists a true model in the set of competitive models and the goal is to select this with probability approaches one as the sample size tends to infinity. In [Bardet et al. \(2020b\)](#), they studied model selection criteria regarding consistency in a large class of time series, which is the interest of this paper. The leading criterion obtained in this framework is the BIC; with a relatively heavy penalty, it ensures the selection of quite simple models. Moreover, several papers have established the consistency property in particular settings. For instance, [Hannan and Quinn \(1979\)](#) shows that the Hannan and Quinn (HQ) penalty  $c \log \log n$  with  $c > 2$  leads to a consistent choice of the true order in the framework of AR type models. One year later, [Hannan \(1980\)](#) (or [Hannan and Deistler \(2012\)](#)) extended this result for ARMA models.

Also, it has been proven in several contexts, that the BIC criterion [Schwarz \(1978\)](#) enjoys the consistency property: [Shibata \(1986\)](#) in the density estimation using hypothesis testing for autoregressive moving average models, [Lebarbier and Mary-Huard \(2004\)](#) in density estimation for independent observations, [Bardet et al. \(2020b\)](#) for a general class of time series, to name a few.

Compare to HQ penalty, the BIC penalty does not have the slowest rate of increase and then it can very often choose very simple models possible wrongs for small samples [Hannan and Quinn \(1979\)](#). Moreover, the HQ criterion has been derived for linear time series: AR models in [Hannan and Quinn \(1979\)](#), ARMA models in [Hannan \(1980\)](#) and [Hannan and Deistler \(2012\)](#). Is the HQ penalty still strongly consistent for heteroscedastic nonlinear models such as GARCH, APARCH or ARMA-GARCH? And what about a general class including linear and non linear models as well?

That raises a challenging question of designing robust penalties for most classical time

series models enjoying the model selection consistency. This is the issue we want to address in this chapter for a general class of times series called affine causal and defined below.

**Class  $\mathcal{AC}(M, f)$  :** A process  $X = (X_t)_{t \in \mathbb{Z}}$  belongs to  $\mathcal{AC}(M, f)$  if it satisfies:

$$X_t = M((X_{t-i})_{i \in \mathbb{N}^*}) \xi_t + f((X_{t-i})_{i \in \mathbb{N}^*}) \quad \text{for any } t \in \mathbb{Z}. \quad (3.1.1)$$

where  $(\xi_t)_{t \in \mathbb{Z}}$  is a sequence of zero-mean independent identically distributed random vectors (i.i.d.r.v) satisfying  $\mathbb{E}(|\xi_0|^r) < \infty$  with  $r \geq 1$  and  $M, f : \mathbb{R}^\infty \rightarrow \mathbb{R}$  are two measurable functions.

For instance,

- if  $M((X_{t-i})_{i \in \mathbb{N}^*}) = \sigma$  and  $f((X_{t-i})_{i \in \mathbb{N}^*}) = \sum_{i=1}^{\infty} \phi_i X_{t-i}$ , then  $(X_t)_{t \in \mathbb{Z}}$  is an  $\text{AR}(\infty)$  process;
- if  $M((X_{t-i})_{i \in \mathbb{N}^*}) = \sqrt{a_0 + a_1 X_{t-1}^2 + \dots + a_p X_{t-p}^2}$  and  $f((X_{t-i})_{i \in \mathbb{N}^*}) = 0$ , then  $(X_t)_{t \in \mathbb{Z}}$  is an  $\text{ARCH}(p)$  process.

Note that, numerous classical time series models such as  $\text{ARMA}(p, q)$ ,  $\text{GARCH}(p, q)$ ,  $\text{ARMA}(p, q)$ - $\text{GARCH}(p, q)$  (see [Ding et al. \(1993\)](#) and [Ling and McAleer \(2003\)](#)) or  $\text{APARCH}(\delta, p, q)$  processes (see [Ding et al. \(1993\)](#)) belongs to  $\mathcal{AC}(M, f)$ .

The study of this type of process more often requires the classical regularity conditions on the functions  $M$  and  $f$ , which are not restrictive at all and remain valid in various time serie models. Let us recall these conditions for  $\Psi_\theta = f_\theta$  or  $M_\theta$  and  $\Theta$  a compact set.

**Hypothesis  $\mathbf{A}(\Psi_\theta, \Theta)$ :** Assume that  $\|\Psi_\theta(0)\|_\Theta < \infty$  and there exists a sequence of non-negative real numbers  $(\alpha_k(\Psi_\theta, \Theta))_{k \geq 1}$  such that  $\sum_{k=1}^{\infty} \alpha_k(\Psi_\theta, \Theta) < \infty$  satisfying:

$$\|\Psi_\theta(x) - \Psi_\theta(y)\|_\Theta \leq \sum_{k=1}^{\infty} \alpha_k(\Psi_\theta, \Theta) |x_k - y_k| \quad \text{for all } x, y \in \mathbb{R}^\infty.$$

In addition, if the noise  $\xi_0$  admits  $r$ -order moments (for  $r \geq 1$ ), let us define:

$$\Theta(r) = \left\{ \theta \in \mathbb{R}^d, A(f_\theta, \{\theta\}) \text{ and } A(M_\theta, \{\theta\}) \text{ hold with } \sum_{k=1}^{\infty} \alpha_k(f_\theta, \{\theta\}) + \|\xi_0\|_r \sum_{k=1}^{\infty} \alpha_k(M_\theta, \{\theta\}) < 1 \right\}. \quad (3.1.2)$$

Under this assumption, [Doukhan and Wintenberger \(2008\)](#) showed that there exists a stationary and ergodic solution to (3.1.1) with  $r$ -order moment for any  $\theta \in \Theta(r)$ . Moreover, [Bardet and Wintenberger \(2009\)](#) studied the consistency and the asymptotic normality of the QMLE of  $\theta^*$  for  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$ .

The main contribution of this chapter is the generalization of the HQ criterion to affine causal class: we provide a minimal multiplicative penalty term  $c_{min}$  so that all penalties of the form  $2c \log \log n D_m$  with  $c \geq c_{min}$  ensure the strong consistency property for affine causal models under some mild conditions on the Lipschitz coefficients of functions  $M_\theta, f_\theta$  ( $D_m$  denotes the size of the model  $m$ ). Monte Carlo experiments have been conducted in order to attest the accuracy of our new criteria.

The chapter is organized as follows. The model selection consistency along with notations and assumptions are described in Section 3.2. Numerical results are presented in Section 3.3 and Section 3.4 contains the proofs.

## 3.2 MODEL SELECTION CONSISTENCY

### 3.2.1 Model Selection Procedure

Let assume  $(X_1, \dots, X_n)$  be a trajectory of a stationary affine causal process  $m^* := \mathcal{AC}(M_{\theta^*}, f_{\theta^*})$ , where  $\theta^*$  is unknown. The goal of the consistency property is to come up with this true model given a set of candidate model  $\mathcal{M}$  such that  $m^* \in \mathcal{M}$ .

A  $D_m$ -dimensional model  $m \in \mathcal{M}$  can be viewed as a set of causal functions  $(M_\theta, f_\theta)$  with  $\theta \in \Theta(m) \subset \mathbb{R}^{D_m}$ .  $\Theta(m)$  is the parameter set of the model  $m$ .

The consistency property will be studied using quasi likelihood functions since assumption on the distribution of the noise is not required.

The Gaussian quasi log-likelihood is derived from the conditional (with respect to the filtration  $\sigma(X_t, t \leq 0)$ ) log-likelihood of  $(X_1, \dots, X_n)$  when  $(\xi_t)$  is supposed to be a Gaussian standard white noise. From (3.1.1), one deduce that the log density of  $X_t$  given  $\sigma(X_i, i < t)$  is

$$-\frac{1}{2} \left[ \frac{(X_t - f_{\theta^*}^t)^2}{H_{\theta^*}^t} + \log(H_{\theta^*}^t) \right].$$

Therefore the conditional log density of  $(X_1, \dots, X_n)$  given  $\sigma(X_t, t \leq 0)$  is

$$-\frac{1}{2} \sum_{t=1}^n \left[ \frac{(X_t - f_{\theta^*}^t)^2}{H_{\theta^*}^t} + \log(H_{\theta^*}^t) \right].$$

From now on, we drop the Gaussian assumption of the noise. The conditional log-density inspires to define for all  $\theta \in \Theta$

$$L_n(\theta) := -\frac{1}{2} \sum_{t=1}^n q_t(\theta), \text{ with } q_t(\theta) := \frac{(X_t - f_\theta^t)^2}{H_\theta^t} + \log(H_\theta^t) \quad (3.2.1)$$

where  $f_\theta^t := f_\theta(X_{t-1}, X_{t-2}, \dots)$ ,  $M_\theta^t := M_\theta(X_{t-1}, X_{t-2}, \dots)$  and  $H_\theta^t = (M_\theta^t)^2$ . The quasi likelihood function  $L_n$  is not computable since it depends on the past  $(X_{-j})_{j \in \mathbb{N}}$  that is unknown. However, the sequence  $(L_n(\cdot))_n$  enjoys very nice asymptotic properties such that the Uniform Law of Large Numbers (see [Bardet and Wintenberger \(2009\)](#)).

Let  $\mathcal{M}$  a finite family of candidate models containing the true one  $m^*$ . According to Proposition 1 in [Bardet et al. \(2020b\)](#), all these models can be included into a big one with parameter space  $\Theta$ . For each specific model  $m \in \mathcal{M}$ , we define the Gaussian QMLE  $\hat{\theta}(m)$  with respect to  $L_n$  as

$$\hat{\theta}(m) = \operatorname{argmax}_{\theta \in \Theta(m)} L_n(\theta). \quad (3.2.2)$$

To select the true model  $m \in \mathcal{M}$ , we consider a penalized contrast  $C(m)$  ensuring a trade-off between  $-2$  times the maximized log-likelihood, which decreases with the size of the model, and a penalty increasing with the size of the model. Therefore, the choice of the "best" model  $\hat{m}$  among the estimated can be performed by minimizing the following criteria

$$\hat{m} = \operatorname{argmin}_{m \in \mathcal{M}} C(m) \quad \text{with} \quad C(m) = -2 L_n(\hat{\theta}(m)) + \kappa_n(m) \quad (3.2.3)$$

where  $(\kappa_n)_n$  an increasing sequence depending on the number of observations  $n$  and the dimension  $D_m$ . There exist several possible choices of  $\kappa_n(m)$  including

- $\kappa_n(m) = 2c D_m \log \log n$  with  $c > 1$ , we retrieve the HQ criterion [Hannan and Quinn \(1979\)](#);

- $\kappa_n = D_m \log n$ ,  $C$  yields to BIC criterion [Schwarz \(1978\)](#);
- $\kappa_n = 2 D_m$ ,  $C$  is the AIC criterion [Akaike \(1973\)](#).

Basically the principle is that by increasing the size of the model, the likelihood increases also. The question is whether this increase in complexity is offset by a sufficient increase in likelihood. If the answer is no, then the least complex model is used, even if it is less likely. If the answer is affirmative, then we accept to work with a more complex model. Of course, all the difficulty lies in the choice of weights between likelihood and complexity, and thus ultimately in the specification of the penalty multiplicative term  $\kappa_n$ .

What is the better weighting term of the model complexity? The aim here is by leveraging the increasing rate of the likelihood, to propose a data driven  $\kappa_n$  in order to guarantee the strong consistency property to our model selection procedure i.e.

$$\hat{m} \xrightarrow[n \rightarrow \infty]{a.s.} m^*. \quad (3.2.4)$$

### 3.2.2 Assumptions

Some mild conditions will be required to prove the consistency of the considered model selection criteria.

The following assumption is well-known as the identifiability one and is always required in order to guarantee the unicity of the global maximum of the MLE at the true parameter  $\theta^*$ . That is:

**Assumption A1:** For all  $\theta, \theta' \in \Theta_m$ ,  $(f_\theta^0 = f_{\theta'}^0)$  and  $(M_\theta^0 = M_{\theta'}^0) \implies \theta = \theta'$ .

Another required assumption concerns the differentiability of  $\Psi_\theta = f_\theta$  or  $M_\theta$  on  $\Theta$ . This type of assumption has already been considered in order to apply the QMLE procedure (see [Bardet and Wintenberger \(2009\)](#), [Straumann and Mikosch \(2006\)](#), [White \(1982\)](#)).

The following condition provides the invertibility of the Fisher's information matrix of  $(X_1, \dots, X_n)$  and was used to prove the asymptotic normality of the QMLE (see [Bardet and Wintenberger \(2009\)](#)).

**Assumption A2:** For any  $x \in \mathbb{R}^\infty$ , the functions  $\theta \rightarrow M_\theta$  and  $\theta \rightarrow f_\theta$  are  $\mathcal{C}^2(\Theta)$  and one of the families  $(\partial f_\theta^t / \partial \theta^{(i)})_{1 \leq i \leq D_{m^*}}$  or  $(\partial H_\theta^t / \partial \theta^{(i)})_{1 \leq i \leq D_{m^*}}$  is a.e. linearly independent.

Note that the definitions of the conditional log-likelihood requires that their denominators do not vanish. Hence, we will suppose in the sequel that the lower bound of  $H_\theta(\cdot) = (M_\theta(\cdot))^2$  (which is reached since  $\Theta$  is compact) is strictly positive:

**Assumption A3:**  $\exists \underline{h} > 0$  such that  $\inf_{\theta \in \Theta} (H_\theta(x)) \geq \underline{h}$  for all  $x \in \mathbb{R}^\infty$ .

Next we assume the existence of the eighth order moment of the noise.

**Assumption A4:**  $\mathbb{E}[\xi_0^8] < \infty$ .

We end the list of assumptions by assuming a suitable relation between the Fisher Information matrix  $G(\theta_m^*)$  and the limiting Hessian matrix of the log-likelihood  $F(\theta_m^*)$  defined as follows

$$(F(\theta_m^*))_{i,j} = \mathbb{E} \left[ \frac{\partial^2 q_0(\theta_m^*)}{\partial \theta_i \partial \theta_j} \right] \quad \text{and} \quad (G(\theta_m^*))_{i,j} = \mathbb{E} \left[ \frac{\partial q_0(\theta_m^*)}{\partial \theta_i} \frac{\partial q_0(\theta_m^*)}{\partial \theta_j} \right],$$



with  $\theta_m^* := (\theta^*, 0, \dots, 0)^\top \in \Theta(m)$ .

**Assumption A5:** There exist absolute constants  $\alpha_1$  and  $\alpha_2$  such that for any  $m \in \mathcal{M}$  verifying  $m^* \subset m$ ,

$$\mathbf{1}_m^\top \Sigma_{\theta_m^*} \mathbf{1}_m = \alpha_1 D_m^1 + \alpha_2 D_m^2 \quad (3.2.5)$$

where  $D_m^1$  and  $D_m^2$  are two integers such that  $D_m^1 + D_m^2 = D_m$ ,  $\mathbf{1}_m := (1, 1, \dots, 1)^\top \in \mathbb{R}^{D_m}$ ,  $\Sigma_{\theta_m^*} := G(\theta_m^*)^{1/2} F(\theta_m^*)^{-1} G(\theta_m^*)^{1/2}$ .

For most classical affine causal models, **A5** is verified (see Proposition 3.2). However, for more complex models such as ARMA-GARCH with  $\mu_4 \neq 3$ ,  $\Sigma_{\theta_m^*}$  is hard to handle.

### 3.2.3 Consistency Result

Before stating the main result of this section, we give important intermediate results. All proof of the results stated in this subsection can be found in Section 3.4.

The following Proposition suggests the existence of a term that will be the keystone of this work.

**Proposition 3.1.** *Let  $m^*$  any affine causal model. For any model  $m$  verifying  $m^* \subset m$ , and under **A1-A5**, it holds*

$$\limsup_{n \rightarrow \infty} \frac{L_n(\hat{\theta}(m)) - L_n(\theta_m^*)}{2 \log \log n} = \frac{1}{4} (\alpha_1 D_m^1 + \alpha_2 D_m^2) \quad a.s. \quad (3.2.6)$$

where the constants  $\alpha_1, \alpha_2, D_m^1, D_m^2$  are specified in assumption **A5**.

For every  $m \in \mathcal{M}$ , let us denote by  $c_{min}(m)$  the following term that will be used several times

$$c_{min}(m) := \frac{1}{4} (\alpha_1 D_m^1 + \alpha_2 D_m^2) \quad (3.2.7)$$

Now we state a result which provides the values of both  $\alpha_1$  and  $\alpha_2$  for most classical affine causal models.

**Proposition 3.2.** *Under the assumptions and notation of Proposition 3.1, we have*

- If  $\mu_4 = \mathbb{E}[\xi_0^4] = 3$  (for instance for Gaussian noise), then  $\alpha_1 = 2, \alpha_2 = 2$  and  $c_{min}(m) = \frac{1}{2} D_m$ ;
- If the parameter  $\theta$  identifying an affine causal model  $X_t = M_\theta^t \xi_t + f_\theta^t$  can be decomposed as  $\theta = (\theta_1, \theta_2)'$  with  $f_\theta^t = \tilde{f}_{\theta_1}^t$  and  $M_\theta^t = \tilde{M}_{\theta_2}^t$ , then  $\alpha_1 = 2, \alpha_2 = \mu_4 - 1$  and

$$c_{min}(m) = \frac{1}{2} D_m^1 + \frac{\mu_4 - 1}{4} D_m^2$$

The second configuration in Proposition 3.2 includes classical time series

- GARCH( $p, q$ ), APARCH( $\delta, p, q$ ) type models and related ones,  $c_{min}(m) = \frac{\mu_4 - 1}{4} D_m$ ;
- ARMA( $p, q$ ) models,  $c_{min}(m) = \frac{D_m}{2}$  if the variance of the noise is known and  $c_{min}(m) = \frac{D_m - 1}{2} + \frac{\mu_4 - 1}{4}$  otherwise.

We can now state the first main result of this paper.

**Theorem 3.1.** *Let  $(X_1, \dots, X_n)$  be an observed trajectory of an affine causal process  $X$  belonging to  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$  where  $\theta^*$  is an unknown vector belonging to  $\Theta(r) \subset \mathbb{R}^{D_{m^*}}$ . Let also  $\mathcal{M}$  be a finite family of candidate models such that  $m^* \in \mathcal{M}$ . If assumptions **A1-A5** hold, then with  $c_{\min} := \max(\frac{\alpha_1}{2}, \frac{\alpha_2}{2})$ , it holds*

for any  $\kappa_n(m) = 2c D_m \log \log n$  with

$$c \geq c_{\min} \quad (3.2.8)$$

it holds for the selected model  $\hat{m}$  according to (3.2.3)

$$\hat{m} \xrightarrow[n \rightarrow \infty]{a.s.} m^*, \quad (3.2.9)$$

**Remark 3.1.** 1. For classical configurations as seen in Proposition 3.2, this result gives a generalization of Hannan and Quinn criterion.

2. For more complex models, the values of  $\alpha_1$  and  $\alpha_2$  are unknowns (at least until a better relationship between matrix  $F(\theta_m^*)$  and  $G(\theta_m^*)$  is found) and so  $c_{\min}$  is also unknown. In these cases, we propose to use adaptive methods such as slope heuristic algorithm or dimension jump [Arlot and Massart \(2009\)](#) to calibrate  $c_{\min}$ .

Let us mention that our result generalizes the strong consistency obtained by [Hannan and Quinn \(1979\)](#) for AR models as the affine causal class also contains GARCH models. It furthermore generalizes the result [Hannan and Deistler \(2012\)](#) for ARMA models.

Theorem 3.1 gives a theoretical guarantee on the consistency of the model selection procedure. However, it does not say anything about the convergence (and its rate) of the parameter estimate resulting from the model selection  $\hat{m}$ . The following results shows that the final estimate  $\hat{\theta}_{\hat{m}}$  is consistent and verifies a CLT.

**Theorem 3.2.** *Under the assumptions of Theorem 3.1, it holds*

$$\hat{\theta}(\hat{m}) \xrightarrow[n \rightarrow \infty]{\mathcal{P}} \theta^*. \quad (3.2.10)$$

Moreover,

$$\sqrt{n} \left( (\hat{\theta}(\hat{m}))_i - (\theta^*)_i \right)_{i \in m^*} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}_{|m^*|}(0, \Sigma_{\theta^*, m^*}) \quad (3.2.11)$$

with  $\Sigma_{\theta^*, m^*} := F(\theta^*)^{-1} G(\theta^*) F(\theta^*)^{-1}$

The proof of this result is identical to the proofs of theorems 3.1 and 3.2 of our previous paper [Bardet et al. \(2020b\)](#). It will therefore not be carried out for the sake of brevity.

In this subsection, we have used the QMLE contrast without any distribution assumption on the noise to derive a consistency property. However, the contrast  $L_n$  as in (3.2.1) depends on all the past values of the process  $X$ , which are unobserved. In the next subsection, we will propose an extension of Theorem 3.1 based on QMLE which does not require knowledge of the initial values of the process.

### 3.2.4 Another Quasi-Likelihood Function

The goal of this subsection is to sharpen the conditions on the sequence  $\kappa_n$  found in [Bardet et al. \(2020b\)](#). Before stating the result, let recall a little bit some definitions and notations used in [Bardet et al. \(2020b\)](#). Following the derivation of  $L_n$ , we define its computable version  $\widehat{L}_n$  as follows:

$$\widehat{L}_n(\theta) := -\frac{1}{2} \sum_{t=1}^n \widehat{q}_t(\theta), \text{ with } \widehat{q}_t(\theta) := \frac{(X_t - \widehat{f}_\theta^t)^2}{\widehat{H}_\theta^t} + \log(\widehat{H}_\theta^t) \quad (3.2.12)$$

where  $\widehat{f}_\theta^t := f_\theta(X_{t-1}, X_{t-2}, \dots, X_1, 0, \dots, 0)$ ,  $\widehat{M}_\theta^t := M_\theta(X_{t-1}, X_{t-2}, \dots, X_1, 0, \dots, 0)$  and  $\widehat{H}_\theta^t = (\widehat{M}_\theta^t)^2$ .

Therefore for every model  $m \in \mathcal{M}$ , we define the Gaussian QMLE  $\widetilde{\theta}(m)$  (with respect to  $\widehat{L}_n$ ) as

$$\widetilde{\theta}(m) = \underset{\theta \in \Theta(m)}{\operatorname{argmax}} \widehat{L}_n(\theta). \quad (3.2.13)$$

Once the estimation in each model in the family  $\mathcal{M}$  has been performed, we select the best model as follows

$$\widehat{m} = \underset{m \in \mathcal{M}}{\operatorname{argmin}} \widehat{C}(m) \quad \text{with} \quad \widehat{C}(m) = -2 \widehat{L}_n(\widetilde{\theta}(m)) + \kappa_n(m). \quad (3.2.14)$$

In this framework, we do not consider long memory process and then we define the class  $\mathcal{H}(M_\theta, f_\theta)$  a subset of  $\mathcal{AC}(M_\theta, f_\theta)$  in which every process has Lipschitz coefficients satisfying the following conditions

$$\alpha_j(f_\theta) + \alpha_j(M_\theta) + \alpha_j(\partial_\theta f_\theta) + \alpha_j(\partial_\theta M_\theta) = O(j^{-\gamma}) \quad \text{with} \quad \gamma > 2.$$

It is then straightforward to see that every process in the class  $\mathcal{H}(M_\theta, f_\theta)$  verifies the following condition

**Condition  $K(\Theta)$ :**

$$\sum_{k \geq e} \frac{1}{\log \log k} \sum_{j \geq k} \alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) + \alpha_j(\partial_\theta f_\theta, \Theta) + \alpha_j(\partial_\theta M_\theta, \Theta) < \infty.$$

This finding allows to propose a sharpen generalization of both Theorem 3.1 in our previous paper [Bardet et al. \(2020b\)](#) and a similar result in [Kengne \(2020\)](#).

Before stating the consistency result, it is important as in [Bardet et al. \(2020b\)](#) or [Bardet and Wintenberger \(2009\)](#), to distinguish the special case of NLARCH( $\infty$ ) processes which includes for instance GARCH( $p, q$ ) or ARCH( $\infty$ ) processes. In such case, let us define the class:

**Class  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$ :** A process  $X = (X_t)_{t \in \mathbb{Z}}$  belongs to  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$  if it satisfies:

$$X_t = \xi_t \sqrt{\widetilde{H}_\theta((X_{t-i}^2)_{i \in \mathbb{N}^*})} \quad \text{for any } t \in \mathbb{Z}. \quad (3.2.15)$$

Therefore, if  $M_\theta^2((X_{t-i})_{i \in \mathbb{N}^*}) = H_\theta((X_{t-i})_{i \in \mathbb{N}^*}) = \widetilde{H}_\theta((X_{t-i}^2)_{i \in \mathbb{N}^*})$  then,  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta) = \mathcal{AC}(M_\theta, 0)$ . In case of the class  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$ , we will use the assumption  $A(\widetilde{H}_\theta, \Theta)$ . The new set of stationary solutions is for  $r \geq 2$ :

$$\widetilde{\Theta}(r) = \left\{ \theta \in \mathbb{R}^d, A(\widetilde{H}_\theta, \{\theta\}) \text{ holds with } (\|\xi_0\|_r)^2 \sum_{k=1}^{\infty} \alpha_k(\widetilde{H}_\theta, \{\theta\}) < 1 \right\}. \quad (3.2.16)$$

Finally, we propose to restrict class  $\widetilde{\mathcal{AC}}(\widetilde{H}_\theta)$  to  $\widetilde{\mathcal{H}}(\widetilde{H}_\theta)$  as done with  $\mathcal{H}(M_\theta, f_\theta)$  by considering all the process checking the condition

$$\alpha_j(H_\theta) + \alpha_j(\partial_\theta H_\theta) = O(j^{-\gamma}) \quad \text{with } \gamma > 2$$

so that

$$\sum_{k \geq e} \frac{1}{\log \log k} \sum_{j \geq k} \alpha_j(H_\theta) + \alpha_j(\partial_\theta H_\theta) < \infty.$$

We can now state the second main result.

**Theorem 3.3.** *Let  $(X_1, \dots, X_n)$  be an observed trajectory of an affine causal process  $X$  belonging to  $\mathcal{H}(M_{\theta^*}, f_{\theta^*})$  (or  $\widetilde{\mathcal{H}}(\widetilde{H}_{\theta^*})$ ) where  $\theta^*$  is an unknown vector belonging to  $\Theta(r) \subset \mathbb{R}^{D_{m^*}}$  (or  $\widetilde{\Theta}(r) \subset \mathbb{R}^{D_{m^*}}$ ). Let also  $\mathcal{M}$  be a finite family of candidate models such that  $m^* \in \mathcal{M}$ . If assumptions **A1-A5** hold, then with  $c_{\min} := \max(\frac{\alpha_1}{2}, \frac{\alpha_2}{2})$ , it holds*

for any  $\kappa_n(m) = 2c D_m \log \log n$  with

$$c \geq c_{\min} \tag{3.2.17}$$

it holds for the selected model  $\widehat{m}$  according to (3.2.14)

$$\widehat{m} \xrightarrow[n \rightarrow \infty]{a.s.} m^*. \tag{3.2.18}$$

All the comments made about the Theorem 3.1 remain valid here. Moreover, recently, [Kengne \(2020\)](#) requires heavy penalties to ensure the strong consistency for the process in the class  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$ . Indeed, according to [Kengne \(2020\)](#), it is necessary that  $\kappa_n$  verified  $\kappa_n / \log \log n \xrightarrow[n \rightarrow \infty]{} \infty$  to obtain (3.2.18) which is a stronger condition since the HQ criterion does not fulfill this condition and it is well known that HQ criterion is strongly consistent (see for instance [Hannan and Quinn \(1979\)](#)). Moreover, the new penalties found in this paper does not satisfy this condition, yet we ensure strongly consistency.

Also, let us mention that for small samples, heavy penalties such as those in [Kengne \(2020\)](#) can very often choose very simple models possible wrongs [Hannan and Quinn \(1979\)](#).

**Remark 3.2.** *Our results show that, asymptotically heavy penalties such as BIC penalty will ensure the consistency property. However, in practise, these penalties are used for fixed  $n$  and most for small samples and it is important to point out their drawbacks. Very often, in the family of competitive models  $\mathcal{M}$ , there are misspecified and underfitted models ( $\mathcal{M}'$ ). Since the difference  $-2\widehat{L}_n(\widehat{\theta}_m) + 2\widehat{L}_n(\widehat{\theta}_{m^*}) > 0$  for every  $m \in \mathcal{M}'$ , making the penalty, heavy could offset this positivity and can lead to the selection of some underfitted models and then wrong models. To be more convincing of that, see the simulations (DGP III) experiments in Section 3.3.*

### 3.2.5 Algorithm of Calibration of the minimal constant

There exist several ways to calibrate the minimal constant  $c_{\min}$  including the dimension jump (presented below) and the data-driven slope estimation. Indeed, once an estimation of  $c_{\min}$  is obtained, many studies advocates the choice of  $2\widehat{c}_{\min}$  which turn out to be optimal ([Massart \(2007\)](#), [Arlot and Massart \(2009\)](#) among others). Now we present the dimension jump algorithm.

**Algorithm 3.1. Dimension Jump Arlot and Massart (2009)**

1. Compute the selected model  $\hat{m}(c)$  as a function of  $c > 0$

$$\hat{m}(c) = \operatorname{argmin}_{m \in \mathcal{M}} \hat{L}_n(\hat{\theta}_m) + c \operatorname{pen}_{\text{shape}}(m); \quad \operatorname{pen}_{\text{shape}}(m) = D_m \log \log n$$

2. Find  $\hat{c}_{\min}$  such that  $D_{\hat{m}(c)}$  is "huge" for  $c < \hat{c}_{\min}$  and "reasonnably small" for  $c \geq \hat{c}_{\min}$ ;

3. Select the model  $\hat{m} := \hat{m}(2\hat{c}_{\min})$ .

This algorithm has been implemented in Baudry et al. (2012) which gives several details including the grid for  $c$  values.

Let us notice that, in view of obtaining penalties, there is no need to calibrate the  $c_{\min}$  constant for most classical time series models. However, since the fourth order moment of the noise is unknown, a consistent estimate of this term is required. To do that, we proceed as in the estimation of the variance of the noise as in the Mallows Cp.

A consistent estimator  $\hat{\mu}_{\bar{m},4}$  of  $\mu_4 = \mathbb{E}[\xi_0^4]$  can be :

$$\hat{\mu}_{\bar{m},4} := \frac{1}{n} \sum_{t=1}^n (\hat{\xi}_t(\bar{m}))^4 \quad \text{with} \quad \hat{\xi}_t(\bar{m}) := (\widehat{M}_{\bar{m}}^t)^{-1} (X_t - \hat{f}_{\bar{m}}^t)$$

where we suppose that  $\bar{m}$  is the "largest" model in the family  $\mathcal{M}$ , typically the largest order of a family of time series. As a result an estimator of the  $c_{\min}$  constant to consider in the penalty  $\kappa_n$  is

- $\frac{\hat{\mu}_{\bar{m},4}-1}{2}$  for GARCH family and related ones;
- 1 for ARMA family with known variance.

### 3.3 NUMERICAL EXPERIMENTS

In this section, several numerical experiments are conducted to assess the consistency property (Section 3.2) of our new criteria.

#### 3.3.1 Monte Carlo: Consistency

This subsection studies the performance of the model selection criteria found in Section 3.2. We have considered three different Data Generating Process (DGP):

$$\begin{aligned} \text{DGP I} \quad & X_t = 0.5 X_{t-1} + 0.2 X_{t-2} + \xi_t, \\ \text{DGP II} \quad & X_t = (0.2 + 0.4 X_{t-1}^2 + 0.2 X_{t-2}^2)^{1/2} \xi_t, \\ \text{DGP III} \quad & X_t = 0.1 (X_{t-1} + X_{t-2} + \dots + X_{t-6}) + \xi_t, \end{aligned}$$

where  $(\xi_t)$  will be a white Gaussian noise with variance one at first and a Student with 5 degrees of freedom on the other hand. For the first and the second DGP, we considered as competitive models all the models in the family  $\mathcal{M}$  defined by:

$$\mathcal{M} = \{ \text{ARMA}(p, q) \text{ or GARCH}(p', q') \text{ processes with } 0 \leq p, q, p' \leq 5, 1 \leq q' \leq 5 \}.$$

Therefore, there are 66 candidate models as in Bardet et al. (2020b). The goal is to compare the ability of selecting the true model for  $\kappa_n^3 = \log n D_m$ ,  $\kappa_n^2 = 2\hat{c}_{\min} \log \log n D_m$

( in accordance with the condition (1.3.12) and  $\kappa_n^2 = 2 \times 2 \hat{c}_{min} \log \log n D_m$ . Moreover, from Proposition 1.4,  $\hat{c}_{min}$  does not need to be estimated and worth one for Gaussian noise. But for Student noise,  $\hat{c}_{min} = \max(1, \frac{\hat{\mu}_{m,4}-1}{2})$  for the DGP I and  $\hat{c}_{min} = \frac{\hat{\mu}_{m,4}-1}{2}$  for DGP II. The Table 2.1 presents the results of the selection procedure. As we can notice, the three penalties have a good consistency property. Moreover, for  $n$  relatively small, the penalty  $\kappa_n^1$  is better than both others. For larger  $n$ ,  $\kappa_n^2$  is the best the penalty to consider.

**Table 3.1:** Percentage of selected order based on 500 independent replications depending on sample's length for the penalty  $\kappa_n^1 = 2 \hat{c}_{min} \log \log n D_m$ ,  $\kappa_n^2 = 4 \hat{c}_{min} \log \log n D_m$  and  $\kappa_n^3 = \log n D_m$ , where W, T, O refers to wrong, true and overfitted selection.

$n$		100			500			1000			2000		
		$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$
DGP I Gaussian	W	58.2	80.8	94.6	23.6	26	36.2	18.8	17.6	17.8	4.2	7.4	7.8
	T	32.6	19	5.4	64.2	73.4	63.8	71.6	82.2	82.0	89.4	92.2	92.2
	O	9.2	0.2	0	12.2	0.6	0	9.6	0.2	0.2	6.4	0.5	0
DGP II Gaussian	W	70.0	84.8	94.0	17.4	21.6	35.4	16.4	15.8	15.8	3.6	3.6	5
	T	29.8	15.2	6.0	81.0	78.4	64.6	83.4	84.2	84.2	96.4	96.4	95
	O	0.2	0	0	1.6	0	0	0.2	0	0.0	0	0	0
DGP I Student	W	70.1	86.8	97.5	38.3	53.1	35.0	16.0	19.5	16.3	18	22.0	18
	T	19.8	13.2	2.5	59.7	46.9	65.9	83.9	80.5	82.7	81.4	78.0	82.0
	O	10.1	0	0	2.0	0	0.1	0.1	0	0	0.6	0	0
DGP II Student	W	75.0	84.2	88.8	45.6	57.2	55.4	25.6	32.6	26.2	14	16.0	13.0
	T	22.4	15.0	11.2	49.0	41.2	44.6	71.0	66.8	73.8	85.0	84.0	87.0
	O	2.6	0.8	0.0	4.4	1.6	1.0	3.4	0.6	0.0	1.0	0	0

For DGP III, as we want to exhibit the possible "non consistency" of BIC for small samples, we have considered as the competitive set, the hierarchical family of AR models

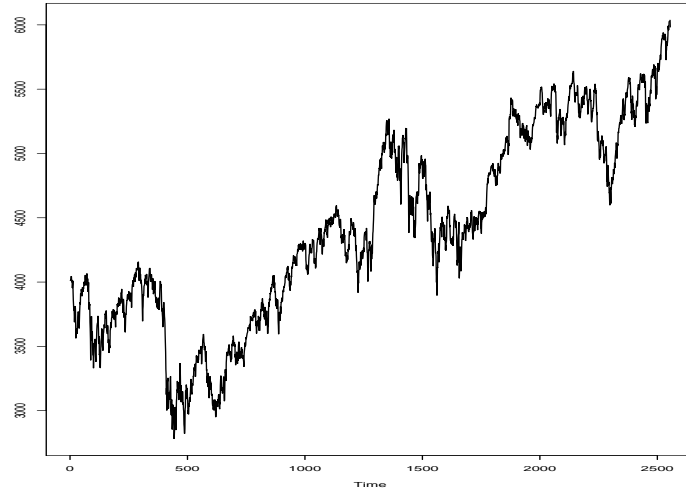
$$\mathcal{M}' = \{AR(1), \dots, AR(15)\}$$

for  $n = 100, 200, 400, 500, 1000, 2000$ . In Table 3.2 below the percentage of selected order based on 1000 independent replications are presented for the above three penalties.

**Table 3.2:** Percentage of selected order based on 1000 independent replications depending on sample's length using penalty terms  $\kappa_n^1 = \hat{c}_{min}$ ,  $\kappa_n^2 = 2 \hat{c}_{min}$  and  $\kappa_n^3 = \log n$ , for DGP III.

$n$	100			200			400			500			1000			2000		
	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$	$\kappa_n^1$	$\kappa_n^2$	$\kappa_n^3$
$p < 6$	61.5	96	99	46	83.5	97	29.5	59	87.5	19.5	45.5	76.5	3.0	9.5	38.5	0	0.5	5.0
$p = 6$	15	3.0	1.0	25	13	2.5	44	37.5	12.5	50.5	49.5	22.5	68	86.5	61	72.5	96	94.5
$p > 6$	23.5	1.0	0	29	3.5	0.5	26.5	3.5	0	30	5.0	1.0	29	4.0	0.5	27.5	3.5	0.5

These results invite us to be cautious when using the BIC for small sample sizes, whereas the proposed adaptive penalty is more robust, as it at least allow us to recover an overfitted model that is less harmful than a wrong model most often chosen by the BIC.



**Figure 3.1:** Daily closing CAC 40 index (January 1st, 2010 to December 31 st, 2019).

### 3.3.2 Real Data Analysis: financial time series

CAC 40 is a benchmark French stock market index and is highly regarded in many statistical studies . Let consider the daily closing prices of the CAC 40 index from January 1st, 2010 to December 31st, 2019 plotted in Figure 3.1. Over the period under review, the CAC 40 increased.

To analyze this type of data, it is common to consider the returns (see Figure 3.2). We can see that the return values display some small auto-correlations. Also, from Figure 3.3, the squared returns of CAC 40 are strongly auto-correlated. These facts suggest that the strong white noise assumption cannot be sustained for this log-returns sequence of the CAC 40 index.

Hence, let consider the competitive set of models  $\mathcal{M}$  used in the previous subsection in order to propose the best suitable model for these data:

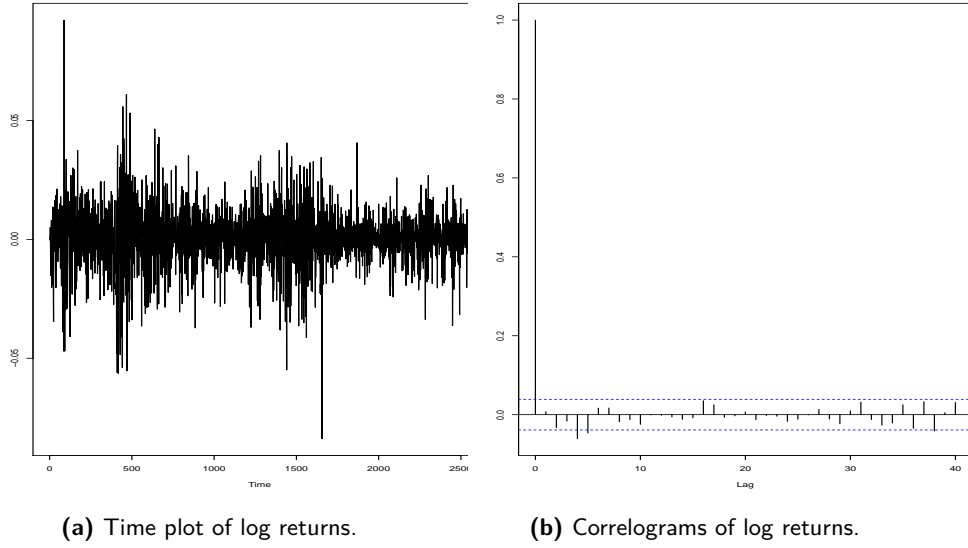
$$\mathcal{M} = \{ \text{ARMA}(p, q) \text{ or GARCH}(p', q') \text{ processes with } 0 \leq p, q, p' \leq 5, 1 \leq q' \leq 5 \}.$$

Using the adaptive penalty and the BIC criterion, we find out that the GARCH(1, 1) is the best model over  $\mathcal{M}$  with respect to both criteria. This fact is in accordance with [Francq and Zakoian \(2010\)](#) which found the same result using the returns of the CAC 40 index from March 2, 1990 to December 29, 2006.

## 3.4 Proofs

### 3.4.1 Proof of Theorem 3.1

Usually, this proof is divided into two parts: one has to show that as  $n$  tends to infinity the probability of overfitting goes to zero and so is the probability of misspecification.



**Figure 3.2:** Daily closing CAC 40 index (January 4th, 2010 to December 31 st, 2019).

### 3.4.1.a Overfitting Case

Let  $m \in \mathcal{M}$  such as  $m^* \subset m$ . We want to show that  $C(m^*) \leq C(m)$  a.s. asymptotically. We have

$$\begin{aligned}
 C(m^*) \leq C(m) &\iff -2L_n(\hat{\theta}(m^*)) + 2c \log \log n D_{m^*} \leq -2L_n(\hat{\theta}(m)) + 2c \log \log n D_m \\
 &\iff \frac{L_n(\hat{\theta}(m)) - L_n(\theta^*)}{\log \log n} \leq \frac{L_n(\hat{\theta}(m^*)) - L_n(\theta^*)}{\log \log n} + c(D_m - D_{m^*}) \quad (3.4.1)
 \end{aligned}$$

therefore, a necessary and sufficient condition to avoid overfitting can be stated by taking  $\limsup_{n \rightarrow \infty}$  on both sides of the inequality (3.4.1); that is by virtue of definition (1.3.11)

$$2c_{\min}(m) - 2c_{\min}(m^*) \leq c(D_m - D_{m^*}) \quad \text{for } m^* \subset m, \quad (3.4.2)$$

i.e.,

$$\frac{\alpha_1}{2} (D_m^1 - D_{m^*}^1) + \frac{\alpha_2}{2} (D_m^2 - D_{m^*}^2) \leq c(D_m - D_{m^*})$$

which is fulfilled for any constant  $c$  such as in (3.2.8). Indeed,  $c_{\min} = \max(\frac{\alpha_1}{2}, \frac{\alpha_2}{2})$  and

$$\begin{aligned}
 \frac{\alpha_1}{2} (D_m^1 - D_{m^*}^1) + \frac{\alpha_2}{2} (D_m^2 - D_{m^*}^2) &\leq c_{\min} (D_m^1 - D_{m^*}^1 + D_m^2 - D_{m^*}^2) \\
 &= c_{\min} (D_m - D_{m^*})
 \end{aligned}$$

where the inequality holds since  $m^* \subset m$  that implies  $D_m^1 \geq D_{m^*}^1$  and  $D_m^2 \geq D_{m^*}^2$ . Hence the associated criterion  $\kappa_n$  will avoid overfitting.

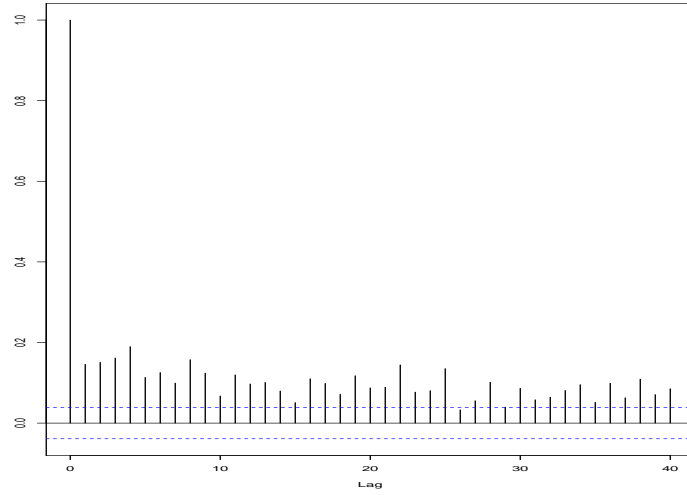
### 3.4.1.b Misspecification/Underfitting Case

All misspecified/underfitted models are contained in the set

$$\mathcal{M}' = \left\{ m \in \mathcal{M} : (m^* \not\subset m) \cup (m \subset m^*) \right\}.$$

The proof is exactly as the one done in [Bardet et al. \(2020b\)](#). But for the sake of completeness, we give here some important steps of the proof.





**Figure 3.3:** Sample autocorrelations of squared returns of the CAC 40 index (January 1st, 2010 to December 31 st, 2019).

The goal is to show that for every  $m \in \mathcal{M}'$ , it holds

$$C(m^*) < C(m) \quad a.s. \quad (3.4.3)$$

Let  $m \in \mathcal{M}'$ . From Proposition 2 in [Bardet et al. \(2020b\)](#) and using continuous mapping Theorem

$$\frac{1}{n} \left[ L_n(\hat{\theta}(m^*)) - L_n(\hat{\theta}(m)) \right] = A(m_0) + o_{a.s.}(1) \quad (3.4.4)$$

where

$$A(m) := L(\theta^*) - L(\theta_m^*) \quad \text{with} \quad L(\theta) = -\frac{1}{2} \mathbb{E}[q_0(\theta)].$$

Let us denote by  $\mathcal{F}_t := \sigma(X_{t-1}, X_{t-2}, \dots)$ . Using conditional expectation, we obtain

$$L(\theta^*) - L(\theta) = \frac{1}{2} \mathbb{E} \left[ \mathbb{E}[q_0(\theta) - q_0(\theta^*) \mid \mathcal{F}_0] \right]. \quad (3.4.5)$$

But,

$$\begin{aligned} \mathbb{E}[q_0(\theta) - q_0(\theta^*) \mid \mathcal{F}_0] &= \mathbb{E} \left[ \frac{(X_0 - f_\theta^0)^2}{H_\theta^0} + \log(H_\theta^0) - \frac{(X_0 - f_{\theta^*}^0)^2}{H_{\theta^*}^0} - \log(H_{\theta^*}^0) \mid \mathcal{F}_0 \right] \\ &= \log \left( \frac{H_\theta^0}{H_{\theta^*}^0} \right) + \frac{\mathbb{E}[(X_0 - f_\theta^0)^2 \mid \mathcal{F}_0]}{H_\theta^0} - \frac{\mathbb{E}[(X_0 - f_{\theta^*}^0)^2 \mid \mathcal{F}_0]}{H_{\theta^*}^0} \\ &= \log \left( \frac{H_\theta^0}{H_{\theta^*}^0} \right) - 1 + \frac{\mathbb{E}[(X_0 - f_{\theta^*}^0 + f_{\theta^*}^0 - f_\theta^0)^2 \mid \mathcal{F}_0]}{H_\theta^0} \\ &= \frac{H_{\theta^*}^0}{H_\theta^0} - \log \left( \frac{H_{\theta^*}^0}{H_\theta^0} \right) - 1 + \frac{(f_{\theta^*}^0 - f_\theta^0)^2}{H_\theta^0} \end{aligned}$$

Thus from (3.4.5),

$$\begin{aligned} 2 A(m) &= \mathbb{E} \left[ \frac{H_{\theta^*}^0}{H_{\theta_m^*}^0} - \log \left( \frac{H_{\theta^*}^0}{H_{\theta_m^*}^0} \right) - 1 + \frac{(f_{\theta^*}^0 - f_{\theta_m^*}^0)^2}{H_{\theta_m^*}^0} \right] \\ &\geq \mathbb{E} \left[ \frac{H_{\theta^*}^0}{H_{\theta_m^*}^0} \right] - \log \left( \mathbb{E} \left[ \frac{H_{\theta^*}^0}{H_{\theta_m^*}^0} \right] \right) - 1 + \mathbb{E} \left[ \frac{(f_{\theta^*}^0 - f_{\theta_m^*}^0)^2}{H_{\theta_m^*}^0} \right] \quad \text{by Jensen Inequality.} \end{aligned}$$

Since  $x - \log(x) - 1 > 0$  for any  $x > 0$ ,  $x \neq 1$  and  $x - \log(x) - 1 = 0$  for  $x = 1$ , we deduce that

- If  $f_{\theta^*}^0 \neq f_{\theta_m^*}^0$ , then  $\mathbb{E}\left[\frac{(f_{\theta^*}^0 - f_{\theta_m^*}^0)^2}{H_{\theta_m^*}^0}\right] > 0$  and  $A(m) > 0$ .

- Otherwise, then

$$2A(m) = \mathbb{E}\left[\frac{H_{\theta^*}^0}{H_{\theta_m^*}^0} - \log\left(\frac{H_{\theta^*}^0}{H_{\theta_m^*}^0}\right) - 1\right],$$

and from assumption **A1**, since  $\theta^* \notin \Theta(m)$  and  $f_{\theta^*}^0 = f_{\theta_m^*}^0$ , we necessarily have  $H_{\theta^*}^0 \neq H_{\theta_m^*}^0$  so that  $\frac{H_{\theta^*}^0}{H_{\theta_m^*}^0} \neq 1$ . Then  $A(m) > 0$ .

As a consequence,

$$\frac{C(m) - C(m^*)}{n} = A(m) + \frac{2c\kappa_n}{n}(D_m - D_{m^*}) + o_{a.s.}(1).$$

That establishes (3.4.3) by virtue of (3.2.8) and as all the considered models are finite dimensional.  $\square$

In the sequel, we state and prove several lemmas to which we referred to when proving above main results.

#### 3.4.1.c Proof of Theorem 3.3

The proof follows mutatis mutandis from the Theorem 3.1's proof after replacing line by line  $L_n$ ,  $\hat{\theta}(m)$  and the criterion  $C$  by their equivalent  $\hat{L}_n$ ,  $\hat{\theta}(m)$  and  $\hat{C}$  respectively and applying Lemma 3.2 instead of Proposition 3.1.  $\square$

#### 3.4.1.d Proof of Proposition 3.1

*Proof.* Applying a second order Taylor expansion of  $L_n$  around  $\hat{\theta}(m)$  for  $n$  sufficiently large such that  $\bar{\theta}(m) \in \Theta(m)$  which are between  $\theta_m^* := (\theta^*, 0, \dots, 0)^\top$  and  $\hat{\theta}(m)$  yields (as  $\partial L_n(\hat{\theta}(m)) = 0$  since  $\hat{\theta}(m)$  is a local extremum):

$$\begin{aligned} L_n(\theta^*) - L_n(\hat{\theta}(m)) &= L_n(\theta_m^*) - L_n(\hat{\theta}(m)) \\ &= \frac{1}{2}(\hat{\theta}(m) - \theta_m^*)^\top \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta^2} (\hat{\theta}(m) - \theta_m^*) := I_1(m). \end{aligned}$$

From the mean value Theorem, and for large  $n$ , there exists  $\bar{\theta}_{m,i}$  between  $(\theta_m^*)_i$  and  $(\hat{\theta}(m))_i$  such that,  $1 \leq i \leq D_m$ :

$$0 = \frac{\partial L_n(\hat{\theta}_m)}{\partial \theta_i} = \frac{\partial L_n(\theta_m^*)}{\partial \theta_i} + \frac{\partial^2 L_n(\bar{\theta}_{m,i})}{\partial \theta \partial \theta_i} (\hat{\theta}_m - \theta_m^*) \quad (3.4.6)$$

Also, using Lemma 4 of [Bardet and Wintenberger \(2009\)](#) and continuous mapping Theorem, we deduce that:

$$F_n := -\left(\frac{2}{n} \frac{\partial^2 L_n(\bar{\theta}_{m,i})}{\partial \theta \partial \theta_i}\right)_{1 \leq i \leq D_m} \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta_m^*) = \mathbb{E}\left[\frac{\partial^2 q_0(\theta_m^*)}{\partial \theta^2}\right]. \quad (3.4.7)$$

On the other hand, under **A2** condition,  $F(\theta_m^*)$  is an invertible matrix and there exists  $n$  sufficiently large such that  $F_n$  is invertible. Therefore, from (3.4.6), it follows

$$\begin{aligned}
\frac{I_1(m)}{2 \log \log n} &= \frac{1}{4 \log \log n} (\hat{\theta}(m) - \theta_m^*)^\top \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta^2} (\hat{\theta}(m) - \theta_m^*) \\
&= \frac{1}{4 \log \log n} \left( \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)^\top \left( -\frac{2}{n} F_n^{-1} \right) \frac{\partial^2 L_n(\bar{\theta}(m))}{\partial \theta^2} \left( -\frac{2}{n} F_n^{-1} \right) \frac{\partial L_n(\theta_m^*)}{\partial \theta} \\
&= -\frac{1}{2n \log \log n} \left( \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)^\top \times F(\theta_m^*)^{-1} \times \frac{\partial L_n(\theta_m^*)}{\partial \theta} (1 + o(1)) \quad \text{a.s.}
\end{aligned}$$

The next step of the proof consists in handling the quadratic form  $\left( \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)^\top \times F(\theta_m^*)^{-1} \times \left( \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)$  by applying the law of iterated logarithm (LIL).

We claim that

$$\limsup_{n \rightarrow \infty} \frac{1}{\sqrt{2n \log \log n}} 2 G(\theta_m^*)^{-1/2} \frac{\partial L_n(\theta_m^*)}{\partial \theta} = (1, \dots, 1)^\top. \quad (3.4.8)$$

**Proof of the claim:** First, since the covariance matrix of  $2 \frac{\partial L_n(\theta_m^*)}{\partial \theta}$  is  $G(\theta_m^*)$ , it follows that the covariance matrix of the vector  $Z_m := 2 G(\theta_m^*)^{-1/2} \frac{\partial L_n(\theta_m^*)}{\partial \theta}$  is the  $D_m \times D_m$  identity matrix. Moreover, as

$$\frac{\partial L_n(\theta_m^*)}{\partial \theta} = -\frac{1}{2} \sum_{t=1}^n \frac{\partial q_t(\theta_m^*)}{\partial \theta} = -\frac{1}{2} \sum_{t=1}^n \frac{\partial q_t(\theta^*)}{\partial \theta},$$

where

$$\mathbb{E} \left[ \frac{\partial q_t(\theta^*)}{\partial \theta} \middle| \sigma(X_{t-1}, X_{t-2}, \dots, X_1) \right] = 0 \quad (3.4.9)$$

and

$$\mathbb{E} \left[ \left( \frac{\partial q_1(\theta^*)}{\partial \theta_i} \right)^2 \right] < \infty \quad (3.4.10)$$

hold from [Bardet and Wintenberger \(2009\)](#) under **A4**. Finally, one can see that the  $i^{th}$  element of  $Z_m$  can be rewritten as

$$\sum_{j=1}^{D_m} (2 G(\theta_m^*)^{-1/2})_{i,j} \frac{\partial L_n(\theta_m^*)}{\partial \theta_j} = \sum_{t=1}^n \zeta_t^i$$

where  $\zeta_t^i = -\sum_{j=1}^{D_m} (G(\theta_m^*)^{-1/2})_{i,j} \frac{\partial q_t(\theta^*)}{\partial \theta_j}$ . By virtue of (3.4.9), we have

$$\mathbb{E} \left[ \zeta_t^i \middle| \sigma(X_{t-1}, X_{t-2}, \dots, X_1) \right] = 0.$$

Hence, any component of  $Z_m$  verifies the LIL. That is for any  $i = 1, \dots, D_m$ ,

$$\limsup_{n \rightarrow \infty} \frac{1}{\sqrt{2n \log \log n}} \left( 2 G(\theta_m^*)^{-1/2} \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)_i = 1.$$

This fact concludes the proof of the claim (3.4.8).

Hence writing

$$\begin{aligned}
&\frac{1}{2n \log \log n} \frac{\partial L_n(\theta_m^*)}{\partial \theta}^\top F(\theta_m^*)^{-1} \frac{\partial L_n(\theta_m^*)}{\partial \theta} = \\
&\left( \frac{1}{\sqrt{2n \log \log n}} 2 G(\theta_m^*)^{-1/2} \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)^\top \frac{G(\theta_m^*)^{1/2} F(\theta_m^*)^{-1} G(\theta_m^*)^{1/2}}{4} \left( \frac{1}{\sqrt{2n \log \log n}} 2 G(\theta_m^*)^{-1/2} \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)
\end{aligned}$$

it follows

$$\limsup_{n \rightarrow \infty} \frac{L_n(\hat{\theta}(m)) - L_n(\theta_m^*)}{2 \log \log n} = \mathbf{1}_m^\top \Sigma_{\theta_m^*} \mathbf{1}_m.$$

□

### 3.4.1.e Proof of Proposition 3.2

It is sufficient to show that

$$\Sigma_{\theta_m^*} := G(\theta_m^*)^{1/2} F(\theta_m^*)^{-1} G(\theta_m^*)^{1/2} = \begin{pmatrix} 2\mathbf{I}_{D_m^1, D_m^1} & \mathbf{O}_{D_m^1, D_m^2} \\ \mathbf{O}_{D_m^2, D_m^1} & (\mu_4 - 1)\mathbf{I}_{D_m^2, D_m^2} \end{pmatrix}, \quad (3.4.11)$$

where  $\mathbf{I}$  is the identity matrix and  $\mathbf{O}$  the null matrix. From [Bardet and Wintenberger \(2009\)](#), we have for  $m^* \subset m$  and  $i, j \in m$ :

$$\begin{aligned} (G(\theta_m^*))_{i,j} &= \mathbb{E} \left[ \frac{\partial q_0(\theta_m^*)}{\partial \theta_i} \frac{\partial q_0(\theta_m^*)}{\partial \theta_j} \right] = \mathbb{E} \left[ 4(H_{\theta_m^*}^0)^{-1} \frac{\partial f_{\theta_m^*}^0}{\partial \theta_i} \frac{\partial f_{\theta_m^*}^0}{\partial \theta_j} + (\mu_4 - 1)(H_{\theta_m^*}^0)^{-2} \frac{\partial H_{\theta_m^*}^0}{\partial \theta_i} \frac{\partial H_{\theta_m^*}^0}{\partial \theta_j} \right] \\ (F(\theta_m^*))_{i,j} &= \mathbb{E} \left[ \frac{\partial^2 q_0(\theta_m^*)}{\partial \theta_i \partial \theta_j} \right] = \mathbb{E} \left[ 2(H_{\theta_m^*}^0)^{-1} \frac{\partial f_{\theta_m^*}^0}{\partial \theta_i} \frac{\partial f_{\theta_m^*}^0}{\partial \theta_j} + (H_{\theta_m^*}^0)^{-2} \frac{\partial H_{\theta_m^*}^0}{\partial \theta_i} \frac{\partial H_{\theta_m^*}^0}{\partial \theta_j} \right] \end{aligned} \quad (3.4.12)$$

1/ If  $\mu_4 = 3$ , then  $G(\theta_m^*) = 2F(\theta_m^*)$  and the result is straightforward.

2/ For the second configuration: from (4.6.1), it is clear that all the terms  $F(\theta_m^*)_{i,j}$  and  $G(\theta_m^*)_{i,j}$  are equals to zero for  $i = 1, \dots, D_m^1$  and  $j = 1, \dots, D_m^2$  implying

$$F(\theta_m^*) = \begin{pmatrix} A_{1, D_m^1} & \mathbf{O}_{D_m^1, D_m^2} \\ \mathbf{O}_{D_m^2, D_m^1} & B_{D_m^1+1, D_m^1+D_m^2} \end{pmatrix} \text{ and } G(\theta_m^*) = \begin{pmatrix} 2A_{1, D_m^1} & \mathbf{O}_{D_m^1, D_m^2} \\ \mathbf{O}_{D_m^2, D_m^1} & (\mu_4 - 1)B_{D_m^1+1, D_m^1+D_m^2} \end{pmatrix}$$

and

$$\begin{aligned} A_{1, D_m^1} &= \left( \mathbb{E} \left[ 2(H_{\theta_m^*}^t)^{-1} \frac{\partial f_{\theta_m^*}^t}{\partial \theta_i} \frac{\partial f_{\theta_m^*}^t}{\partial \theta_j} \right] \right)_{1 \leq i, j \leq D_m^1} \\ B_{D_m^1+1, D_m^1+D_m^2} &= \left( \mathbb{E} \left[ (H_{\theta_m^*}^t)^{-2} \frac{\partial H_{\theta_m^*}^t}{\partial \theta_i} \frac{\partial H_{\theta_m^*}^t}{\partial \theta_j} \right] \right)_{D_m^1+1 \leq i, j \leq D_m^1+D_m^2}. \end{aligned}$$

As a consequence,

$$\begin{aligned} G(\theta_m^*) F(\theta_m^*)^{-1} &= \text{Diag}(2A_{1, D_m^1}, (\mu_4 - 1)B_{D_m^1+1, D_m^1+D_m^2}) \text{Diag}(A_{1, D_m^1}^{-1}, B_{D_m^1+1, D_m^1+D_m^2}^{-1}) \\ &= \text{Diag}(2\mathbf{I}_{D_m^1}, (\mu_4 - 1)\mathbf{I}_{D_m^2}). \end{aligned}$$

As a covariance matrix,  $G(\theta_m^*)$  is positive definite. Therefore the square root  $G(\theta_m^*)^{1/2}$  of  $G(\theta_m^*)$  is unique and blocks diagonal. Thus,

$$\begin{aligned} \Sigma_{\theta_m^*} &= G(\theta_m^*)^{-1/2} \left( G(\theta_m^*) F(\theta_m^*)^{-1} \right) G(\theta_m^*)^{1/2} \\ &= G(\theta_m^*) F(\theta_m^*)^{-1}, \end{aligned}$$

which gives (3.4.11). □

### 3.4.2 Technical Lemmas

**Lemma 3.1.** *Let  $X \in \mathcal{H}(M_{\theta^*}, f_{\theta^*})$  (or  $\tilde{\mathcal{H}}(H_{\theta^*})$ ) and  $\Theta \subseteq \Theta(r)$  (or  $\Theta \subseteq \tilde{\Theta}(r)$ ) with  $r \geq 4$ . Assume that assumption **A3** holds. Then for  $i = 0, 1, 2$ , it holds*

$$\frac{1}{\log \log n} \left\| \frac{\partial^{(i)} \hat{L}_n(\theta)}{\partial \theta^i} - \frac{\partial^{(i)} L_n(\theta)}{\partial \theta^i} \right\|_{\Theta} \xrightarrow[n \rightarrow +\infty]{a.s.} 0. \quad (3.4.13)$$

*Proof.* This Lemma has already been proved in Bardet et al. (2020b) in a more general framework. Let us prove the result for  $i = 0$ . Other cases can be deduced by using a similar reasoning.

We have for any  $\theta \in \Theta$ ,  $|\hat{L}_n(\theta) - L_n(\theta)| \leq \sum_{t=1}^n |\hat{q}_t(\theta) - q_t(\theta)|$ . Then,

$$\frac{1}{\log \log n} \|\hat{L}_n(\theta) - L_n(\theta)\|_{\Theta} \leq \frac{1}{\log \log n} \sum_{t=1}^n \|\hat{q}_t(\theta) - q_t(\theta)\|_{\Theta}.$$

By Corollary 1 of Kounias and Weng (1969), (3.4.13) is established when:

$$\sum_{k \geq 1} \frac{1}{\log \log k} \mathbb{E} \|\hat{q}_k(\theta) - q_k(\theta)\|_{\Theta} < \infty. \quad (3.4.14)$$

From Bardet and Wintenberger (2009) and Bardet et al. (2020b), there exists a constant  $C$  such that

1/ If  $X \in \mathcal{H}(M_{\theta}, f_{\theta})$ , we deduce

$$\mathbb{E} [\|\hat{q}_t(\theta) - q_t(\theta)\|_{\Theta}] \leq C \left( \sum_{j \geq t} \alpha_j(f_{\theta}, \Theta) + \sum_{j \geq t} \alpha_j(M_{\theta}, \Theta) \right). \quad (3.4.15)$$

Hence,

$$\sum_{k \geq 1} \frac{1}{\log \log k} \mathbb{E} [\|\hat{q}_k(\theta) - q_k(\theta)\|_{\Theta}] \leq C \sum_{k \geq 1} \frac{1}{\log \log k} \left( \sum_{j \geq k} \alpha_j(f_{\theta}, \Theta) + \alpha_j(M_{\theta}, \Theta) \right),$$

which is finite by definition of the class  $\mathcal{H}(M_{\theta}, f_{\theta})$ , and this achieves the proof.

2/ If  $X \in \tilde{\mathcal{H}}(\tilde{H}_{\theta})$ ,

$$\mathbb{E} [\|\hat{q}_t(\theta) - q_t(\theta)\|_{\Theta}] \leq C \left( \sum_{j \geq t} \alpha_j(H_{\theta}, \Theta) \right). \quad (3.4.16)$$

This fact along with Corollary 1 of Kounias and Weng (1969) enable us to conclude the proof in this case.  $\square$

**Lemma 3.2.** *Under the assumptions of Theorem 3.3, for any model  $m \in \mathcal{M}$  with  $\theta^* \in \overset{\circ}{\Theta(m)}$ , it holds*

$$\limsup_{n \rightarrow \infty} \left( \frac{\hat{L}_n(\tilde{\theta}(m)) - \hat{L}_n(\theta^*)}{2 \log \log n} \right) = c_{min}(m) \quad a.s. \quad (3.4.17)$$

*Proof.* Applying a second order Taylor expansion of  $\tilde{L}_n$  around  $\tilde{\theta}(m)$  for  $n$  sufficiently large such that  $\tilde{\theta}(m) \in \Theta(m)$  which are between  $\theta_m^* := (\theta^*, 0, \dots, 0)^\top$  and  $\tilde{\theta}(m)$  yields (as  $\partial \hat{L}_n(\tilde{\theta}(m)) = 0$  since  $\tilde{\theta}(m)$  is a local extremum):

$$\begin{aligned} \hat{L}_n(\theta_m^*) - \hat{L}_n(\tilde{\theta}(m)) &= \hat{L}_n(\theta_m^*) - \hat{L}_n(\tilde{\theta}(m)) \\ &= \frac{1}{2}(\tilde{\theta}(m) - \theta_m^*)^\top \frac{\partial^2 \hat{L}_n(\tilde{\theta}(m))}{\partial \theta^2} (\tilde{\theta}(m) - \theta_m^*) \\ &:= I_2(m). \end{aligned}$$

But  $I_2(m)$  can be rewritten as

$$\begin{aligned} \frac{I_2(m)}{2 \log \log n} &= \frac{1}{4} (\tilde{\theta}(m) - \theta_m^*)^\top \frac{1}{\log \log n} \left[ \frac{\partial^2 \hat{L}_n(\tilde{\theta}(m))}{\partial \theta^2} - \frac{\partial^2 L_n(\tilde{\theta}(m))}{\partial \theta^2} \right] (\tilde{\theta}(m) - \theta_m^*) \\ &\quad + \frac{1}{4 \log \log n} (\tilde{\theta}(m) - \theta_m^*)^\top \frac{\partial^2 L_n(\tilde{\theta}(m))}{\partial \theta^2} (\tilde{\theta}(m) - \theta_m^*) \\ &=: I_{21}(m) + I_{22}(m). \end{aligned}$$

First, as  $\tilde{\theta}(m) \xrightarrow[n \rightarrow +\infty]{a.s.} \theta_m^*$  along side with Lemma 3.1, it follows

$$I_{21}(m) \xrightarrow[n \rightarrow +\infty]{a.s.} 0. \quad (3.4.18)$$

Writing a counterpart of (3.4.6) using the quasi-functions, we have

$$\tilde{\theta}(m) - \theta_m^* = \left( \partial^2 \hat{L}_n(\tilde{\theta}(m)) \right)^{-1} \frac{\partial \hat{L}_n(\theta_m^*)}{\partial \theta}$$

Hence  $I_{22}(m)$  can be rewritten as

$$\begin{aligned} I_{22}(m) &= \frac{1}{4 \log \log n} (\tilde{\theta}(m) - \theta_m^*)^\top \frac{\partial^2 L_n(\tilde{\theta}(m))}{\partial \theta^2} (\tilde{\theta}(m) - \theta_m^*) \\ &= \frac{1}{4 \log \log n} \left( \frac{\partial \hat{L}_n(\theta_m^*)}{\partial \theta} \right)^\top \left( \partial^2 \hat{L}_n(\tilde{\theta}(m)) \right)^{-1} \left( \partial^2 L_n(\tilde{\theta}(m)) \right) \left( \partial^2 \hat{L}_n(\tilde{\theta}(m)) \right)^{-1} \frac{\partial \hat{L}_n(\theta_m^*)}{\partial \theta} \\ &= -\frac{1}{2n \log \log n} \left( \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)^\top \times F(\theta_m^*)^{-1} \times \frac{\partial L_n(\theta_m^*)}{\partial \theta} (1 + o(1)) \quad \text{a.s.} \end{aligned} \quad (3.4.19)$$

since from (3.4.7), it holds

$$-\frac{n}{2} \left( \partial^2 L_n(\tilde{\theta}(m)) \right)^{-1} \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta_m^*)^{-1}$$

and along with Lemma 3.1, it also holds

$$-\frac{n}{2} \log \log n \left( \partial^2 \hat{L}_n(\tilde{\theta}(m)) \right)^{-1} \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta_m^*)^{-1}.$$

As a consequence, the chain of following equalities holds a.s.

$$\begin{aligned} \limsup_{n \rightarrow \infty} \left( \frac{\hat{L}_n(\tilde{\theta}(m)) - \hat{L}_n(\theta_m^*)}{2 \log \log n} \right) &= \limsup_{n \rightarrow \infty} \left( -\frac{1}{2n \log \log n} \left( \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right)^\top \times F(\theta_m^*)^{-1} \times \frac{\partial L_n(\theta_m^*)}{\partial \theta} \right) \\ &= c_{min}(m) \end{aligned}$$

That ends the proof of (3.4.17).  $\square$



# 4

## Efficient and consistent data-driven model selection for time series

---

4.1	Introduction	86
4.2	Model selection framework	87
4.2.1	Finite family $\mathcal{M}$ of parametric affine causal models	87
4.2.2	Maximum Likelihood Estimation	88
4.2.3	Quasi-Maximum Likelihood Estimation	89
4.2.4	The penalization procedure	90
4.3	Asymptotic behavior of the QMLE	91
4.3.1	Notations and main assumptions	91
4.3.2	New asymptotic results satisfied by $\hat{\theta}_m$	93
4.4	Efficient model selection result	93
4.5	From a Bayesian model selection to a data-driven consistent model selection	95
4.5.1	Bayesian model selection	95
4.6	Computations of the ideal penalties	97
4.6.1	Computations of $\text{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right)$	97
4.7	Numerical Studies	100
4.8	Proofs	101
4.8.1	Proofs of Section 4.3	101
4.8.2	Proofs of Section 4.4	107

---

The content of this chapter is based on a work in progress in collaboration with J-M BARDET and W. KENGNE.

### Abstract



This chapter studies the model selection problem in a large class of causal time series models, which includes both the ARMA or  $\text{AR}(\infty)$  processes, as well as the GARCH or ARCH( $\infty$ ), APARCH, ARMA-GARCH and many others processes. On the one hand, by seeking a penalty that minimizes the risk induced by the quasi-likelihood, we establish the asymptotic efficiency of a data-driven penalty generalizing the AIC penalty term. On the other hand, by a Bayesian approach, we also show the asymptotic consistency of a data-driven penalty that generalizes the BIC criterion. Monte Carlo experiments are performed to highlight the obtained results.

## 4.1 Introduction

Model selection is one of the fundamental tasks in scientific research specially in Statistics and Data Science. It aims at providing a model (or an algorithm) that is best, following a criterion, suitable to observed data in order to gain insight on the underlying data generating process or/and to make good forecasts.

Two leading model selection procedures have received a lot of attention in the literature. On one hand, the resampling methods such as hold out or more generally  $V$ -fold cross-validation are widely used in machine learning community. On the other hand, the methods of penalization are also now very popular in the community of applied or theoretical statisticians..

The main challenging task when designing a penalized based criterion is the calibration of the penalty. This is mainly dependent on the goal one would like the final criterion achieves. For instance, the objective could be the *consistency*, the *efficiency* or the *adaptive nature in the minimax sense* to name few.

The consistency property aims at identifying the data generating process with high probability. Hence, it requires the assumption whereby there exists a true model in the set of competitive models and the goal is to select this with probability approaches one as the sample size tends to infinity.

Although the consistency is a convincing mathematical property, this asymptotic property is not always the most interesting when switching to a practical implementation. Indeed the true underlying process is generally unknown and trying to identify the true model for any data is quite ambitious. It is often more plausible to assume that the true data generating process is infinite-dimensional, and that one tries to identify a "good" finite-dimensional model based on the data (Hurvich and Tsai (1989)). Therefore, it is common in this framework to let the dimension of the competitive models to depend on the number of observations in order to obtain better approximation and to reduce the risk of prediction. Hence, the model selection is said to be efficient when its risk is asymptotically equivalent to the risk of the *oracle*.

In this work, we are interested by providing efficient and consistent penalized data-driven criteria for affine causal times series.

**Class  $\mathcal{AC}(M, f)$  :** A process  $X = (X_t)_{t \in \mathbb{Z}}$  belongs to  $\mathcal{AC}(M, f)$  if it satisfies:

$$X_t = M((X_{t-i})_{i \in \mathbb{N}^*}) \xi_t + f((X_{t-i})_{i \in \mathbb{N}^*}) \quad \text{for any } t \in \mathbb{Z}. \quad (4.1.1)$$

where  $(\xi_t)_{t \in T}$  is a sequence of zero-mean independent identically distributed random vectors (i.i.d.r.v) satisfying  $\mathbb{E}(|\xi_0|^r) < \infty$  with  $r \geq 1$  and  $M, f : \mathbb{R}^\infty \rightarrow \mathbb{R}$  are two measurable functions, where  $\mathbb{R}^\infty$  is the set of numeric sequence with finite number of non-zero terms.

For instance,

- if  $M((X_{t-i})_{i \in \mathbb{N}^*}) = \sigma$  and  $f((X_{t-i})_{i \in \mathbb{N}^*}) = \phi_1 X_{t-1} + \dots + \phi_p X_{t-p}$ , then  $(X_t)_{t \in \mathbb{Z}}$  is an AR( $p$ ) process;
- if  $M((X_{t-i})_{i \in \mathbb{N}^*}) = \sqrt{a_0 + a_1 X_{t-1}^2 + \dots + a_p X_{t-p}^2}$  and  $f((X_{t-i})_{i \in \mathbb{N}^*}) = 0$ , then  $(X_t)_{t \in \mathbb{Z}}$  is an ARCH( $p$ ) process.

Note that, numerous classical time series models such as ARMA( $p, q$ ), GARCH( $p, q$ ), ARMA( $p, q$ )-GARCH( $p, q$ ) (see [Ding et al. \(1993\)](#) and [Ling and McAleer \(2003\)](#)) or APARCH( $\delta, p, q$ ) processes (see [Ding et al. \(1993\)](#)) belongs to  $\mathcal{AC}(M, f)$ .

The study of this type of process more often requires the classical regularity conditions on the functions  $M$  and  $f$  that are not really restrictive and remain valid for many time series.

In this semi-parametric framework, we consider  $(f_\theta)_{\theta \in \Theta}$  and  $(M_\theta)_{\theta \in \Theta}$  two families of functions such as for  $\theta \in \Theta$ ,  $f_\theta : \mathbb{R}^\infty \rightarrow \mathbb{R}$  and  $M_\theta : \mathbb{R}^\infty \rightarrow [0, \infty)$  are known.

There already exist several important contributions devoted to the model selection for time series ; we refer to the book of [McQuarrie and Tsai \(1998\)](#) and the references therein for an overview on this topic. Also, the time series model selection literature is very extensive and still growing ; we refer to the monograph of [Rao et al. \(2001\)](#), which provided an excellent summary of existing model selection procedure, including the case of time series models as well as the recent review paper of [Ding et al. \(2018\)](#).

The asymptotically efficient selection property have already been tackled: by [Shibata \(1980\)](#), and recently by [Hsu et al. \(2019b\)](#)). But this was done for linear process type AR( $\infty$ ). In this paper, we focus on the class of models (4.1.1), and addressed this questions: what regularity conditions are sufficient to build a efficient model selection criteria? Are these obtained criteria data-driven? Is the classic criterion such as AIC, still have efficient property for choosing a model among the collection  $\mathcal{M}$ ?

These questions have not yet been answered for the affine causal class and the framework considered here. This new contribution provides theoretical and numerical response of these issues. Also, following the derivation technique of the BIC criterion, we propose a data driven criterion allowing to obtain better results in consistency.

The paper is organized as follows. The model selection framework along with notations and assumptions are described in Section 4.2. Efficient criteria and the asymptotic efficiency are studied in Section 4.4. In Section 4.5, the derivation of a consistent data-driven BIC type criteria for affine causal class is studied. Section 4.7 provides some simulation experiments and finally, Section 4.8 contains the proofs.

## 4.2 Model selection framework

### 4.2.1 Finite family $\mathcal{M}$ of parametric affine causal models

Assume a trajectory  $(X_1, \dots, X_n)$  is observed from a causal stationary solution of (4.1.1) where  $M$  and  $f$  are two known functions indexed by an unknown finite dimensional vector of parameters  $\theta^*$ . Also  $(X_1, \dots, X_n)$  is supposed to be sampled according to a joint

distribution  $\mathbb{P}_{(X_1, \dots, X_n)}$ .

Now consider a finite family  $\mathcal{M}$  of models belonging to parametric affine causal models. In Proposition 1 of [Bardet et al. \(2020b\)](#), due to the linearity of such models, it was established that it is always possible to find a dimension  $d \in \mathbb{N}^*$  and a unique couple of known functions  $(M_\theta, f_\theta)$  with  $\theta \in \mathbb{R}^d$  in such a way that any model  $m \in \mathcal{M}$  belongs to the class  $\mathcal{AC}(M_\theta, f_\theta)$ . More precisely, there is a one-to-one correspondence between each model  $m \in \mathcal{M}$  and a linear subspace  $\Theta_m \subset \mathbb{R}^d$  and  $\dim(\Theta_m) = |m|$  the number of unknown parameters of the model  $m$ . As a consequence, if we denote  $m^*$  the "true" model corresponding to  $\mathcal{AC}(M_{\theta^*}, f_{\theta^*})$ , we will say:

- if  $m \in \mathcal{M}$  is such that  $\Theta_{m^*} \subset \Theta_m$  and  $\Theta_{m^*} \neq \Theta_m$  (also denoted  $m^* \subset m$  and  $m^* \neq m$ ), this is an overfitting's case;
- if  $m \in \mathcal{M}$  is such that  $\Theta_{m^*} \not\subset \Theta_m$  (also denoted  $m^* \not\subset m$ ), this is a misspecified case.

For example, if  $m^*$  corresponds to a AR(2) process and if  $\mathcal{M}$  contains AR( $p_{\max}$ ) processes and ARCH( $q_{\max}$ ), we have  $d = 1 + p_{\max} + q_{\max}$  and for  $\theta = (\theta_i)_{0 \leq i \leq d}$ ,

$$f_\theta((X_{t-k})_{k \geq 1}) = \sum_{i=1}^{p_{\max}} \theta_i X_{t-i} \quad \text{and} \quad M_\theta((X_{t-k})_{k \geq 1}) = \left( \theta_0 + \sum_{i=p_{\max}+1}^{p_{\max}+q_{\max}} \theta_i X_{t-i}^2 \right)^{1/2}.$$

Then  $\Theta_{m^*} = \{(\theta_0, \theta_1, \theta_2, 0, \dots, 0), (\theta_0, \theta_1, \theta_2) \in \mathbb{R}^3\}$ , an AR(4) process implies an overfitting, while an AR(1) or an ARCH(2) process implies a misspecified case.

In the sequel, we will always assume that

$$m^* \in \mathcal{M}.$$

But  $m^*$  is supposed to be unknown and our goal is to find a "best" model among a finite family  $\mathcal{M}$  that forecasts with a minimum risk (defined in Subsection 4.2.2) or the most probable model after observing the trajectory  $(X_1, \dots, X_n)$  (see Section 4.5). This section aims at describing the Maximum and Quasi-Maximum Likelihood contrasts that provides our chosen risk and presenting a data-driven penalization procedure for model selection.

## 4.2.2 Maximum Likelihood Estimation

Given a predictor  $\theta$ , we measure its quality by the risk defined as

$$R(\theta) := \mathbb{P}\gamma(\theta) = \mathbb{E}[\gamma(\theta, X_1)]$$

$$\text{with } \gamma(\theta, X_t) := \frac{(X_t - f_\theta^t)^2}{H_\theta^t} + \log(H_\theta^t), \text{ and } \begin{cases} f_\theta^t &:= f_\theta((X_{t-k})_{k \geq 1}) \\ M_\theta^t &:= M_\theta((X_{t-k})_{k \geq 1}) \\ H_\theta^t &:= (M_\theta^t)^2 \end{cases} \quad (4.2.1)$$

By referring to [Massart \(2007\)](#) or [Francq and Zakoian \(2010\)](#), the contrast  $\gamma(\theta, \cdot)$  is  $-2$  times the Gaussian conditional log-density of  $X_t$ . Moreover, the Gaussian MLE is derived from the conditional (with respect to the filtration  $\sigma\{(X_t)_{t \leq 0}\}$ ) log-likelihood of  $(X_1, \dots, X_n)$  when  $(\xi_t)$  is supposed to be a Gaussian standard white noise. We deduce that this conditional log-likelihood (up to an additional constant)  $L_n$  is defined for a parameter  $\theta$  by:

$$L_n(\theta) := -\frac{1}{2} \sum_{t=1}^n \gamma(\theta, X_t) \quad (4.2.2)$$

As proved in [Bardet and Wintenberger \(2009\)](#), the risk function  $R$  achieves its unique minimum at the "true" parameter  $\theta^*$  over any parameter set  $\Theta$ , when  $\theta^* \in \Theta$

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} R(\theta). \quad (4.2.3)$$

Therefore  $\theta^*$  is considered as an ideal for the model selection procedure and serves as a benchmark to compare predictors. Given a model  $m \in \mathcal{M}$  with  $\Theta_m$  its parameter space that does not necessary contains  $\theta^*$ , let define

$$\theta_m^* = \operatorname{argmin}_{\theta \in \Theta_m} R(\theta). \quad (4.2.4)$$

As a consequence, we have:

$$\theta_{m^*}^* = \theta^* \quad \text{and more generally, if } m^* \subset m, \quad \theta_m^* = \theta^*.$$

Besides of minimizing the risk, we also consider the minimization of the natural associated loss function, which is defined as

$$\ell(\theta, \theta^*) := R(\theta) - R(\theta^*) \geq 0. \quad (4.2.5)$$

This is a well-known measure of separation between the candidate model generated by  $\theta$  and the true one indexed by  $\theta^*$ .

Let us introduce the empirical distribution  $\mathbb{P}_n$  be such that

$$\mathbb{P}_n = \frac{1}{n} \sum_{t=1}^n \delta_{X_t},$$

where  $\delta_{X_t}$  is the Dirac distribution at observation  $X_t$ . Let set by  $\gamma_n$  the associated empirical criterion

$$\gamma_n(\theta) := \mathbb{P}_n \gamma(\theta, \cdot) = \frac{1}{n} \sum_{t=1}^n \gamma(\theta, X_t).$$

We have  $\gamma_n(\theta) = -\frac{2}{n} L_n(\theta)$ , so that maximizing the log-likelihood is equivalent to minimize the empirical criterion  $\gamma_n$ .

### 4.2.3 Quasi-Maximum Likelihood Estimation

Since the the white noise is not necessary a Gaussian one and since the log-likelihood (and then the empirical risk)  $L_n(\theta)$  depends on  $(X_t)_{t \leq 0}$  that are unknown, a quasi-log-likelihood can be used as an approximation of the log-likelihood. It consists of replacing  $\gamma(\theta, X_t)$  by an approximation  $\hat{\gamma}(\theta, X_t)$ . Hence, in our framework, we consider the following conditional quasi log-likelihood (up to an additional constant) given for all  $\theta \in \Theta$  by

$$\begin{aligned} \hat{L}_n(\theta) &:= -\frac{1}{2} \sum_{t=1}^n \hat{\gamma}(\theta, X_t) \\ \text{with } \hat{\gamma}(\theta, X_t) &:= \frac{(X_t - \hat{f}_\theta^t)^2}{\hat{H}_\theta^t} + \log(\hat{H}_\theta^t) \quad \text{and} \quad \begin{cases} \hat{f}_\theta^t &:= f_\theta(X_{t-1}, X_{t-2}, \dots, X_1, u) \\ \hat{M}_\theta^t &:= M_\theta(X_{t-1}, X_{t-2}, \dots, X_1, u) \\ \hat{H}_\theta^t &:= (\hat{M}_\theta^t)^2 \end{cases} \end{aligned} \quad (4.2.6)$$

for any deterministic sequence  $u = (u_n)_{n \in \mathbb{N}}$  with finitely many non-zero values (we will use  $u = 0$  without loss of generality).

In addition the computable empirical risk is then:

$$\hat{\gamma}_n(\theta) = \mathbb{P}_n \hat{\gamma}(\theta, \cdot) = -\frac{2}{n} \hat{L}_n(\theta).$$

Finally, for each specific model  $m \in \mathcal{M}_n$ , we define the Gaussian Quasi-Maximum Likelihood Estimator (QMLE)  $\hat{\theta}_m$  as

$$\hat{\theta}_m = \underset{\theta \in \Theta_m}{\operatorname{argmin}} \hat{\gamma}_n(\theta). \quad (4.2.7)$$

The estimator  $\hat{\theta}_m$  is commonly called the *Empirical Risk Minimizer* (ERM).

#### 4.2.4 The penalization procedure

For  $m \in \mathcal{M}$ , the ERM provides an estimator in  $\Theta_m$ . The goal is to come up with the model that minimizes the excess loss over  $\mathcal{M}$

$$\inf_{m \in \mathcal{M}} \ell(\hat{\theta}_m, \theta^*). \quad (4.2.8)$$

This model is unknown since (4.2.8) depends on  $\theta^*$  and the distribution  $P_{(X_1, \dots, X_n)}$  that are unknown.

A classical way to solve (4.2.8) problem is to design for every  $m \in \mathcal{M}$  an estimator of  $R(\hat{\theta}_m)$  and we naturally choose  $\hat{\gamma}_n(\hat{\theta}_m)$ . First, it is well known that the empirical criterion  $\hat{\gamma}_n(\hat{\theta}_m)$  is an optimist version of  $R(\hat{\theta}_m)$  and decreases with the dimension of the model. Therefore, it is common to add a penalty term to counteract this bias. As a consequence, define a function pen, which is called the penalty function, possibly data-dependent, such as pen:  $m \in \mathcal{M} \mapsto \text{pen}(m) \in \mathbb{R}^+$ . Then define the penalized contrast and the model selected by it:

$$\hat{m}_{\text{pen}} = \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \hat{C}_{\text{pen}}(m) \} \quad \text{with} \quad \hat{C}_{\text{pen}}(m) := \hat{\gamma}_n(\hat{\theta}_m) + \text{pen}(m). \quad (4.2.9)$$

In order to achieve (4.2.8), the *ideal penalty* to consider in (4.2.9) is

$$\text{pen}_{id}(m) = R(\hat{\theta}_m) - \hat{\gamma}_n(\hat{\theta}_m). \quad (4.2.10)$$

Using its definition we have:

$$\hat{m}_{id} := \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \ell(\hat{\theta}_m, \theta^*) \} = \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ R(\hat{\theta}_m) \} = \underset{m \in \mathcal{M}}{\operatorname{argmin}} \{ \hat{C}_{\text{pen}_{id}}(m) \}. \quad (4.2.11)$$

However, the function  $R$  is unknown (as it depends on the unknown distribution of the process) therefore  $\text{pen}_{id}$  cannot generally be used directly.

The question is how to choose the penalty in (4.2.9) so that  $\hat{m}_{\text{pen}}$  mimics the *oracle*, which is the model associated with the minimum risk *i.e.* (4.2.8). Hence, we would like our final estimator  $\hat{\theta}_{\hat{m}_{\text{pen}}}$  to behave asymptotically like the oracle. That is to satisfy:

$$\mathbb{P} \left( \ell(\hat{\theta}_{\hat{m}_{\text{pen}}}, \theta^*) \leq \min_{m \in \mathcal{M}} \{ \ell(\hat{\theta}_m, \theta^*) \} + \frac{C}{n} \right) \xrightarrow[n \rightarrow \infty]{} 1 \quad (4.2.12)$$

and/or for any  $n \geq n_0$

$$\mathbb{E}[\ell(\hat{\theta}_{\hat{m}_{\text{pen}}}, \theta^*)] \leq \min_{m \in \mathcal{M}} \left\{ \mathbb{E}[\ell(\hat{\theta}_m, \theta^*)] \right\} + \frac{C}{n}. \quad (4.2.13)$$

The aim of this paper is to find a good choice of  $\text{pen}(m)$  in order to obtain the *asymptotic optimality* (4.2.13), that we prefer to (4.2.12), which means that the expectation of the risk of the chosen estimator is equivalent to the expectation of the risk of the oracle when the sample size tends to infinity.

### 4.3 Asymptotic behavior of the QMLE

Before considering the problem of model selection, we establish a central limit theorem satisfied by  $\hat{\theta}_m$  for any model  $m \in \mathcal{M}$ , *i.e.* as well if  $m$  is an overfitted or a misspecified model. Before this, some notations and assumptions have to be precised.

#### 4.3.1 Notations and main assumptions

In the sequel, we will consider a subset  $\Theta$  of  $\mathbb{R}^d$  which is compact. We will use the following norms:

- $\|\cdot\|$  denotes the usual Euclidean norm on  $\mathbb{R}^\nu$ , with  $\nu \geq 1$ ;
- for a matrix  $A$ , denote  $\|A\|$  the subordinate matrix norm such that  $\|A\| = \sup_{v \neq 0} \frac{\|Av\|}{\|v\|}$ ;
- if  $X$  is a  $\mathbb{R}^\nu$ -random variable and  $r \geq 1$ , we set  $\|X\|_r = (\mathbb{E}[\|X\|^r])^{1/r} \in [0, \infty]$ ;
- for  $\theta \in \Theta \subset \mathbb{R}^d$ , if  $\Psi_\theta : \mathbb{R}^\infty \rightarrow E$  where  $E = \mathbb{R}^\nu$  or  $E$  is a set of square matrix, denote  $\|\Psi_\theta(\cdot)\|_\Theta = \sup_{\theta \in \Theta} \{\|\Psi_\theta(\cdot)\|\}$ ;
- for  $\theta \in \Theta \subset \mathbb{R}^d$ , if  $\Psi_\theta : \mathbb{R}^\infty \rightarrow \mathbb{R}$  is a  $\mathcal{C}^2(\Theta \times \mathbb{R}^\infty)$  function, we will denote

$$\partial_\theta \Psi_\theta(\cdot) = \left( \frac{\partial}{\partial \theta_i} \Psi_\theta(\cdot) \right)_{1 \leq i \leq d} = (\partial_{\theta_i} \Psi_\theta(\cdot))_{1 \leq i \leq d} \quad \text{and} \quad \partial_{\theta^2}^2 \Psi_\theta(\cdot) = \left( \frac{\partial^2}{\partial \theta_i \partial \theta_j} \Psi_\theta(\cdot) \right)_{1 \leq i, j \leq d};$$

- consider  $\Psi_\theta : \mathbb{R}^\infty \rightarrow \mathbb{R}$  for any  $\theta \in \Theta \subset \mathbb{R}^d$ . Then, we define:

**A**( $\Psi_\theta, \Theta$ ):  $\|\Psi_\theta(0)\|_\Theta < \infty$  and there exists a sequence of non-negative real numbers  $(\alpha_k(\Psi_\theta, \Theta))_{k \geq 1}$  such that  $\sum_{k=1}^{\infty} \alpha_k(\Psi_\theta, \Theta) < \infty$  satisfying:

$$\|\Psi_\theta(x) - \Psi_\theta(y)\|_\Theta \leq \sum_{k=1}^{\infty} \alpha_k(\Psi_\theta, \Theta) |x_k - y_k| \quad \text{for all } x, y \in \mathbb{R}^\infty.$$

Several assumptions on the AC class will be considered thereafter:

**Assumption A0:** The process  $X \in \mathcal{AC}(M_{\theta^*}, f_{\theta^*})$  where  $\theta^* \in \Theta$  is defined in (4.1.1) where:

- the white noise  $(\xi_t)_t$  is such as  $\|\xi_0\|_r < \infty$  with  $8 < r$ ;
- for any  $x \in \mathbb{R}^\infty$ , the functions  $\theta \rightarrow M_\theta$  and  $\theta \rightarrow f_\theta$  are  $\mathcal{C}^2(\Theta)$  functions:

- $\Theta \in \mathbb{R}^d$  is a compact set such as

$$\Theta \subset \left\{ \theta \in \mathbb{R}^d, A(f_\theta, \{\theta\}) \text{ and } A(M_\theta, \{\theta\}) \text{ hold with} \right. \\ \left. \sum_{k=1}^{\infty} \alpha_k(f_\theta, \{\theta\}) + \|\xi_0\|_r \sum_{k=1}^{\infty} \alpha_k(M_\theta, \{\theta\}) < 1 \right\}. \quad (4.3.1)$$

Under this assumption, Doukhan and Wintenberger (2008) showed that there exists a stationary causal (*i.e.*  $X_t$  is depending only on  $(X_{t-k})_{k \in \mathbb{N}}$  for any  $t \in \mathbb{Z}$ ) and ergodic solution of (4.1.1) with 8-order moment for any  $\theta \in \Theta$ .

Now assumption **A0** holds. We will also add several assumptions required for insuring the strong consistency and the asymptotic normality of the QMLE:

The following classical assumption ensures the identifiability of the  $\theta^*$ .

**Assumption A1:** For all  $\theta, \theta' \in \Theta$ ,  $(f_\theta^0 = f_{\theta'}^0 \text{ and } M_\theta^0 = M_{\theta'}^0) \text{ a.s.} \implies \theta = \theta'$ .

Next, the following Assumption ensures the invertibility of the "Fisher's information matrix" and is necessary to prove the asymptotic normality of the QMLE.

**Assumption A2:**  $\langle \alpha, \partial_\theta f_\theta^0 \rangle = 0 \implies \alpha = 0 \text{ a.s. or } \langle \alpha, \partial_\theta H_\theta^0 \rangle = 0 \implies \alpha = 0 \text{ a.s.}$

The definition of the computable empirical risk and require that its denominators do not vanish. Hence, we are going to assume throughout this paper that the lower bound of  $H_\theta(\cdot) = (M_\theta(\cdot))^2$  is strictly positive:

**Assumption A3:**  $\exists \underline{h} > 0$  such that  $\inf_{\theta \in \Theta} (H_\theta(x)) \geq \underline{h}$  for all  $x \in \mathbb{R}^\infty$ .

The following assumption is a technical classical condition (see Lv and Liu (2014)).

**Assumption A4:** For every  $m \in \mathcal{M}$ , if  $(\bar{\theta}_{m,n})$  is a sequence of  $\Theta_m$  satisfying  $\bar{\theta}_{m,n} \xrightarrow[n \rightarrow +\infty]{a.s.} \theta^*$ , then

$$\limsup_{n \rightarrow \infty} \left\{ \mathbb{E} \left[ \left\| \frac{1}{n} (\partial_{\bar{\theta}_i \theta_j}^2 L_n(\bar{\theta}_m))_{i,j \in m} \right\|^{-1} \right]^8 \right\} < \infty. \quad (4.3.2)$$

**Remark 4.1.** Note that under assumption **A0**, if  $\bar{\theta}_{m,n} \xrightarrow[n \rightarrow +\infty]{a.s.} \theta_m^*$  then

$$\left\| \left( \frac{1}{n} (\partial_{\bar{\theta}_i \theta_j}^2 L_n(\bar{\theta}_{m,n}))_{i,j \in m} \right)^{-1} \right\|^8 \xrightarrow[n \rightarrow +\infty]{a.s.} \left\| \left( -\frac{1}{2} \partial_{\bar{\theta}_i \theta_j}^2 \gamma(\theta_m^*)_{i,j \in m} \right)^{-1} \right\|^8$$

Thus, from the Egorov's Theorem, we can find an event  $\tilde{\Omega}$  with sufficiently large probability such that the relation (4.3.2) in the assumption **A4** holds if the expectation is taken on the event  $\tilde{\Omega}$ . For the particular case of the linear processes, the assumption **A4** holds true under a mild condition on the distribution of  $X$ , see for instance Papangelou (1994) and Findley and Wei (2002).

Finally, the decrease rates of  $(\alpha_j(f_\theta, \Theta))_j$ ,  $(\alpha_j(M_\theta, \Theta))_j$ ,  $(\alpha_j(\partial_\theta f_\theta, \Theta))_j$  and  $(\alpha_j(\partial_\theta M_\theta, \Theta))_j$  have to be fast enough for insuring the strong consistency and the asymptotic normality of the QMLE:

**Assumption A5:** Conditions **A**( $f_\theta, \Theta$ ), **A**( $M_\theta, \Theta$ ), **A**( $\partial_\theta f_\theta, \Theta$ ), **A**( $\partial_\theta M_\theta, \Theta$ ), **A**( $\partial_{\theta^2}^2 f_\theta, \Theta$ ) and **A**( $\partial_{\theta^2}^2 M_\theta, \Theta$ ) hold with

$$\alpha_j(f_\theta, \Theta) + \alpha_j(M_\theta, \Theta) + \alpha_j(\partial_\theta f_\theta, \Theta) + \alpha_j(\partial_\theta M_\theta, \Theta) = O(j^{-\delta}) \quad \text{where } \delta > 7/2.$$

Note that Assumption A5 does not allow to consider long-range dependent processes, but usual short memory causal time series satisfy this assumption.

### 4.3.2 New asymptotic results satisfied by $\hat{\theta}_m$

The asymptotic normality of  $\hat{\theta}_m$  has been already established in [Bardet and Wintenberger \(2009\)](#) when  $m = m^*$  and in [Bardet et al. \(2020b\)](#) when  $m^* \subset m$  (overfitting). This property can also be extended in the case of misspecified model, *i.e.* when  $m^* \not\subset m$ . Before this, we define the following key matrix for  $\theta_m \in \Theta_m$  with  $m \in \mathcal{M}$ ,

$$F(\theta_m) := \left( \mathbb{E} \left[ \partial_{\theta_i}^2 \gamma(\theta_m, X_0) \right] \right)_{i,j \in m}. \quad (4.3.3)$$

First, the following corollary can be established as a particular case of more general result, Proposition 4.3, which is stated in Section 4.8 devoted to the proofs.

**Corollary 4.1.** *For  $m \in \mathcal{M}$ , let  $\theta_m \in \Theta_m$  and denote*

$$\begin{aligned} G(\theta_m) &:= \frac{1}{4} \left( \sum_{t \in \mathbb{Z}} \text{Cov}(\partial_{\theta_i} \gamma(\theta_m, X_0), \partial_{\theta_j} \gamma(\theta_m, X_t)) \right)_{i,j \in m} \\ \implies G(\theta_m^*) &= \frac{1}{4} \left( \text{Cov}(\partial_{\theta_i} \gamma(\theta_m^*, X_0), \partial_{\theta_j} \gamma(\theta_m^*, X_0)) \right)_{i,j \in m} \quad \text{if } m^* \subset m. \end{aligned} \quad (4.3.4)$$

Then, under assumptions **A0-A5**, with  $\theta_m^*$  defined in (4.2.4),

$$\frac{1}{\sqrt{n}} \left( \partial_{\theta_j} L_n(\theta_m^*) \right)_{j \in m} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, G(\theta_m^*)). \quad (4.3.5)$$

Using mainly this new result, we also obtain:

**Theorem 4.1.** *Under assumptions **A0-A5**, for any  $m \in \mathcal{M}$ ,*

$$\sqrt{n} \left( (\hat{\theta}_m)_i - (\theta_m^*)_i \right)_{i \in m} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N} \left( 0, (F(\theta_m^*))^{-1} G(\theta_m^*) (F(\theta_m^*))^{-1} \right), \quad (4.3.6)$$

with  $F$  defined in (4.3.3) and  $G$  in (4.3.4).

Hence, even in the misspecified case,  $\hat{\theta}_m$  satisfies a central limit theorem. We will use this result several times, in particular to prove that the probability of selecting a misspecified model tends quickly enough towards 0. Another technical result will also be useful for the sequel:

**Proposition 4.1.** *Under assumptions **A0-A5**, with  $8/3 < r' \leq (8+r)/6$  and  $r' < 2(\delta-1)$  where  $\delta > 7/2$  is given in Assumption A5 and for any  $m \in \mathcal{M}$ , then we have*

$$\sup_{n \in \mathbb{N}^*} \left\| \sqrt{n} \left( (\hat{\theta}_m)_i - (\theta_m^*)_i \right)_{i \in m} \right\|_{r'} < \infty. \quad (4.3.7)$$

Note that we also have  $\sup_{n \in \mathbb{N}^*} \left\| \sqrt{n} \left( (\hat{\theta}_m)_i - (\theta_m^*)_i \right)_{i \in m} \right\|_2 < \infty$ . This result will be essential for establishing the asymptotic behavior of the expectation of the ideal penalty.

## 4.4 Efficient model selection result

The expectation of the ideal penalty (4.2.10) has been computed (or asymptotically approximated) in several frameworks (see [Mallows \(1973\)](#), [Akaike \(1973\)](#), [Schwarz \(1978\)](#), [Hurvich and Tsai \(1989\)](#), [Cavanaugh \(1997\)](#); etc) and it is most often proportional to the dimension of the model (denoted  $D_m$  in the sequel).



- [Mallows \(1973\)](#) in the regression setting and the penalty is  $2D_m \sigma^2/n$ ;
- [Akaike \(1973\)](#) in the density estimation framework where the penalty is  $D_m/n$ ;
- [Lv and Liu \(2014\)](#) in misspecified density estimation and the penalty is  $\text{Trace}(B_n A_n^{-1})/n$  where  $A_n$  is the opposite of the Hessian matrix of the log-likelihood and  $B_n$  the Fisher Information matrix.

In order to approximate (4.2.10) in this framework, let first provide a decomposition of this term in order to facilitate the computation. For any model  $m \in \mathcal{M}$ , write

$$\text{pen}_{id}(m) := R(\hat{\theta}_m) - \hat{\gamma}_n(\hat{\theta}_m) = I_1(m) + I_2(m) + I_3(m), \quad (4.4.1)$$

$$\text{with } \begin{cases} I_1(m) &:= R(\hat{\theta}_m) - R(\theta_m^*) \\ I_2(m) &:= \hat{\gamma}_n(\theta_m^*) - \hat{\gamma}_n(\hat{\theta}_m) \\ I_3(m) &:= R(\theta_m^*) - \hat{\gamma}_n(\theta_m^*) \end{cases}.$$

Next we provide a preliminary result about the asymptotic behavior of the terms  $I_1(m)$  and  $I_2(m)$ . Then we obtain:

**Lemma 4.1.** *Under assumptions A0-A5, for any model  $m \in \mathcal{M}$ , there exists a probability distribution  $U^*(m)$  such that*

$$\begin{aligned} 1. \quad n I_1(m) &= n (R(\hat{\theta}_m) - R(\theta_m^*)) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} U^*(m) \\ \text{and } \mathbb{E}[n I_1(m)] &\xrightarrow[n \rightarrow \infty]{} \mathbb{E}[U^*(m)] = \frac{1}{2} \text{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right). \end{aligned} \quad (4.4.2)$$

$$\begin{aligned} 2. \quad n I_2(m) &= n (\hat{\gamma}_n(\theta_m^*) - \hat{\gamma}_n(\hat{\theta}_m)) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} U^*(m) \\ \text{and } \mathbb{E}[n I_2(m)] &\xrightarrow[n \rightarrow \infty]{} \frac{1}{2} \text{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right). \end{aligned} \quad (4.4.3)$$

The proof of this lemma, as well as all the other proofs, can be found in Section 4.8. This result leads to our first main result devoted the ideal penalty defined in (4.2.10).

**Proposition 4.2.** *Under assumptions A0-A5 and for any  $m \in \mathcal{M}$ , there exists a bounded sequence  $(v_n^*)_{n \in \mathbb{N}^*}$  not depending on  $m$  satisfying*

$$\mathbb{E}[\text{pen}_{id}(m)] \underset{n \rightarrow \infty}{\sim} 2 \left( \frac{1}{2n} \text{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right) \right) + \frac{v_n^*}{n}. \quad (4.4.4)$$

Note that the Slope Heuristic Procedure which allows to estimate a so called *minimal penalty* (see [Arlot and Massart \(2009\)](#)) consists in evaluating the slope of a linear regression of  $\hat{\gamma}_n(\hat{\theta}_m)$  onto  $D_m$  for  $m^* \subset m$  and this is equivalent to estimating the slope of  $-\frac{1}{2n} \text{Trace}(G(\theta_m^*)F(\theta_m^*)^{-1})$  onto  $D_m$ . We will see that  $\text{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right)$  behaves as a linear function of  $D_m$  in many cases, which also gives legitimacy to this approach in the case of time series after having obtained it in the case of linear regression. The minimal penalty is then  $-2 \times$  the estimated slope and this finally corresponds to an estimation of  $\mathbb{E}[\text{pen}_{id}(m)]$ .

**Remark 4.2.** *The trace of the matrix mentioned above is easily computable in some cases using the explicit forms of the matrices  $F(\theta^*)$ ,  $G(\theta^*)$  in [Bardet and Wintenberger \(2009\)](#). As showed in Section 4.6, this trace is proportionnal to the dimension of the model  $D_m$  in some cases but could be more complex functions of  $D_m$ .*

We can now state the main results of this paper.

**Theorem 4.2.** *Under assumptions A0-A5 and if for any  $\varepsilon > 0$ ,*

$$n \mathbb{P}(\text{pen}(m) \geq \varepsilon) \xrightarrow{n \rightarrow \infty} 0 \quad \text{for any } m \in \mathcal{M}. \quad (4.4.5)$$

*Then,*

$$n \mathbb{P}(m^* \notin \hat{m}_{\text{pen}}) \xrightarrow{n \rightarrow \infty} 0. \quad (4.4.6)$$

Theorem 4.2 says that if the penalty asymptotically decreases to 0 in probability, then the criterion  $\hat{C}_{\text{pen}}$  does not select a misspecified model asymptotically.

Now, we can specify the convergence rate of pen to obtain an excess loss close to the minimal one over  $\mathcal{M}$ :

**Theorem 4.3.** *Under assumptions A0-A5, and if for any  $\varepsilon > 0$  there exists  $K_\varepsilon > 0$  such as*

$$\limsup_{n \rightarrow \infty} \max_{m \in \mathcal{M}} \mathbb{P}(n \text{pen}(m) \geq K_\varepsilon) \leq \varepsilon. \quad (4.4.7)$$

*Then for any  $\varepsilon > 0$ , there exists  $M_\varepsilon > 0$  and  $N_\varepsilon \in \mathbb{N}^*$  such as for any  $n \geq N_\varepsilon$ ,*

$$\mathbb{P}\left(\ell(\hat{\theta}_{\hat{m}_{\text{pen}}}, \theta^*) \leq \inf_{m \in \mathcal{M}} \{\ell(\hat{\theta}_m, \theta^*)\} + \frac{M_\varepsilon}{n}\right) \geq 1 - \varepsilon. \quad (4.4.8)$$

**Remark 4.3.** *Let notice that this asymptotic optimality is quite a bit different from the classical one about asymptotic efficiency, where both the cardinal of the collection  $\mathcal{M}$  and the dimension of competitive models are allowed to depend on  $n$ . However, this is done in the framework where the parameter  $\theta^*$  is infinite-dimensional (see for example Shibata (1980), Li (1987), Hsu et al. (2019b)).*

## 4.5 From a Bayesian model selection to a data-driven consistent model selection

### 4.5.1 Bayesian model selection

Another classical paradigm for model selection is the Bayesian one, leading typically to the BIC criterion (see Schwarz (1978)). In this approach, the construction of the model selection criterion is first done by assuming that the parameter vector  $\theta^*$  is a random vector. Let recall the hierarchical prior sampling scheme in the Bayesian setting: given the finite family of models  $\mathcal{M}$ , a model  $m$  is drawn according to a prior distribution  $(\pi_m)_{m \in \mathcal{M}}$  (generally a uniform distribution) and then, conditionally on  $m$ ,  $\theta$  is sampled according to some prior distribution  $\mu_m(\theta)$ .

The goal of this model selection procedure is to choose the most probable model after observing the trajectory  $X := (X_1, \dots, X_n)$ , i.e.

$$\hat{m}_B = \operatorname{argmax}_{m \in \mathcal{M}} \{\mathbb{P}(m | X)\}. \quad (4.5.1)$$

Using Bayes Formula, we can write  $\mathbb{P}(m | X) = \frac{\pi_m \mathbb{P}(X | m)}{\mathbb{P}(X)}$ . Moreover, we have:

$$\mathbb{P}(X | m) = \int_{\Theta_m} \mathbb{P}(X | \theta, m) d\mu_m(\theta).$$

In addition, since  $\mathbb{P}(X)$  does not depend on  $m$ , and  $\mathbb{P}(X | \theta, m)$  is the likelihood of  $X$  given  $\theta \in \Theta_m$  and  $m \in \mathcal{M}$ , maximizing  $\mathbb{P}(m | X)$  is equivalent to maximizing

$$S_n(m, X) := \log(\pi_m \mathbb{P}(X | m)) = \log\left(\int_{\Theta_m} \pi_m \exp(L_n(\theta)) d\mu_m(\theta)\right).$$

From now on, we will assume that  $\pi_m = 1/|\mathcal{M}|$  for any  $m \in \mathcal{M}$ , a priori uniform distribution of the models in the family  $\mathcal{M}$ . We can also assume that there exists a non-negative Borel function  $\theta \rightarrow b_m(\theta)$  such as  $d\mu_m(\theta) = b_m(\theta) d\theta$ . Then we have:

$$S_n(m, X) = -\log(|\mathcal{M}|) + \log\left(\int_{\Theta_m} b_m(\theta) \exp(L_n(\theta)) d\theta\right) \quad (4.5.2)$$

and by replacing  $L_n$  by the quasi version  $\hat{L}_n$ , we introduce

$$\hat{S}_n(m, X) = -\log(|\mathcal{M}|) + \log\left(\int_{\Theta_m} b_m(\theta) \exp(\hat{L}_n(\theta)) d\theta\right). \quad (4.5.3)$$

Let us give an asymptotic expansion of the *a posteriori* probability  $\hat{S}_n(m, X)$  in order to derive a BIC type criterion that is coherent with our framework where the observed trajectory is that of a causal affine process. This could be obtained from a Laplace approximation leading to the following theorem:

**Theorem 4.4.** *Under assumptions A0, A1, A2, A3, A5 and if for any  $x \in \mathbb{R}^\infty$ , the functions  $\theta \rightarrow M_\theta$  and  $\theta \rightarrow f_\theta$  are  $\mathcal{C}^6(\Theta)$  functions satisfying  $\mathbf{A}(\partial_{\theta^k}^k f_\theta, \Theta)$  and  $\mathbf{A}(\partial_{\theta^k}^k M_\theta, \Theta)$  for any  $0 \leq k \leq 2$ . Then,*

$$\begin{aligned} \hat{S}_n(m, X) &= \hat{L}_n(\hat{\theta}_m) - \frac{\log(n)}{2} |m| + \log(b_m(\hat{\theta}_m)) \\ &\quad + \frac{\log(2\pi)}{2} |m| - \frac{1}{2} \log\left[\det(-\hat{F}_n(m))\right] - \log(|\mathcal{M}|) + O(n^{-1}) \quad a.s. \end{aligned} \quad (4.5.4)$$

where  $\hat{F}_n(m) := \left(\partial_{\theta_i \theta_j}^2 \hat{L}_n(\hat{\theta}_m)\right)_{i,j \in m}$

In the above equation, it is clear that  $-2\hat{S}_n(m, X) \simeq -2\hat{L}_n(\hat{\theta}_m) + \log(n)|m|$  a.s. This gives legitimacy to the usual BIC criterion within the framework of causal affine processes since:

$$\hat{m}_{BIC} = \operatorname{argmin}_{m \in \mathcal{M}} \left\{ -2\hat{L}_n(\hat{\theta}_m) + \log(n)|m| \right\},$$

and we see that  $\hat{m}_{BIC}$  maximizes the main terms of  $\hat{S}_n(m, X)$ .

From the relation (4.5.4), considering certain second order terms of the asymptotic expansion of  $\hat{S}_n(m, X)$ , we also obtain the Kashyap criterion (see Kashyap (1982), Sclove (1987), Bozdogan (1987)), defined for all  $m \in \mathcal{M}$  by

$$\begin{aligned} \widehat{KC}(m) &:= -2\hat{L}_n(\hat{\theta}_m) + \log(n)|m| + \log\left[\det(-\hat{F}_n(m))\right] \\ \text{and } \hat{m}_{KC} &= \operatorname{argmin}_{m \in \mathcal{M}} \{\widehat{KC}(m)\}. \end{aligned} \quad (4.5.5)$$

Therefore the term  $\log\left[\det(\hat{F}_n(m))\right]$  is added to the usual BIC criterion. Several example of computations of this term, generally equal to  $c|m|$  but not always, are provided in the forthcoming Section 4.6. It is clear that  $\hat{m}_{KC}$  can be more interesting than  $\hat{m}_{BIC}$  in terms of consistency only for non asymptotic framework (typically for  $n$  of the order of a hundred

or several hundred). Note also that the data-driven criteria  $\widehat{KC}$  that is "optimal" in the sense of the a posteriori probability (see Kashyap [Kashyap \(1982\)](#)) is also asymptotically consistent under the assumption **A5**, see Bardet *et al.* [Bardet et al. \(2020b\)](#).

However this choice of second order terms of the asymptotic expansion of  $\widehat{S}_n(m, X)$  is somewhere arbitrary. A criterion taking account of all the second order terms could also be defined. For this, we could define a uniform distribution  $b_m$  on a compact set included in  $\Theta_m$ . As a consequence, using condition (4.3.1) of Assumption **A0**, there always exists  $0 < C_1 \leq C_2$  such as  $\frac{C_1}{m} \leq b_m(\theta_m) \leq \frac{C_2}{m}$ . As a consequence, we could define a new data-driven consistent criterion, called  $\widehat{KC}'$ , such as for any  $m \in \mathcal{M}$

$$\widehat{KC}'(m) := -2\widehat{L}_n(\widehat{\theta}_m) + (\log(n) - \log(2\pi))|m| + \log \left[ \det(-\widehat{F}_n(m)) \right] + 2\log(|m|)$$

and  $\widehat{m}_{KC'} = \operatorname{argmin}_{m \in \mathcal{M}} \{\widehat{KC}'(m)\}$ . (4.5.6)

**Remark 4.4.** We also know that under assumptions **A0**, **A1**, **A2**, **A3**, **A5**,  $\widehat{F}_n(m) \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta_m^*)$ . Therefore  $\log(\det(F(\theta_m^*)))$  can also replace  $\log(\det(\widehat{F}_n(m)))$  in the expression of  $\widehat{KC}'(m)$ .

**Corollary 4.2.** Under assumptions **A0**, **A1**, **A2**, **A3** and **A5**, from [Bardet et al. \(2020b\)](#),  $\widehat{m}_{BIC}$ ,  $\widehat{m}_{KC}$  and  $\widehat{m}_{KC'}$  are consistent criteria, i.e.

$$\widehat{m}_{BIC} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} m^*, \quad \widehat{m}_{KC} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} m^* \quad \text{and} \quad \widehat{m}_{KC'} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} m^*.$$

## 4.6 Computations of the ideal penalties

### 4.6.1 Computations of $\operatorname{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right)$

From [Bardet and Wintenberger \(2009\)](#), with  $\mu_4 = \mathbb{E}[\xi_0^4]$ ,  $f_\theta^0$  and  $H_\theta^0$  defined in (1.2.3), we have for  $m^* \subset m$  and  $i, j \in m$ :

$$\begin{aligned} (G(\theta_m^*))_{i,j} &= \mathbb{E} \left[ \frac{\partial_{\theta_i} f_{\theta_m^*}^0 \partial_{\theta_j} f_{\theta_m^*}^0}{H_{\theta_m^*}^0} + \frac{(\mu_4 - 1)}{4} \frac{\partial_{\theta_i} H_{\theta_m^*}^0 \partial_{\theta_j} H_{\theta_m^*}^0}{(H_{\theta_m^*}^0)^2} \right] \\ (F(\theta_m^*))_{i,j} &= -\mathbb{E} \left[ \frac{\partial_{\theta_i} f_{\theta_m^*}^0 \partial_{\theta_j} f_{\theta_m^*}^0}{H_{\theta_m^*}^0} + \frac{1}{2} \frac{\partial_{\theta_i} H_{\theta_m^*}^0 \partial_{\theta_j} H_{\theta_m^*}^0}{(H_{\theta_m^*}^0)^2} \right], \end{aligned} \quad (4.6.1)$$

Here there are 3 frameworks where  $\operatorname{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right)$  can be computed for  $m^* \subset m$ :

1/ A first and well-known case is the Gaussian case. Indeed, when  $(\xi_t)$  is a Gaussian white noise, then  $\mu_4 = 3$  and then from (4.6.1), for any  $i, j \in m$ ,

$$(G(\theta_m^*))_{i,j} = -(F(\theta_m^*))_{i,j} \implies -2\operatorname{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right) = 2\operatorname{Trace}(I_{|m|}) = 2|m|,$$

with  $I_\ell$  the identity matrix of size  $\ell \in \mathbb{N}^*$ . As a consequence, in the Gaussian framework, for  $m^* \subset m$ , the expectation of the ideal penalty is exactly the classical Akaike Criterion (AIC).

2/ A frequent case is when the parameter  $\theta$  identifying an affine causal model  $X_t = M_\theta^t \xi_t + f_\theta^t$  can be decomposed as  $\theta = (\theta_1, \theta_2)'$  with  $f_\theta^t = \widehat{f}_{\theta_1}^t$  and  $M_\theta^t = \widehat{M}_{\theta_2}^t$ . Let  $p_1, p_2$

such that  $p_1 = |\theta_1|$ ,  $p_2 = |\theta_2|$  and  $|m| = p_1 + p_2$ .

In such a case, from (4.6.1), it is clear that all the terms  $F(\theta_m^*)_{i,j}$  and  $G(\theta_m^*)_{i,j}$  are equals to zero for  $i = 1, \dots, p_1$  and  $j = 1, \dots, p_2$  implying

$$F(\theta_m^*) = - \begin{pmatrix} A_{1,p_1} & O_{p_1,p_2} \\ O_{p_2,p_1} & B_{p_1+1,p_1+p_2} \end{pmatrix} \text{ and } G(\theta_m^*) = \begin{pmatrix} A_{1,p_1} & O_{p_1,p_2} \\ O_{p_2,p_1} & \frac{(\mu_4-1)}{2} B_{p_1+1,p_1+p_2} \end{pmatrix}$$

where  $O$  is the null matrix and from the expressions of matrix  $G(\theta_m^*)$  and  $F(\theta_m^*)$  in (4.6.1),

$$A_{1,p_1} = \left( \mathbb{E} \left[ \frac{\partial_{\theta_i} f_{\theta_m^*}^0 \partial_{\theta_j} f_{\theta_m^*}^0}{H_{\theta_m^*}^0} \right] \right)_{1 \leq i,j \leq p_1} \text{ and } B_{p_1+1,p_1+p_2} = \left( \frac{1}{2} \mathbb{E} \left[ \frac{\partial_{\theta_i} H_{\theta_m^*}^0 \partial_{\theta_j} H_{\theta_m^*}^0}{(H_{\theta_m^*}^0)^2} \right] \right)_{p_1+1 \leq i,j \leq p_1+p_2}.$$

As a consequence,

$$\begin{aligned} G(\theta_m^*) F(\theta_m^*)^{-1} &= -\text{Diag} \left( A_{1,p_1}, \frac{(\mu_4-1)}{2} B_{p_1+1,p_1+p_2} \right) \times \text{Diag} (A_{1,p_1}^{-1}, B_{p_1+1,p_1+p_2}^{-1}) \\ &= -\text{Diag} \left( I_{p_1}, \frac{(\mu_4-1)}{2} I_{p_2} \right) \end{aligned}$$

and we obtain

$$-2 \text{Trace} \left( (F(\theta_m^*))^{-1} G(\theta_m^*) \right) = 2 p_1 + (\mu_4 - 1) p_2. \quad (4.6.2)$$

This setting includes many classical times series:

- For ARMA( $p, q$ ) processes, we have  $X_t = f_{\theta}^t + \sigma \xi_t$  since  $X_t + a_1 X_{t-1} + \dots + a_p X_{t-p} = \sigma (\xi_t + b_1 \xi_{t-1} + \dots + b_q \xi_{t-q})$  for all  $t \in \mathbb{Z}$ . Then  $\theta_1 = (a_1, \dots, a_p, b_1, \dots, b_q)$  and  $\theta_2 = \sigma$ . The penalty term is slightly different according to  $\sigma$  is known or not:

- (a) if  $\sigma$  is known, then  $\theta = \theta_1$  and  $G(\theta^*) = -F(\theta^*)$ , so that we recover exactly the AIC penalty term:

$$-2 \text{Trace} (G(\theta_m^*) F(\theta_m^*)^{-1}) = 2 |m| = 2 (p + q);$$

- (b) Otherwise,  $\theta = (\theta_1, \sigma)$  and simple computations lead to

$$F(\theta^*) = \begin{pmatrix} (F(\theta^*))_{1 \leq i,j \leq |m|-1} & 0 \\ 0 & -\frac{1}{2\sigma^4} \end{pmatrix} \text{ and } G(\theta^*) = \begin{pmatrix} (G(\theta^*))_{1 \leq i,j \leq |m|-1} & 0 \\ 0 & \frac{(\mu_4-1)}{4\sigma^4} \end{pmatrix}$$

where  $(G(\theta^*))_{1 \leq i,j \leq |m|-1} = -(F(\theta^*))_{1 \leq i,j \leq |m|-1}$ .

Thus, we obtain  $G(\theta^*) F(\theta^*)^{-1} = - \begin{pmatrix} I_{1 \leq i,j \leq |m|-1} & 0 \\ 0 & \frac{\mu_4-1}{2} \end{pmatrix}$  and therefore, with  $|m| = p + q + 1$  in this case,

$$-2 \text{Trace} (G(\theta_m^*) F(\theta_m^*)^{-1}) = 2 |m| + (\mu_4 - 3) = 2(p + q) + (\mu_4 - 1),$$

and therefore once again the expectation of the ideal penalty leads to the AIC model selection.

- For GARCH( $p, q$ ) processes (see [Francq and Zakoian \(2010\)](#)), we have  $f_{\theta} = 0$  and  $X_t = M_{\theta}^t \xi_t$  since for any  $t \in \mathbb{Z}$ ,

$$\begin{cases} X_t &= \sigma_t \xi_t \\ \sigma_t^2 &= \omega_0 + a_1 X_{t-1}^2 + \dots + a_p X_{t-p}^2 + b_1 \sigma_{t-1}^2 + \dots + b_q \sigma_{t-q}^2 \end{cases}.$$

Denote  $\theta = \theta_2 = (\omega_0, a_1, \dots, a_p, b_1, \dots, b_q)$ .

Then we have  $A_{p_1} = 0$  and therefore  $G(\theta^*) = -\frac{(\mu_4-1)}{2} F(\theta^*)$ . As a result:

$$-2 \text{Trace} (G(\theta_m^*) F(\theta_m^*)^{-1}) = (\mu_4 - 1) |m| = (\mu_4 - 1) (p + q + 1).$$

- For APARCH( $\delta, p, q$ ) processes (see Ding et al. (1993)), we also have  $f_\theta = 0$  and  $X_t = M_\theta^t \xi_t$  since for any  $t \in \mathbb{Z}$ ,

$$\begin{cases} X_t &= \sigma_t \xi_t \\ \sigma_t^\delta &= \omega_0 + a_1 (X_{t-1} - \gamma_1 |X_{t-1}|)^\delta + \cdots + a_p (X_{t-p} - \gamma_p |X_{t-p}|)^\delta \\ &\quad + b_1 \sigma_{t-1}^\delta + \cdots + b_q \sigma_{t-q}^\delta \end{cases}.$$

For such a process,  $\theta = \theta_2 = (\omega_0, a_1, \dots, a_p, \gamma_1, \dots, \gamma_p, b_1, \dots, b_q)$  when we assume that  $\delta$  is known, and, *mutatis mutandis*, the result is the same than for GARCH processes:

$$-2 \text{Trace}(G(\theta_m^*) F(\theta_m^*)^{-1}) = (\mu_4 - 1) |m| = (\mu_4 - 1) (2p + q + 1).$$

3/ Otherwise, the computations are no longer easy. Let us see the example of the family of  $AR(1) - ARCH(p)$  processes. Then for any  $t \in \mathbb{Z}$  we have  $X_t = \phi X_{t-1} + Z_t$  where  $Z_t = \xi_t (\alpha_0 + \alpha_1 Z_{t-1}^2 + \cdots + \alpha_p Z_{t-p}^2)^{1/2}$ . As a consequence, with  $\theta = (\phi, \alpha_0, \dots, \alpha_p)'$ , we obtain for any  $t \in \mathbb{Z}$ ,

$$X_t = f_\theta(X_{t-1}) + M_\theta(X_{t-1}, \dots, X_{t-p-1}) \xi_t$$

with 
$$\begin{cases} f_\theta(X_{t-1}) &= \phi X_{t-1} \\ M_\theta(X_{t-1}, \dots, X_{t-p}) &= (\alpha_0 + \sum_{i=1}^p \alpha_i (X_{t-i} - \phi X_{t-i-1})^2)^{1/2} \end{cases}.$$

Thus the parameter  $\phi$  is present in  $f_\theta$  as well as in  $M_\theta$ . From (4.6.1), and with the notations of 1/, we obtain:

$$F(\theta_m^*) = - \begin{pmatrix} A_{1,1} & O_{1,p+1} \\ O_{p+1,1} & O_{p+1,p+1} \end{pmatrix} - B_{1,p+2}$$

and 
$$G(\theta_m^*) = \begin{pmatrix} A_{1,1} & O_{1,p+1} \\ O_{p+1,1} & O_{p+1,p+1} \end{pmatrix} + \frac{(\mu_4 - 1)}{2} B_{1,p+2}.$$

As a consequence,

$$G(\theta_m^*) = -\frac{(\mu_4 - 1)}{2} F(\theta_m^*) + \frac{(\mu_4 - 3)}{2} \begin{pmatrix} A_{1,1} & O_{1,p+1} \\ O_{p+1,1} & O_{p+1,p+1} \end{pmatrix}.$$

Thus, with  $|m| = p + 2$ ,

$$G(\theta_m^*) F^{-1}(\theta_m^*) = -\frac{(\mu_4 - 1)}{2} I_{|m|} + \frac{(\mu_4 - 3)}{2} \begin{pmatrix} A_{1,1} & O_{1,p+1} \\ O_{p+1,1} & O_{p+1,p+1} \end{pmatrix} F^{-1}(\theta_m^*).$$

Whatever the matrix  $F^{-1}(\theta_m^*)$ , we have  $\begin{pmatrix} A_{1,1} & O_{1,p+1} \\ O_{p+1,1} & O_{p+1,p+1} \end{pmatrix} F^{-1}(\theta_m^*) = \begin{pmatrix} c(\theta_m^*) & O_{1,p+1} \\ O_{p+1,1} & O_{p+1,p+1} \end{pmatrix}$  with  $c(\theta_m^*) = c(\theta^*) \in \mathbb{R}$  since  $m^* \subset m$ . Then for all  $m^* \subset m$ ,

$$-2 \text{Trace}(G(\theta_m^*) F(\theta_m^*)^{-1}) = -2 c(\theta^*) + (\mu_4 - 1) |m|,$$

where  $-2 c(\theta^*)$  does not depend on  $m$ .

## 4.7 Numerical Studies

This section aims to investigate how well are our new model selection criteria based on the theoretical results obtained in both Section 4.4 and Section 4.5.

To do that, three Data Generating Process (DGP) have been considered:

$$\begin{array}{lll} \text{DGP I} & \text{AR}(2) & X_t = 0.4 X_{t-1} + 0.4 X_{t-2} + \xi_t, \\ \text{DGP II} & \text{ARMA}(1, 1) & X_t - 0.5 X_{t-1} = \xi_t + 0.6 \xi_{t-1}, \\ \text{DGP III} & \text{GARCH}(1, 1) & X_t = \sigma_t \xi_t \quad \text{with } \sigma_t^2 = (1 + 0.35 X_{t-1}^2 + 0.4 \sigma_{t-1}^2)^{1/2}, \end{array}$$

where  $(\xi_t)_t$  is a Gaussian white noise with variance unity.

In order to identify possible discrepancies between asymptotically expected results and those obtained at finite distance, we consider  $n$  belongs to 200, 500, 1000, 2000. We will compare the performance of the AIC, BIC and KC'.

For the efficiency property (Theorem 1.2), we will be interested in the difference

$$ME := n \left( \ell(\hat{\theta}_{\hat{m}}, \theta^*) - \inf_{m \in \mathcal{M}} \left\{ \ell(\hat{\theta}_m, \theta^*) \right\} \right)$$

and the idea would be to see the decay rate of this residual term.

The considered family of competitive models is the same for the three DGP

$$\mathcal{M} = \{ \text{ARMA}(p, q) \text{ or } \text{GARCH}(p, q) \quad \text{with} \quad 0 \leq p, q \leq 6 \}.$$

The empirical estimates of the ratio ME are obtained based on 500 replications. Here is how we empirically estimated ME.

Given a DGP and a list of candidate models, the empirical estimates of  $\ell(\hat{\theta}_{\hat{m}}, \theta^*)$  have been obtained in the following way:

1. For any model  $m$ , the observations  $(X_1, \dots, X_n)$  from DGP have been used to estimated the QMLE  $\hat{\theta}_m$  and then we considered as an estimate of  $R(\hat{\theta}_m)$  the quantity  $\hat{\gamma}_n(\hat{\theta}_m)$  which is the empirical risk of the estimator  $\hat{\theta}_m$  but computed over observations  $(Y_1, \dots, Y_N)$  (with  $N = 10^6$ ) arisen from the DGP and independent to  $(X_1, \dots, X_n)$ .
2. Selection of  $\hat{m}$  according to (4.2.9) and then we retrieve the estimation of  $R(\hat{\theta}_{\hat{m}})$  in the list of estimates obtained in the first point;
3. Then take the difference of the results obtained in 1 and 2. The experiment was repeated 500 times and the average of these differences times  $n$  is our estimate of ME.

For the consistency property, the Table 4.1 shows the frequency of selecting respectively the true model versus a model other than the true one (called here "wrong").

**Table 4.1:** Percentage of selected order based on 500 replications depending on sample's length for DGP I, DGP II and DGP III respectively.

Sample length	200			500			1000			2000		
	AIC	BIC	KC'	AIC	BIC	KC'	AIC	BIC	KC'	AIC	BIC	KC'
True	17.2	36.2	35.6	30.4	73.2	78.2	36.4	87.4	92.2	32.4	96.2	98.4
Wrong	82.8	63.8	64.4	69.6	26.8	21.8	63.6	13.6	7.8	67.6	3.8	1.6
True	27.8	80.8	92.0	30.6	88.4	96.6	31.0	89.1	97.5	33.3	95.2	99.9
Wrong	72.2	19.2	8.0	69.7	11.6	3.4	69.0	10.9	2.5	66.7	4.8	0.1
True	0.4	10.8	14.8	1.4	32.2	55.8	1.0	54.8	82.0	2.0	75.8	93.8
Wrong	99.6	89.2	85.2	98.6	67.8	44.2	99.0	45.2	18.0	98.0	24.2	6.2

From these results, it follows that KC' outperforms BIC when dealing both with small and larges samples. These results confirm that it is important to consider also the neglected terms in the derivation of the BIC criterion.

Moreover, the percentage of selecting the true model using both BIC and KC' approaches 100 with increasing  $n$ . This is consistent with our asymptotic result (Corollary 4.2). Hence, KC' is more robust to the sample size and thus improves BIC.

**Table 4.2:** ME estimates based on 500 replications depending on sample's length.

Sample length	200			500			1000			2000		
	AIC	BIC	KC'	AIC	BIC	KC'	AIC	BIC	KC'	AIC	BIC	KC'
DGP I	4.91	2.59	5.35	3.46	1.11	1.18	3.08	0.98	0.75	3.05	0.38	0.29
DGP II	3.66	0.87	0.54	3.37	0.42	0.11	2.62	0.15	0.05	2.5	0.10	0.04
DGP III	2.39	4.63	13.16	2.53	4.08	9.54	2.69	2.96	2.52	3.21	2.06	0.76

In view of the results of Table 4.2, we notice a decrease of the residual term  $\widehat{ME}$  for all the criteria. This decrease is much faster and tends towards 0 for the consistent criteria while it is much less for the AIC. This allows us to conclude that Theorem 4.3 is verified by AIC criterion. Moreover, Theorem 4.3 is also verified by both criteria BIC and KC' despite the fact that these criteria do not satisfy the conditions of the theorem. One might think that the validity of (4.4.8) is due to the fact that since  $m^* \in \mathcal{M}$ ,  $\inf_{m \in \mathcal{M}} \{\ell(\widehat{\theta}_m, \theta^*)\}$  is reached for  $m = m^*$  in most cases and to the consistency of the BIC and KC' criteria. Thus, one could say that in configurations where the true model is part of the candidate models, the consistent criteria are also efficient.

## 4.8 Proofs

### 4.8.1 Proofs of Section 4.3

The asymptotic normality of  $(\frac{1}{n}(\partial_{\theta_i} L_n(\theta_m^*)))_{i \in m}$  was established in [Bardet and Wintenberger \(2009\)](#) and [Bardet et al. \(2020b\)](#) when  $m^* \subset m$  using a central limit theorem for stationary martingale difference. Here we extend this result to any  $m \in \mathcal{M}$ :



**Proposition 4.3.** *Under Assumption A0-A5, for any  $\theta \in \Theta$ , we have*

$$\sqrt{n} \left( \frac{1}{n} \partial_{\theta} L_n(\theta) + \frac{1}{2} \mathbb{E}[\partial_{\theta} \gamma(\theta, X_0)] \right) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \Sigma(\theta))$$

$$\text{with } \Sigma(\theta) := \frac{1}{4} \left( \sum_{t \in \mathbb{Z}} \text{Cov}(\partial_{\theta_i} \gamma(\theta, X_0), \partial_{\theta_j} \gamma(\theta, X_t)) \right)_{1 \leq i, j \leq d}. \quad (4.8.1)$$

The main tool we use here for establishing Theorem 4.3 is the notion of  $\tau$ -dependence for stationary time series. More precisely, the  $\tau$ -dependence coefficients, which are a version of the coupling coefficients introduced in [Dedecker and Prieur \(2004\)](#) and used for stationary infinite memory chains. The reader is deferred to the lecture notes [Dedecker et al. \(2007\)](#) for complements and details on coupling, based on the Wasserstein distance between probabilities defined as below. Its stationary version is: Let  $(\Omega, \mathcal{C}, \mathbb{P})$  be a probability space,  $\mathcal{M}$  a  $\sigma$ -subalgebra of  $\mathcal{C}$  and  $Z$  a random variable with values in  $E$ . Assume that  $\|Z\|_p < \infty$  and define the coefficient  $\tau^{(p)}$  as

$$\tau^{(p)}(\mathcal{M}, Z) = \left\| \sup_{f \in \Lambda_1(E)} \left\{ \left| \int f(x) \mathbb{P}_{Z|\mathcal{M}}(dx) - \int f(x) \mathbb{P}_Z(dx) \right| \right\} \right\|_p.$$

Using the definition of  $\tau$ , the dependence between the past of the sequence  $(Z_t)_{t \in \mathbb{Z}}$  and its future  $k$ -tuples may be assessed: consider the norm  $\|x - y\| = \|x_1 - y_1\| + \dots + \|x_k - y_k\|$  on  $E^k$ , set  $\mathcal{M}_p = \sigma(Z_t, t \leq p)$  and define

$$\tau_Z^{(p)}(s) = \sup_{k > 0} \left\{ \max_{1 \leq l \leq k} \frac{1}{l} \sup \left\{ \tau^{(p)}(\mathcal{M}_p, (Z_{j_1}, \dots, Z_{j_l})) \text{ with } p + s \leq j_1 < \dots < j_l \right\} \right\}.$$

Finally, the time series  $(Z_t)_{t \in \mathbb{Z}}$  is  $\tau_Z^{(p)}$ -weakly dependent when its coefficients  $\tau_Z^{(p)}(s)$  tend to 0 as  $s$  tends to infinity.

**Lemma 4.2.** *Under Assumption A0, then for  $p \leq r$  and  $b_k^{(p)} = \|\xi_0\|_p \alpha_k(M_{\theta}, \Theta) + \alpha_k(f_{\theta}, \Theta)$  for any  $j \in \mathbb{N}^*$ ,*

$$\tau_X^{(p)}(s) \leq C \lambda_s \quad \text{with} \quad \lambda_s = \inf_{1 \leq r \leq s} \left\{ \left( \sum_{k=1}^{\infty} b_k^{(p)} \right)^{s/r} + \sum_{t=r+1}^{\infty} b_t^{(p)} \right\} \quad \text{for } s \geq 1. \quad (4.8.2)$$

*Proof of Lemma 4.2.* This Lemma can be directly deduced from Proposition 3.1 of [Doukhan and Wintenberger \(2008\)](#) where  $F(x, \xi_0) = M_{\theta}(x) \xi_0 + f_{\theta}(x)$  for any  $x \in \mathbb{R}^{\infty}$  and therefore

$$\|F(x, \xi_0) - F(y, \xi_0)\|_p \leq \|\xi_0\|_p |M_{\theta}(x) - M_{\theta}(y)| + |f_{\theta}(x) - f_{\theta}(y)|$$

inducing  $\|F(x, \xi_0) - F(y, \xi_0)\|_p \leq \sum_{k=1}^{\infty} b_k^{(p)}$  with  $b_k^{(p)} = \|\xi_0\|_p \alpha_k(M_{\theta}, \Theta) + \alpha_k(f_{\theta}, \Theta)$ .  $\square$

**Remark 4.5.** *Using Assumption A0 and A5, we deduce that  $b_t^{(p)} = O(t^{-\delta})$  with  $\delta > 7/2$ , and therefore  $\tau_X^{(p)}(s) \leq \lambda_s = O((s^{-1} \log s)^{\delta-1})$ .*

Now, under Assumption A0, since  $X$  is a causal time series, define for any  $j = 1, \dots, d$  and  $\theta \in \Theta$ ,

$$\phi_{\theta}^{(j)}((X_{t-k})_{k \geq 0}) := \partial_{\theta_j} \gamma(\theta, X_t) = -2 \partial_{\theta_j} M_{\theta}^t \frac{(X_t - f_{\theta}^t)^2}{(M_{\theta}^t)^3} - 2 \partial_{\theta_j} f_{\theta}^t \frac{X_t - f_{\theta}^t}{(M_{\theta}^t)^2} + 2 \frac{\partial_{\theta_j} M_{\theta}^t}{M_{\theta}^t}.$$

Then we have:

**Lemma 4.3.** *Under Assumption A0-A5, for any  $j = 1, \dots, d$ , for any  $\theta \in \Theta$ , the sequence  $(\phi_\theta^{(j)}((X_{t-k})_{k \geq 0}))_{t \in \mathbb{Z}}$  is a causal stationary sequence that is  $\tau_{\phi_\theta^{(j)}}^{(p)}$ -weakly dependent where its coefficients  $\tau_{\phi_\theta^{(j)}}^{(1)}(s)$  satisfies:*

$$\begin{aligned} \tau_{\phi_\theta^{(j)}}^{(1)}(s) &\leq C \left( \sum_{\ell=1}^s (\alpha_\ell(f_\theta, \Theta) + \alpha_\ell(M_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} M_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} f_\theta, \Theta)) \lambda_{s+1-\ell} \right. \\ &\quad \left. + \sum_{\ell=s+1}^{\infty} (\alpha_\ell(f_\theta, \Theta) + \alpha_\ell(M_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} M_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} f_\theta, \Theta)) \right), \end{aligned} \quad (4.8.3)$$

for any  $s \geq 0$  where  $(\lambda_s)$  is defined in (4.8.2).

*Proof of Lemma 4.3.* In the proof of Proposition 4.1 of Bardet et al. (2020a), it has been proven for  $U = (U_i)_{i \geq 1}$  and  $V = (V_i)_{i \geq 1}$  such as  $\sup_{i \geq 1} \{\|U_i\|_4 \vee \|V_i\|_4\} < \infty$  that there exists  $C > 0$  satisfying

$$\begin{aligned} \mathbb{E} \left[ \sup_{\theta \in \Theta} \|\phi_\theta^{(j)}(U) - \phi_\theta^{(j)}(V)\| \right] &\leq C \left( \|U_1 - V_1\|_4 \right. \\ &\quad \left. + \sum_{i=2}^{\infty} (\alpha_i(f_\theta, \Theta) + \alpha_i(M_\theta, \Theta) + \alpha_i(\partial_{\theta_j} f_\theta, \Theta) + \alpha_i(\partial_{\theta_j} M_\theta, \Theta)) \|U_i - V_i\|_4 \right). \end{aligned} \quad (4.8.4)$$

Using coupling techniques, if  $(\tilde{\xi}_t)_{t \in \mathbb{Z}}$  is an independent replication of  $(\xi_t)_{t \in \mathbb{Z}}$ , define also  $(\tilde{X}_t)_{t \in \mathbb{Z}}$  satisfying Assumption with  $(\tilde{\xi}_t)_{t \in \mathbb{Z}}$  instead of  $(\xi_t)_{t \in \mathbb{Z}}$  and  $(\phi_\theta^{(j)}((\tilde{X}_{t-k})_{k \geq 0}))_{t \in \mathbb{Z}}$ . Then for  $s \geq 0$ , using (4.8.4),

$$\begin{aligned} \tau_{\phi_\theta^{(j)}}^{(1)}(s) &\leq \|\phi_\theta^{(j)}((X_{s-k})_{k \geq 0}) - \phi_\theta^{(j)}((\tilde{X}_{s-k})_{k \geq 0})\|_1 \\ &\leq C \left( \|X_1 - \tilde{X}_1\|_4 \right. \\ &\quad \left. + \sum_{i=2}^{\infty} (\alpha_i(f_\theta, \Theta) + \alpha_i(M_\theta, \Theta) + \alpha_i(\partial_{\theta_j} f_\theta, \Theta) + \alpha_i(\partial_{\theta_j} M_\theta, \Theta)) \|X_i - \tilde{X}_i\|_4 \right) \\ &\leq C \sum_{\ell=1}^{\infty} (\alpha_\ell(f_\theta, \Theta) + \alpha_\ell(M_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} f_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} M_\theta, \Theta)) \lambda_{s+1-\ell}, \end{aligned}$$

that implies (4.8.3).  $\square$

**Remark 4.6.** *Under Assumption A0 and A5, and therefore with  $\lambda_s = O(s^{1-\delta} \log s)$  with  $\delta > 7/2$ , we also deduce that  $\tau_{\phi_\theta^{(j)}}^{(1)}(s) = O(s^{1-\delta} \log s)$ .*

*Proof of Proposition 4.3.* If  $Z$  is a  $\tau_Z$ -dependent centered stationary time series satisfying  $\mathbb{E}[|Z_0|^\kappa] < \infty$  with  $\kappa > 2$ , and  $\sum_{s=1}^{\infty} s^{1/(\kappa-2)} \tau_Z(s) < \infty$ , we deduce from Lemma 2, point 2. of Dedecker and Doukhan (2003) that condition D(2,  $\theta/2$ ,  $X$ ) is satisfied as  $\theta$ -weakly dependent coefficients are smaller than  $\tau$ -weakly dependent coefficients, see (2.2.13) p.16 of Dedecker et al. (2007), and  $0 < \sum_{t \in \mathbb{Z}} |\mathbb{E}[Z_0 Z_t]| < \infty$  from Proposition 2 of Dedecker and Doukhan (2003). Then,

$$\frac{1}{\sqrt{n}} \sum_{t=1}^n Z_t \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}\left(0, \sum_{t \in \mathbb{Z}} \mathbb{E}[Z_0 Z_t]\right).$$

We can apply this central limit theorem to

$$Z_t := \sum_{j=1}^d c_j (\phi_\theta^{(j)}((X_{t-k})_{k \geq 0}) - \mathbb{E}[\phi_\theta^{(j)}((X_{t-k})_{k \geq 0})]) \quad \text{with } (c_j)_{1 \leq j \leq d} \in \mathbb{R}^d.$$

Indeed, using Lemma 4.3, we easily obtain for  $s \geq 0$

$$\tau_Z(s) \leq C \left( \sum_{j=1}^d |c_j| \right) \sum_{\ell=1}^{\infty} (\alpha_\ell(f_\theta, \Theta) + \alpha_\ell(M_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} f_\theta, \Theta) + \alpha_\ell(\partial_{\theta_j} M_\theta, \Theta)) \lambda_{s+1-\ell},$$

and therefore under Assumption A0 and A5,  $\tau_Z(s) = O(s^{1-\delta} \log s)$ . Moreover, using Lemma 4.6, we deduce  $\mathbb{E}[|Z_0|^{8/3}] < \infty$ . Then with  $\kappa = 8/3$ ,  $\sum_{s=1}^{\infty} s^{1/(\kappa-2)} \tau_Z(s) = \sum_{s=1}^{\infty} s^{3/2} \tau_Z(s) < \infty$  is satisfied since  $\delta > 7/2$ .

Therefore, we deduce for any  $\theta \in \Theta$ ,

$$\begin{aligned} \sqrt{n} \sum_{j=1}^d c_j \left( \frac{1}{n} \partial_{\theta_j} L_n(\theta) + \frac{1}{2} \mathbb{E}[\partial_{\theta_j} \gamma(\theta, X_0)] \right) \\ \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}\left(0, \frac{1}{4} \sum_{i=1}^d \sum_{j=1}^d c_i c_j \sum_{t \in \mathbb{Z}} \text{Cov}(\partial_{\theta_i} \gamma(\theta, X_0), \partial_{\theta_j} \gamma(\theta, X_t))\right), \end{aligned}$$

which implies the multidimensional central limit theorem (4.8.1).  $\square$

*Proof of Corollary 4.1.* Firstly, it was already established in Bardet et al. (2020b) that if  $m^* \subset m$  then  $(\partial_{\theta_i} \gamma(\theta_m, X_t))_{t \in \mathbb{Z}}$  is a stationary martingale increments with respect to  $\mathcal{F}_t = \sigma((X_{t-k})_{k \in \mathbb{N}})$ . As a consequence  $\text{Cov}(\partial_{\theta_i} \gamma(\theta, X_0), \partial_{\theta_j} \gamma(\theta, X_t)) = 0$  if  $t \neq 0$ .

Secondly, for all  $m \in \mathcal{M}$ , from the definition of  $\theta_m^*$  as a local minimum of  $R$  on  $\Theta_m$ , and from Assumption A0-A5, then  $\partial_{\theta_j} R(\theta_m^*) = \mathbb{E}[\partial_{\theta_j} \gamma(\theta_m^*, X_0)] = 0$  for all  $j \in m$ .  $\square$

*Proof of Theorem 4.1.* We use here the standard proof allowing to show the asymptotic normality of the QMLE and already used in Bardet and Wintenberger (2009).

Firstly, it was established in Bardet et al. (2020b) that  $\hat{\theta}_m \xrightarrow[n \rightarrow +\infty]{a.s.} \theta_m^*$ .

Secondly, a Taylor-Lagrange expansion is applied to  $(\partial_{\theta_j} L_n(\hat{\theta}_m))_{j \in m}$  around  $\theta_m^*$ :

$$\begin{aligned} \frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\hat{\theta}_m))_{j \in m} &= \frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\theta_m^*))_{j \in m} \\ &\quad + \left( \frac{1}{n} \partial_{\theta_i \theta_j}^2 L_n(\bar{\theta}_m) \right)_{i,j \in m} \times \sqrt{n} ((\hat{\theta}_m)_i - (\theta_m^*)_i)_{i \in m} \end{aligned} \quad (4.8.5)$$

with  $\bar{\theta}_m = c \hat{\theta}_m + (1-c) \theta_m^*$  and  $0 < c < 1$ .

Using  $\hat{\theta}_m \xrightarrow[n \rightarrow +\infty]{a.s.} \theta_m^*$  and the ergodic theorem  $\frac{1}{n} \partial_{\theta_i \theta_j}^2 L_n(\theta_m) \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta_m)$  for any  $\theta_m \in \Theta_m$  since  $\mathbb{E}[\|\partial_{\theta^2}^2 \gamma(\theta, X_0)\|_\Theta] < \infty$ , we obtain:

$$\left( \frac{1}{n} \partial_{\theta_i \theta_j}^2 L_n(\bar{\theta}_m) \right)_{i,j \in m} \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta_m^*). \quad (4.8.6)$$

Finally, by definition of  $\hat{\theta}_m$ ,  $\partial_{\theta_j} \hat{L}_n(\hat{\theta}_m) = 0$  for any  $j \in m$ . As a consequence,

$$\frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\hat{\theta}_m))_{j \in m} \xrightarrow[n \rightarrow \infty]{\mathcal{P}} 0, \quad (4.8.7)$$

using a Markov Inequality and  $\mathbb{E}\left[\frac{1}{\sqrt{n}} \|\partial_\theta \widehat{L}_n(\theta) - \partial_\theta L_n(\theta)\|_\Theta\right] \xrightarrow{n \rightarrow \infty} 0$  established in (5.11) of [Bardet and Wintenberger \(2009\)](#). Considering (4.8.5), (4.8.6) and (4.8.7), and with the central limit theorem satisfied by  $\frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\theta_m^*))_{j \in m}$  provided in Corollary 4.1, this achieves the proof.  $\square$

Now, before establishing Proposition 4.1, three technical lemmas can be stated:

**Lemma 4.4.** *Under the assumptions **A0-A5**, with  $8/3 < r' \leq r/3$  and  $r' < 2(\delta - 1)$  where  $\delta > 7/2$  is given in Assumption A5, for any  $m \in \mathcal{M}$ , there exists  $C > 0$  such as for any  $n \in \mathbb{N}^*$*

$$\left\| \left( \frac{1}{\sqrt{n}} \partial_{\theta_j} L_n(\theta_m^*) \right)_{j \in m} \right\|_{r'} \leq C. \quad (4.8.8)$$

*Proof.* First, for any  $m \in \mathcal{M}$  and  $n \in \mathbb{N}^*$ ,

$$\begin{aligned} \left\| (\partial_{\theta_j} L_n(\theta_m^*))_{j \in m} \right\|^{r'} &\leq |m|^{r'/2-1} \sum_{j \in m} |\partial_{\theta_j} L_n(\theta_m^*)|^{r'} \\ &\leq \frac{|m|^{r'/2-1}}{2^{r'}} \sum_{j \in m} \left| \sum_{t=1}^n \partial_{\theta_j} \gamma(\theta_m^*, X_t) \right|^{r'}. \end{aligned} \quad (4.8.9)$$

Now, for all  $j \in m$ ,  $(\partial_{\theta_j} \gamma(\theta_m^*, X_t))_{t \in \mathbb{Z}}$  is a centered (from the proof of Corollary 4.1) stationary  $\tau_{\phi_{\theta}^{(j)}}^{(p)}$ -weakly dependent where its coefficients  $(\tau_{\phi_{\theta}^{(j)}}^{(1)}(s))_s$  satisfies (4.8.3) (see Lemma 4.3). Moreover, from the proof of Proposition 4.3,  $\tau_{\phi_{\theta}^{(j)}}^{(1)}(s) = O(s^{1-\delta} \log s)$ .

In Proposition 5.5 of [Dedecker et al. \(2007\)](#), since  $\mathbb{E}[|\partial_{\theta_j} \gamma(\theta_m^*, X_0)|^{r'}] < \infty$  from Lemma 4.6, it has been established that:

$$\begin{aligned} \mathbb{E}\left[\left|\sum_{t=1}^n \partial_{\theta_j} \gamma(\theta_m^*, X_t)\right|^{r'}\right] &\leq C_{r'} (M_{r',n} + M_{2,n}^{r'/2}) \\ \text{where } M_{m,n} &:= 2n \sum_{i=0}^{n-1} (i+1)^{m-2} \tau_{\phi_{\theta}^{(j)}}^{(1)}(i). \end{aligned}$$

Using  $8/3 < r' < 2(\delta - 1)$  with  $\delta > 7/2$ , we obtain that

$$M_{r',n} \leq C n \sum_{i=1}^n i^{r'-1-\delta} \log(i) \leq C' n^{1+r'-\delta} \log(n) = O(n^{r'/2})$$

and  $M_{2,n} \leq C n \sum_{i=1}^n i^{1-\delta} \log(i) \leq C'' n$ . As a consequence, there exists  $C > 0$  such as for any  $n \in \mathbb{N}^*$ ,

$$\mathbb{E}\left[\left|\sum_{t=1}^n \partial_{\theta_j} \gamma(\theta_m^*, X_t)\right|^{r'}\right] \leq C n^{r'/2}. \quad (4.8.10)$$

Then, using (4.8.9) and (4.8.10), the proof is established.  $\square$

**Lemma 4.5.** *Under the assumptions **A0-A5**, then for any  $m \in \mathcal{M}$ , there exists  $C > 0$  such as for any  $n \in \mathbb{N}^*$ ,*

$$\left\| \left( \frac{1}{\sqrt{n}} \partial_{\theta_j} L_n(\widehat{\theta}_m) \right)_{j \in m} \right\|_{r/3} \leq C.$$

*Proof.* First, from the definition of  $\hat{\theta}_m$ , we have  $\partial_{\theta_j} \hat{L}_n(\hat{\theta}_m) = 0$  for any  $j \in m$ . Then,

$$\begin{aligned} \left\| \left( \frac{1}{\sqrt{n}} \partial_{\theta_j} L_n(\hat{\theta}_m) \right)_{j \in m} \right\|_{r/3} &= \left\| \left( \frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\hat{\theta}_m) - \partial_{\theta_j} \hat{L}_n(\hat{\theta}_m)) \right)_{j \in m} \right\|_{r/3} \\ &\leq \frac{|m|^{(r-6)/2r}}{\sqrt{n}} \sum_{j \in m} \left\| \partial_{\theta_j} L_n(\hat{\theta}_m) - \partial_{\theta_j} \hat{L}_n(\hat{\theta}_m) \right\|_{r/3} \\ &\leq \frac{|m|^{1/2}}{2\sqrt{n}} \sum_{j \in m} \sum_{t=1}^n \left\| \partial_{\theta_j} \gamma(\hat{\theta}_m, X_t) - \partial_{\theta_j} \hat{\gamma}(\hat{\theta}_m, X_t) \right\|_{r/3}. \end{aligned}$$

From the proof of Lemma 2 in [Bardet et al. \(2020b\)](#), there exists  $C > 0$  such as

$$\begin{aligned} \mathbb{E} \left[ \sup_{\theta \in \Theta} \left\| \partial_{\theta_j} \gamma(\theta, X_t) - \partial_{\theta_j} \hat{\gamma}(\theta, X_t) \right\|^{r/3} \right] \\ \leq C \left( \sum_{k \geq t} \alpha_k(f_\theta, \Theta) + \alpha_k(M_\theta, \Theta) + \alpha_k(\partial f_\theta, \Theta) + \alpha_k(\partial M_\theta, \Theta) \right)^{r/3}. \end{aligned}$$

Therefore,

$$\begin{aligned} \left\| \left( \frac{1}{\sqrt{n}} \partial_{\theta_j} L_n(\hat{\theta}_m) \right)_{j \in m} \right\|_{r/3} &\leq C \frac{|m|^{3/2}}{2\sqrt{n}} \sum_{t=1}^n \sum_{k \geq t} (\alpha_k(f_\theta, \Theta) + \alpha_k(M_\theta, \Theta) + \alpha_k(\partial f_\theta, \Theta) + \alpha_k(\partial M_\theta, \Theta)) \\ &\leq \frac{C'}{\sqrt{n}} \sum_{t=1}^n \sum_{j \geq t} j^{-\delta} \leq \frac{C'''}{\sqrt{n}} \sum_{t=1}^n t^{1-\delta} \leq C''', \end{aligned}$$

with  $C' > 0$ ,  $C'' > 0$  and  $C''' > 0$  and where the last inequality holds since  $\delta > 7/2$  under Assumption A5.  $\square$

**Lemma 4.6.** *Under the assumptions A0-A5, for any  $m \in \mathcal{M}$  and any  $\theta \in \Theta_m$ ,*

$$\left\| (\partial_{\theta_j} \gamma(\theta, X_0))_{j \in m} \right\|_{r/3} < \infty. \quad (4.8.11)$$

*Proof.* For  $j \in m$ , we have for any  $\theta \in \Theta_m$ ,

$$\partial_{\theta_j} \gamma(\theta, X_0) = -2(M_\theta^0)^{-2}(X_0 - f_\theta^0) \partial_{\theta_j} f_\theta^0 - 2(M_\theta^0)^{-3}(X_0 - f_\theta^0)^2 \partial_{\theta_j} M_\theta^0 + 2(M_\theta^0)^{-1} \partial_{\theta_j} M_\theta^0.$$

Therefore with Assumption A3 and Minkowski Inequality,

$$\begin{aligned} \left\| (\partial_{\theta_j} \gamma(\theta, X_0))_{j \in m} \right\|_{r/3} &\leq \frac{2}{h^{3/2}} \left( h^{1/2} \left\| (\partial_{\theta_j} f_\theta^0)(X_0 - f_\theta^0) \right\|_{r/3} \right. \\ &\quad \left. + \left\| (\partial_{\theta_j} M_\theta^0) |X_0 - f_\theta^0|^2 \right\|_{r/3} + h \left\| (\partial_{\theta_p} M_\theta^0) \right\|_{r/3} \right). \end{aligned}$$

Now, applying the Hölder Inequality, we obtain that there exists  $C > 0$  such that for any  $\theta \in \Theta_m$ ,

$$\begin{aligned} \mathbb{E} \left[ \left| \partial_{\theta_p} \gamma(\theta^*, X_0) \right|^{r/3} \right] &\leq C \left( \left\| \partial_{\theta_j} f_\theta^0 \right\|_{2r/3} \left\| X_0 - f_\theta^0 \right\|_{2r/3} \right. \\ &\quad \left. + \left\| \partial_{\theta_j} M_\theta^0 \right\|_r \left\| X_0 - f_\theta^0 \right\|_r^2 + \left\| \partial_{\theta_p} M_\theta^0 \right\|_{r/3} \right). \quad (4.8.12) \end{aligned}$$

Using Assumption A0 and A5 and the proof of Lemma 1 in [Bardet and Wintenberger \(2009\)](#), all the right side terms in (4.8.12) are finite for any  $\theta \in \Theta_m$  and this achieves the proof.  $\square$

Then Proposition 4.1 can be established:

*Proof of Proposition 4.1.* From (4.8.5) and (4.8.6) with  $F(\theta_m^*)$  is a positive definite matrix, we know that for  $n$  large enough,

$$\begin{aligned} & \left\| \sqrt{n} ((\hat{\theta}_m)_i - (\theta_m^*)_i)_{i \in m} \right\|_{r'} \\ &= \left\| \left( \left( \frac{1}{n} \partial_{\theta_i \theta_j}^2 L_n(\bar{\theta}_m) \right)_{i,j \in m} \right)^{-1} \times \frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\hat{\theta}_m) - \partial_{\theta_j} L_n(\theta_m^*))_{j \in m} \right\|_{r'}. \end{aligned} \quad (4.8.13)$$

Therefore, using Hölder and Minkowski inequalities, we obtain:

$$\begin{aligned} \left\| \sqrt{n} ((\hat{\theta}_m)_i - (\theta_m^*)_i)_{i \in m} \right\|_{r'} &\leq \left\| \left( \left( \frac{1}{n} \partial_{\theta_i \theta_j}^2 L_n(\bar{\theta}_m) \right)_{i,j \in m} \right)^{-1} \right\|_{\frac{rr'}{r-3r'}} \\ &\quad \times \left\| \frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\hat{\theta}_m) - \partial_{\theta_j} L_n(\theta_m^*))_{j \in m} \right\|_{\frac{r}{3}} \\ &\leq \left\| \left( \left( \frac{1}{n} \partial_{\theta_i \theta_j}^2 L_n(\bar{\theta}_m) \right)_{i,j \in m} \right)^{-1} \right\|_{\frac{rr'}{r-3r'}} \\ &\quad \times \left( \left\| \frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\hat{\theta}_m))_{j \in m} \right\|_{\frac{r}{3}} + \left\| \frac{1}{\sqrt{n}} (\partial_{\theta_j} L_n(\theta_m^*))_{j \in m} \right\|_{\frac{r}{3}} \right). \end{aligned}$$

Now using Assumption A4, Lemmas 4.4 and 4.5, we deduce (4.3.7).  $\square$

#### 4.8.2 Proofs of Section 4.4

*Proof of Lemma 4.1.* **1.** From the assumptions, the function  $R : \theta \in \Theta \mapsto R(\theta)$  is a  $\mathcal{C}^2(\Theta)$  function and the Hessian matrix  $\partial_{\theta^2} R = F$  is a definite positive matrix. Therefore, from a Taylor-Lagrange expansion:

$$\begin{aligned} n(R(\hat{\theta}_m) - R(\theta_m^*)) &= n \left( R(\theta_m^*) + (\hat{\theta}_m - \theta_m^*)^\top \partial_{\theta} R(\theta_m^*) \right. \\ &\quad \left. + \frac{1}{2} (\hat{\theta}_m - \theta_m^*)^\top \partial_{\theta^2} R(\bar{\theta}) (\hat{\theta}_m - \theta_m^*) - R(\theta_m^*) \right) \\ &= \frac{1}{2} (\sqrt{n}(\hat{\theta}_m - \theta_m^*))^\top \partial_{\theta^2} R(\bar{\theta}) (\sqrt{n}(\hat{\theta}_m - \theta_m^*)), \end{aligned} \quad (4.8.14)$$

with  $\bar{\theta} = \theta_m^* + c(\hat{\theta}_m - \theta_m^*) \in \Theta_m$  since  $c \in [0, 1]$ . Using Lemma 4 of [Bardet and Wintenberger \(2009\)](#) and continuous mapping Theorem, we deduce that:

$$\partial_{\theta^2} R(\bar{\theta}) = F(\bar{\theta}) \xrightarrow[n \rightarrow \infty]{\mathcal{P}} -2F(\theta_m^*) \quad \text{and} \quad G(\bar{\theta}) \xrightarrow[n \rightarrow \infty]{\mathcal{P}} G(\theta_m^*). \quad (4.8.15)$$

Moreover, using the asymptotic normality of  $\hat{\theta}_m$  established in [Bardet and Wintenberger \(2009\)](#) and [Bardet et al. \(2020b\)](#), we have:

$$\sqrt{n}((\hat{\theta}_m)_i - (\theta_m^*)_i)_{i \in m} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}\left(0, (F(\theta_m^*))^{-1} G(\theta_m^*) (F(\theta_m^*))^{-1}\right). \quad (4.8.16)$$

As a consequence, with  $Z_n = (G(\bar{\theta}))^{-1/2} F(\bar{\theta}) \sqrt{n}((\hat{\theta}_m)_i - (\theta_m^*)_i)_{i \in m} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, I_{|m|})$ , we have

$$\begin{aligned} n(R(\hat{\theta}_m) - R(\theta_m^*)) &= -Z_n^\top (G(\bar{\theta}))^{1/2} (F(\bar{\theta}))^{-1} F(\bar{\theta}) (F(\bar{\theta}))^{-1} (G(\bar{\theta}))^{1/2} Z_n \\ &= -Z_n^\top (G(\bar{\theta}))^{1/2} (F(\bar{\theta}))^{-1} (G(\bar{\theta}))^{1/2} Z_n. \end{aligned}$$

Define  $U^*(m) := -Z^\top (G(\theta_m^*))^{1/2} (F(\theta_m^*))^{-1} (G(\theta_m^*))^{1/2} Z$  where  $Z \stackrel{\mathcal{D}}{\sim} \mathcal{N}(0, I_{|m|})$ . Then using (4.8.15) we obtain

$$n(R(\hat{\theta}_m) - R(\theta_m^*)) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} U^*(m).$$

The computation of the expectation of  $U_m^*$  follows from

$$\begin{aligned} \mathbb{E}[U_m^*] &= \mathbb{E}[\text{Trace}(U_m^*)] = -\text{Trace}\left((G(\theta_m^*))^{1/2} (F(\theta_m^*))^{-1} (G(\theta_m^*))^{1/2}\right) \\ &= -\text{Trace}\left((F(\theta_m^*))^{-1} G(\theta_m^*)\right). \end{aligned}$$

Finally, for establishing  $\mathbb{E}[n(R(\hat{\theta}_m) - R(\theta_m^*))] \xrightarrow[n \rightarrow \infty]{} \mathbb{E}[U_m^*]$ , we have to prove that there exists  $n_0 \in \mathbb{N}$  such as

$$\sup_{n \geq n_0} \mathbb{E}[n|R(\hat{\theta}_m) - R(\theta_m^*)|] < \infty. \quad (4.8.17)$$

Indeed, from (4.8.14), we have:

$$\begin{aligned} n|R(\hat{\theta}_m) - R(\theta_m^*)| &\leq \frac{1}{2} \sup_{\theta \in \Theta} \|\partial_{\theta^2}^2 R(\theta)\| \|\sqrt{n}(\hat{\theta}_m - \theta_m^*)\|^2 \\ \implies \mathbb{E}[n|R(\hat{\theta}_m) - R(\theta_m^*)|] &\leq \frac{\lambda_{\max}}{2} \mathbb{E}[\|\sqrt{n}(\hat{\theta}_m - \theta_m^*)\|^2], \end{aligned} \quad (4.8.18)$$

since there exists  $\lambda_{\max} < \infty$  such as  $\|\partial_{\theta^2}^2 R(\theta)\| \leq \lambda_{\max}$  for any  $\theta \in \Theta$  from Assumption  $A_2(f_\theta, \Theta)$  and  $A_2(M_\theta, \Theta)$  where  $\Theta$  is a compact set.

Using Proposition 4.1, we know that

$$\sup_{n \in \mathbb{N}^*} \|\sqrt{n}(\hat{\theta}_m - \theta_m^*)\|_2 < \infty.$$

Finally using (4.8.18), we deduce (4.8.17).

**2.** As in the proof of **1.**, we use a Taylor-Lagrange expansion of  $\hat{\gamma}_n(\theta_m^*)$  around  $\hat{\theta}_m$  since  $\partial_\theta \hat{\gamma}_n(\hat{\theta}_m) = 0$ . Then,

$$n(\hat{\gamma}_n(\theta_m^*) - \hat{\gamma}_n(\hat{\theta}_m)) = \frac{1}{2} \sqrt{n}(\hat{\theta}_m - \theta_m^*)^\top (\partial_{\theta^2}^2 \hat{\gamma}_n(\bar{\theta}_m)) \sqrt{n}(\hat{\theta}_m - \theta_m^*).$$

But using  $\hat{\theta}_m \xrightarrow[n \rightarrow +\infty]{a.s.} \theta_m^*$  and  $\mathbb{E}\left[\left\|\frac{1}{n} L_n(\theta) - \frac{1}{n} \hat{L}_n(\theta)\right\|_\Theta\right] \xrightarrow[n \rightarrow \infty]{} 0$ , we have

$$(\partial_{\theta^2}^2 \hat{\gamma}_n(\bar{\theta}_m)) \xrightarrow[n \rightarrow \infty]{\mathcal{P}} -2F(\theta_m^*).$$

Therefore, using the same reasoning as in **1.**, we deduce that

$$n(\hat{\gamma}_n(\theta_m^*) - \hat{\gamma}_n(\hat{\theta}_m)) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} U^*(m).$$

With Hölder Inequality and using  $8/3 < r'$  defined in Proposition 4.1, we obtain for  $n$  large enough

$$\begin{aligned} \mathbb{E}[n(\hat{\gamma}_n(\theta_m^*) - \hat{\gamma}_n(\hat{\theta}_m))] &\leq \|(\partial_{\theta^2}^2 \hat{\gamma}_n(\bar{\theta}_m))^{1/2} \sqrt{n}(\hat{\theta}_m - \theta_m^*)\|_2^2 \\ &\leq \|(\partial_{\theta^2}^2 \hat{\gamma}_n(\bar{\theta}_m))\|_{\frac{r'}{r'-2}} \|\sqrt{n}(\hat{\theta}_m - \theta_m^*)\|_{r'}^2. \end{aligned} \quad (4.8.19)$$

Finally, with Proposition 4.1 and Lemma 4 of Bardet and Wintenberger (2009), we have  $\sup_{n \in \mathbb{N}^*} \mathbb{E}[\|\partial_{\theta^2}^2 \hat{\gamma}_n(\bar{\theta}_m)\|_{\Theta}^4] < \infty$  and  $\frac{r'}{r'-2} \leq 4$  since  $r' > 8/3$ , and therefore

$$\sup_{n \in \mathbb{N}^*} \mathbb{E}[n(\hat{\gamma}_n(\theta_m^*) - \hat{\gamma}_n(\bar{\theta}_m))] < \infty,$$

which concludes the proof.  $\square$

*Proof of Proposition 1.5.* The proof of this proposition can be deduced from

$$\mathbb{E}[n I_3(m)] = \mathbb{E}[n(R(\theta_m^*) - \hat{\gamma}_n(\theta_m^*))] = v_n^* \quad (4.8.20)$$

for any  $m \in \mathcal{M}$ . For establishing (4.8.20), we begin by

$$I_3(m) = (R(\theta_m^*) - \gamma_n(\theta_m^*)) + (\gamma_n(\theta_m^*) - \hat{\gamma}_n(\theta_m^*)) := I_{31}(m) + I_{32}(m). \quad (4.8.21)$$

Firstly, since  $\mathbb{E}[\gamma(\theta_m^*, X_0)] = R(\theta_m^*)$  and  $(X_t)_{t \in \mathbb{Z}}$  is a stationary times series, then for any  $n \in \mathbb{N}^*$ ,

$$\mathbb{E}[\gamma_n(\theta_m^*)] = \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\gamma(\theta_m^*, X_t)] = R(\theta_m^*) \implies \mathbb{E}[I_{31}(m)] = 0. \quad (4.8.22)$$

Secondly, from Assumption A0 and Bardet and Wintenberger (2009), there exists  $C > 0$  such that for any  $t \geq 1$

$$\mathbb{E}[\|\gamma(\theta, X_t) - \hat{\gamma}(\theta, X_t)\|_{\Theta}] \leq C \sum_{s \geq t} (\alpha_k(f_{\theta}, \Theta) + \alpha_k(M_{\theta}, \Theta)).$$

Therefore,

$$\begin{aligned} \mathbb{E}[\|\gamma_n(\theta_m^*) - \hat{\gamma}_n(\theta_m^*)\|_{\Theta}] &\leq \frac{C}{n} \sum_{t=1}^n \sum_{s \geq t} (\alpha_k(f_{\theta}, \Theta) + \alpha_k(M_{\theta}, \Theta)) \\ &\leq \frac{C}{n} \sum_{t=1}^n t^{1-\delta} \leq \frac{C'}{n}, \end{aligned} \quad (4.8.23)$$

since  $\delta > 7/2$  from Assumption A5. Moreover, for any  $m \in \mathcal{M}$  such as  $m^* \subset m$ , we have  $\gamma_n(\theta_m^*) - \hat{\gamma}_n(\theta_m^*) = \gamma_n(\theta_m^*) - \hat{\gamma}_n(\theta_m^*)$ . Using this and (4.8.23) we deduce that for any  $m \in \mathcal{M}$ , there exists a bounded sequence  $(v_n^*)_{n \in \mathbb{N}^*}$  not depending on  $m$  satisfying

$$\mathbb{E}[I_{32}(m)] = \frac{v_n^*}{n}. \quad (4.8.24)$$

Using also Lemma 4.1, this implies the asymptotic behavior of  $\mathbb{E}[\text{pen}_{id}(m)]$ .  $\square$

Now we establish a preliminary lemma that is an important step towards the proof of Theorem 4.3.

**Lemma 4.7.** *Let  $\text{pen} : m \in \mathcal{M}_n \mapsto \text{pen}(m) \in \mathbb{R}_+$ . Then*

$$\ell(\hat{\theta}_{\hat{m}_{\text{pen}}}, \theta^*) \leq \min_{m \in \mathcal{M}} \{\ell(\hat{\theta}_m, \theta^*)\} + (\text{pen}(\hat{m}_{id}) - \text{pen}(\hat{m}_{\text{pen}})) - (\text{pen}_{id}(\hat{m}_{id}) - \text{pen}_{id}(\hat{m}_{\text{pen}})). \quad (4.8.25)$$



*Proof.* By definition, for any  $m \in \mathcal{M}$ ,

$$\hat{C}_{\text{pen}_{id}}(m) = R(\hat{\theta}_m) = \ell(\hat{\theta}_m, \theta^*) + R(\theta^*). \quad (4.8.26)$$

As a consequence,

$$\min_{m \in \mathcal{M}} \{\ell(\hat{\theta}_m, \theta^*)\} = \ell(\hat{\theta}_{\hat{m}_{id}}, \theta^*) = \min_{m \in \mathcal{M}} \{\hat{C}_{\text{pen}_{id}}(m)\} + R(\theta^*). \quad (4.8.27)$$

For any  $m \in \mathcal{M}$ , we also have

$$\hat{C}_{\text{pen}}(m) = \hat{C}_{\text{pen}_{id}}(m) + \text{pen}(m) - \text{pen}_{id}(m).$$

By definition of  $\hat{m}_{\text{pen}}$ , we have  $\hat{C}_{\text{pen}}(\hat{m}_{\text{pen}}) \leq \hat{C}_{\text{pen}}(\hat{m}_{\text{pen}_{id}})$ . Therefore,

$$\begin{aligned} \hat{C}_{\text{pen}}(\hat{m}_{\text{pen}}) &\leq \hat{C}_{\text{pen}_{id}}(\hat{m}_{\text{pen}_{id}}) + \text{pen}(\hat{m}_{\text{pen}_{id}}) - \text{pen}_{id}(\hat{m}_{\text{pen}_{id}}) \\ C_{\text{pen}_{id}}(\hat{m}_{\text{pen}}) + \text{pen}(\hat{m}_{\text{pen}}) - \text{pen}_{id}(\hat{m}_{\text{pen}}) &\leq \hat{C}_{\text{pen}_{id}}(\hat{m}_{\text{pen}_{id}}) + \text{pen}(\hat{m}_{\text{pen}_{id}}) - \text{pen}_{id}(\hat{m}_{\text{pen}_{id}}). \end{aligned}$$

By replacing  $\hat{C}_{\text{pen}_{id}}(m)$  by  $\ell(\hat{\theta}_m, \theta^*)$  following (4.8.26) and using (4.8.27), then (4.8.25) is established.  $\square$

*Proof of Theorem 4.2.* Let  $\mathcal{M}^* = \{m \in \mathcal{M}, m^* \subset m\}$  and  $\mathcal{M}' = \mathcal{M} \setminus \mathcal{M}^*$ . Let  $m \in \mathcal{M}'$ . We have:

$$\begin{aligned} \mathbb{P}(\hat{m}_{\text{pen}} = m) &\leq \mathbb{P}(\hat{C}_{\text{pen}}(m) \leq \hat{C}_{\text{pen}}(m^*)) \\ &\leq \mathbb{P}\left\{(\hat{\gamma}_n(\hat{\theta}_m) - \hat{\gamma}_n(\hat{\theta}_{m^*})) \leq \text{pen}(m^*) - \text{pen}(m)\right\} \\ &\leq \mathbb{P}\left\{n(\hat{\gamma}_n(\hat{\theta}_m) - \hat{\gamma}_n(\theta_m^*)) + n(\hat{\gamma}_n(\theta_m^*) - R(\theta_m^*)) + n(R(\theta^*) - \hat{\gamma}_n(\theta^*)) \right. \\ &\quad \left. + n(\hat{\gamma}_n(\theta^*) - \hat{\gamma}_n(\hat{\theta}_{m^*})) \leq n(R(\theta^*) - R(\theta_m^*)) + n(\text{pen}(m^*) - \text{pen}(m))\right\} \\ &\leq \mathbb{P}\left\{Z_1 + Z_2 + Z_3 + Z_4 + Z_5 \leq -2n DK_L(\theta^* \|\theta_m^*)\right\} \end{aligned}$$

with  $Z_5 = n(\text{pen}(m) - \text{pen}(m^*))$  and with  $R(\theta^*) - R(\theta_m^*) = -2 DK_L(\theta^* \|\theta_m^*) < 0$  since  $m \not\subset m^*$  from [Bardet et al. \(2020b\)](#). Now, using  $\mathbb{P}(Z_1 + \dots + Z_5 \leq c) \leq \mathbb{P}(Z_1 \leq c/5) + \dots + \mathbb{P}(Z_5 \leq c/5)$  for any random variables  $Z_i$  and real number  $c$ , we obtain:

$$\mathbb{P}(\hat{m}_{\text{pen}} = m) \leq \sum_{i=1}^5 \mathbb{P}(Z_i \leq c_n), \quad (4.8.28)$$

where  $c_n = -\frac{2}{5} n DK_L(\theta^* \|\theta_m^*)$ .

Let  $Z_1 := n(\hat{\gamma}_n(\hat{\theta}_m) - \hat{\gamma}_n(\theta_m^*))$ . Following the same computations than in (4.8.19), with  $8/3 < r' \leq r/3$  and  $r' < 2(\delta - 1)$  defined in Proposition 4.1, and Hölder Inequality,

$$\begin{aligned} \mathbb{E}[|Z_1|^{\frac{3r'}{8}}] &\leq \|(\partial_{\theta^2}^2 \hat{\gamma}_n(\bar{\theta}_m))^{1/2} \sqrt{n}(\hat{\theta}_m - \theta_m^*)\|_{\frac{3r'}{4}}^{\frac{3r'}{4}} \\ &\leq \|(\partial_{\theta^2}^2 \hat{\gamma}_n(\bar{\theta}_m))\|_{\frac{3r'}{2}} \|\sqrt{n}(\hat{\theta}_m - \theta_m^*)\|_{r'}^{\frac{3r'}{4}}. \end{aligned}$$

Therefore, using Proposition 4.1,

$$\begin{aligned} \mathbb{P}(Z_1 \leq c_n) &\leq \mathbb{P}\left(|Z_1|^{\frac{3r'}{8}} \geq (|c_n|)^{\frac{3r'}{8}}\right) \leq \mathbb{E}[|Z_1|^{\frac{3r'}{8}}] \frac{1}{|c_n|^{\frac{3r'}{8}}} \\ &\implies \mathbb{P}(Z_1 \leq c_n) = O\left(\frac{1}{n^{\frac{3r'}{8}}}\right) = o\left(\frac{1}{n}\right), \quad (4.8.29) \end{aligned}$$

since  $3r'/8 > 1$ . The same kind of computations can also be done for  $Z_4 := n(\hat{\gamma}_n(\theta^*) - \hat{\gamma}_n(\hat{\theta}_m^*))$  and we also obtain  $\mathbb{P}(Z_4 \leq c_n) = o(\frac{1}{n})$ .

Consider now  $Z_2 := n(\hat{\gamma}_n(\theta_m^*) - R(\theta_m^*))$ . Then,

$$\mathbb{E}[|Z_2|^{8/3}] \leq 2^{5/3} \left( \mathbb{E}[\|\hat{L}_n(\theta) - L_n(\theta)\|_{\Theta}^{8/3}] + n^{8/3} \mathbb{E}\left[\left|\sum_{k=1}^n (\gamma(\theta_m^*, X_k) - R(\theta_m^*))\right|^{8/3}\right] \right).$$

Using [Bardet and Wintenberger \(2009\)](#), we know that  $\sup_{n \in \mathbb{N}^*} E[\|\hat{L}_n(\theta) - L_n(\theta)\|_{\Theta}^{8/3}] < \infty$  from Assumption A5 and since  $\delta > 7/2 > 2$ . Now, consider  $Y_k := \gamma(\theta_m^*, X_k) - R(\theta_m^*)$ . Then,  $(Y_k)_{k \in \mathbb{Z}}$  is a stationary time series,  $\tau_Y$ -weakly dependent because, using the same type of arguments as in the proof of Lemma 4.3, we have:

$$\tau_Y(s) \leq \sum_{\ell=1}^{\infty} (\alpha_{\ell}(f_{\theta}, \Theta) + \alpha_{\ell}(M_{\theta}, \Theta)) \lambda_{s+1-\ell},$$

with  $\lambda$  defined in Lemma 4.2. Therefore, using Assumption A5, we also have  $\tau_Y(s) = O(s^{\delta-1} \log(s))$ , with  $\delta > 7/2$ . Now, using the same type of arguments as in the proof of Lemma 4.4,

$$\mathbb{E}\left[\left|\sum_{k=1}^n Y_k\right|^{8/3}\right] \leq C_{8/3} (M_{8/3,n} + M_{2,n}^{4/3}),$$

and  $M_{2,n} \leq Cn$  while  $M_{8/3,n} \leq Cn \sum_{i=1}^n i^{8/3-1-\delta} \log(i) = o(n^{4/3})$ . Therefore, there exists  $C > 0$  such that for any  $n \in \mathbb{N}^*$ ,

$$\mathbb{E}\left[\left|\sum_{k=1}^n Y_k\right|^{8/3}\right] \leq Cn^{4/3}.$$

Finally, we deduce that there exists  $C > 0$  such that for any  $n \in \mathbb{N}^*$ ,

$$\mathbb{E}[|Z_2|^{8/3}] \leq Cn^{4/3}. \quad (4.8.30)$$

This result and Markov Inequality imply,

$$\begin{aligned} \mathbb{P}(Z_2 \leq c_n) &\leq \mathbb{P}(|Z_2|^{8/3} \geq (|c_n|)^{8/3}) \leq \mathbb{E}[|Z_2|^{8/3}] \frac{1}{|c_n|^{8/3}} \\ &\implies \mathbb{P}(Z_2 \leq c_n) = O\left(n^{4/3} \frac{1}{|c_n|^{8/3}}\right) = O\left(\frac{1}{n^{4/3}}\right), \end{aligned} \quad (4.8.31)$$

We obtain the same bound for  $Z_3 := n(R(\theta^*) - \hat{\gamma}_n(\theta^*))$ .

Finally using the assumption (1.3.20), we have:

$$\begin{aligned} n \mathbb{P}(Z_5 \leq c_n) &= \mathbb{P}\left((\text{pen}(m) - \text{pen}(m^*)) \leq -\frac{2}{5} DK_L(\theta^* \parallel \theta_m^*)\right) \\ &\leq \mathbb{P}\left(\text{pen}(m^*) \geq \frac{2}{5} DK_L(\theta^* \parallel \theta_m^*)\right) \xrightarrow{n \rightarrow \infty} 0. \end{aligned} \quad (4.8.32)$$

By this way, (4.4.6) is established.  $\square$

*Proof of Theorem 4.3.* The proof is mainly based on Lemma 4.7. From the proof of Lemma 4.1, we deduce that for any  $m \in \mathcal{M}$ , there exists two positive random variable  $Y(m)$  and

$Z(m)$  such as  $n |I_1(m) + I_2(m)| \leq Y(m)$  and  $n |I_3(m)| \leq Z$  for any  $n \in \mathbb{N}^*$ . Moreover,  $Y(m)$  and  $Z$  have bounded expectations. Therefore, using Markov Inequality and since  $\mathcal{M}$  is supposed to be a finite family of models, for any  $\varepsilon > 0$  there exists  $K'_\varepsilon > 0$  such as

$$\limsup_{n \rightarrow \infty} \max_{m \in \mathcal{M}} \mathbb{P}(n \text{pen}_{id}(m) \geq K'_\varepsilon) \leq \varepsilon.$$

Therefore, using this inequality and (4.4.7), we deduce that for any  $\varepsilon > 0$  there exist  $M_\varepsilon > 0$  and  $N_\varepsilon \in \mathbb{N}^*$  such that for any  $n \geq N_\varepsilon$ ,

$$\mathbb{P}\left(n |(\text{pen}(\widehat{m}_{id}) - \text{pen}(\widehat{m}_{\text{pen}})) - (\text{pen}_{id}(\widehat{m}_{id}) - \text{pen}_{id}(\widehat{m}_{\text{pen}}))| \leq M_\varepsilon\right) \geq 1 - \varepsilon. \quad (4.8.33)$$

The proof of (4.4.8) is now completed from (4.8.25) of Lemma 4.7 and (4.8.33).  $\square$

*Proof of Theorem 4.4.* We first verify conditions (C1) and (C2) of [Chen \(1985\)](#) that are sufficient to imply Conditions (i), (ii) and (iii) of [Kass et al. \(1990\)](#). Condition (C1) requires that  $\widehat{\sigma}_n$  the largest eigenvalue of  $(-\partial_{\theta_i \theta_j}^2 \widehat{L}_n(\widehat{\theta}_m))_{i,j \in m}^{-1}$  satisfies  $\widehat{\sigma}_n \xrightarrow[n \rightarrow +\infty]{a.s.} 0$ , which is satisfied since it was already established that  $\frac{1}{n} (\partial_{\theta_i \theta_j}^2 \widehat{L}_n(\widehat{\theta}_m))_{i,j \in m} \xrightarrow[n \rightarrow +\infty]{a.s.} F(\theta_m^*)$  and  $F(\theta_m^*)$  is assumed to be a negative definite matrix. Moreover, condition (C2) is also satisfied because  $\theta_m \in \Theta_m \mapsto (\partial_{\theta_i \theta_j}^2 \widehat{L}_n(\theta_m))_{i,j \in m}$  and  $\theta_m \in \Theta_m \mapsto ((\partial_{\theta_i \theta_j}^2 \widehat{L}_n(\theta_m))_{i,j \in m})^{-1}$  are continuous functions for  $n$  large enough. Therefore, using  $h_n = -\frac{1}{n} \widehat{L}_n$ , the assumptions of Theorem 1 of [Kass et al. \(1990\)](#) are satisfied and this implies that:

$$\begin{aligned} \int_{\Theta_m} b_m(\theta) \exp(\widehat{L}_n(\theta)) d\theta &= \exp(\widehat{L}_n(\widehat{\theta}_m)) (2\pi)^{|m|/2} \\ &\quad \times \det\left(n\left(-\frac{1}{n} \partial_{\theta_i \theta_j}^2 \widehat{L}_n(\widehat{\theta}_m)\right)_{i,j \in m}\right)^{-1/2} (b_m(\widehat{\theta}_m) + O(n^{-1})) \quad a.s. \end{aligned}$$

As a consequence, we have:

$$\begin{aligned} \widehat{S}_n(m, X) &= -\log(|\mathcal{M}|) + \log \left[ \int_{\Theta_m} b_m(\theta) \exp(\widehat{L}_n(\theta)) d\theta \right] \\ &= \widehat{L}_n(\widehat{\theta}_m) - \frac{\log(n)}{2} |m| + \log(b_m(\widehat{\theta}_m)) \\ &\quad + \frac{\log(2\pi)}{2} |m| - \frac{1}{2} \log \left[ \det(-\widehat{F}_n(m)) \right] - \log(|\mathcal{M}|) + O(n^{-1}) \quad a.s. \end{aligned}$$

$\square$

# 5

## Data driven model selection for same-realization predictions in autoregressive processes

---

5.1	INTRODUCTION .....	114
5.2	MODEL SELECTION APPROACH AND PRELIMINARY RESULTS .....	115
5.2.1	Model Selection Approach .....	116
5.2.2	Notations .....	117
5.2.3	Preliminary Results .....	117
5.3	Bias-Variance Result .....	121
5.4	PROOFS .....	121
5.4.1	Proof of Theorem 5.1 .....	121
5.4.2	Proof of Proposition 5.4 .....	128
5.4.3	Proof of Proposition 5.3 .....	130
5.4.4	Technical Lemmas .....	131
5.5	THEORETICAL TOOLS .....	133

---

The content of this chapter is contained in the submitted preprint, Kare Kamila. "Data Driven Model Selection for Same-Realization Predictions in Autoregressive Processes" <https://hal.archives-ouvertes.fr/hal-03169343/document>.

### Abstract

This paper is about the one-step ahead prediction of the future of observations drawn from an infinite-order autoregressive  $AR(\infty)$  process. It aims to design penalties (fully data driven) ensuring that the selected model verifies the efficiency property but in the non asymptotic framework. We show that the excess risk of the selected estimator enjoys the best bias-variance trade-off over the considered collection. To achieve these results, we needed to overcome the dependence difficulties by following a classical approach which consists in restricting to a set where the empirical covariance matrix is equivalent to the

theoretical one. We show that this event happens with probability larger than  $1 - c_0/n^2$  with  $c_0 > 0$ . The proposed data driven criteria are based on the minimization of the penalized criterion akin to the Mallows's  $C_p$ .

**Key words:** Model selection, oracle inequality, efficiency, autoregressive process, data driven.

## 5.1 INTRODUCTION

Consider observations  $(X_1, X_2, \dots, X_n)$  arising from a trajectory of the process

$$X_t = f^*((X_{t-i})_{i \in \mathbb{N}^*}) + \sigma \xi_t \text{ for any } t \in \mathbb{Z}. \quad (5.1.1)$$

where  $(\xi_t)_{t \in \mathbb{Z}}$  is a sequence of zero-mean independent identically distributed random variables (i.i.d.r.v) satisfying  $\mathbb{E}(|\xi_0|^4) < \infty$  and  $f^* : \mathbb{R}^{\mathbb{N}} \rightarrow \mathbb{R}$  is a measurable function and  $\sigma > 0$  an unknown constant.

The problem is to estimate the function  $f^*$  using these observations. The process (5.1.1) is a particular case of the general class of affine causal process studied in [Doukhan and Wintenberger \(2008\)](#) and [Bardet and Wintenberger \(2009\)](#). The study of this type of process more often requires the classical regularity condition on the function  $f^*$ , which are not restrictive at all and remain valid in various time series models. This condition can be stated as follows:

$$\sum_{k=1}^{\infty} \left( \sup_{x \in \mathbb{R}^{\mathbb{N}}} \left| \frac{\partial}{\partial x_k} f^*(x) \right| \right) < 1, \quad (5.1.2)$$

provided that that  $f^*$  admits partial derivatives on  $\mathbb{R}^{\mathbb{N}}$ . Under (5.1.2) and if the noise  $\xi_0$  admits  $r$ -order moments, [Doukhan and Wintenberger \(2008\)](#) showed that there exists a stationary, mixing and ergodic solution to (5.1.1) admitting  $r$ -order moments.

Moreover, [Bardet and Wintenberger \(2009\)](#) studied the consistency and the asymptotic normality of the QMLE of  $\theta^* = (\theta_i^*)_{i \in \mathbb{N}}$  in the case  $f^* = f_{\theta^*}$ .

In this paper, we will focus only on processes with a linear regression function ( $f_{\theta^*}$ ) with respect to the past and depending on some parameter  $\theta^* \in \mathbb{R}^{\mathbb{N}}$ ; that is

$$f^*(X_{t-1}, X_{t-2}, \dots) = f_{\theta^*}(X_{t-1}, X_{t-2}, \dots) = \sum_{i=1}^{\infty} \theta_i^* X_{t-i}. \quad (5.1.3)$$

For such processes, condition (5.1.2) becomes

$$\mathbf{A1} : \quad \sum_{i=1}^{\infty} |\theta_i^*| < 1.$$

Even if this condition reduces the set of parameters a bit, the class of  $\text{AR}(\infty)$  processes checking the condition **A1** is rich and of practical importance because it contains almost all invertible causal  $\text{ARMA}(p, q)$  processes and it is very useful for prediction given the past. Moreover, contrary to the autocovariance of  $\text{ARMA}(p, q)$  processes which decays exponentially fast,  $\text{AR}(\infty)$  are able to model more complex behaviour such as slower decay of the covariance structure.

Henceforth, let observations  $(X_1, X_2, \dots, X_n)$  be a trajectory of the solution  $X := (X_t)_{t \in \mathbb{Z}}$  of (5.1.1) verifying **A1**. The goal of this paper is to predict the next value  $X_{n+1}$ . In fact, if  $\theta^*$  were known, a simple prediction of  $X_{n+1}$  could be  $f_{\theta^*}(X_n, X_{n-1}, \dots)$  setting  $X_t = 0$  for all  $t < 0$ . However,  $\theta^*$  is generally unknown and it is impossible to provide a direct estimator since its coordinate are infinite. It is classical to identify a 'good' finite-dimensional model based on the data which can be done by sieve estimation where only a finite number of  $\{\theta_i^*\}_{i=1}^K$  is estimated and letting  $K$  grows as the sample size increases. A usual approach to this is model selection and the goal is to provide a model with the prediction error as small as the oracle's one.

This question has already been addressed in the literature. [Shibata \(1980\)](#) was the first to tackle this issue. He proved that Akaike criterion is *asymptotically efficient* in the sense that the selected model achieves a smaller one-step mean squared error of prediction when it is fitted to predict an independent realization of the same process. Following Shibata's asymptotically setting, [Ing and Wei \(2003\)](#) and [Ing et al. \(2005\)](#) extended this result for same realization predictions. Indeed, they argued that the Shibata's idea to fit the model to another independent realization is unrealistic since in practice we only have one data at hand. The common feature of these works is their asymptotic framework. Meanwhile, there were several authors which study this question in non asymptotic regime. [Goldenshluger and Zeevi \(2001\)](#) in the non parametric framework, studied how well a Gaussian process admitting an  $\text{AR}(\infty)$  representation can be approximated by a finite-order AR model.

In [Baraud et al. \(2001a\)](#) and [Baraud et al. \(2001b\)](#), they analyzed similar question, but a little bit different as observations arise from an auto-regressive model of order  $k$ . They proved an oracle inequality under several conditions, for instance the compactly supported base of the regression function. Moreover, they assume that the process is  $\beta$ -mixing which is usually admitted, but quite hard to verify in practice. For linear processes, the  $\tau$ -mixing is more suitable since its coefficients can be easily computed (see [Comte et al. \(2008\)](#)) and be bounded by a function of the model parameter  $\theta^*$  (see [Doukhan and Wintenberger \(2008\)](#)). In this work, we do not assume any mixing property of the process since the condition **A1** implies the  $\tau$ -mixing property (see [Doukhan and Wintenberger \(2008\)](#)) and we will see that the decreasing rate of  $\tau$ -mixing coefficients is bounded by the decreasing rate of the coefficients  $\theta^* = (\theta_i^*)_{i \in \mathbb{N}}$ .

Based on the above and following a model selection approach, our purpose in this work is to design adaptive penalties in such a way that the selected model mimic the oracle when observations arise from  $\text{AR}(\infty)$  under mild conditions, including the existence of the all order moment of the noise, the decreasing rate of the coefficients of  $(\theta_i^*)_{i \in \mathbb{N}}$  so that thanks to a result by [Doukhan and Wintenberger \(2008\)](#), the generating process has nice properties such as stationarity,  $\tau$ -mixing.

The main contribution of this paper is to have proved that the excess risk of the selected estimator enjoys the best bias-variance trade-off over the considered collection.

The paper is organized as follows. The model selection approach along with preliminary results are described in Section 5.2. The main results are presented Section 5.3. Finally, Section 5.4 contains the proofs.

## 5.2 MODEL SELECTION APPROACH AND PRELIMINARY RESULTS .....

### 5.2.1 Model Selection Approach

Let  $S_m$  (shortly  $m$ ) a model for  $f^*$  to be the set of linear function  $f$  from  $\mathbb{R}^{D_m}$  to  $\mathbb{R}$  such that

$$f(x_1, x_2, \dots, x_{D_m}) = \sum_{i=1}^{D_m} \theta_i x_i, \quad (5.2.1)$$

with  $\theta = (\theta_1, \dots, \theta_{D_m}) \in \Theta_m$  and  $\Theta_m$  a compact set of  $\mathbb{R}^{D_m}$  satisfying  $\sup_{\theta \in \Theta_m} \sum_{i=1}^{D_m} |\theta_i| < 1$ .

$S_m$  can be viewed as an  $\text{AR}(D_m)$  model.

Given a predictor  $f_\theta \in S_m$ , its quality is measured by the quadratic loss

$$R(\theta) = \mathbb{E}[(X_{n+1} - f_\theta^{n+1})^2]$$

where  $f_\theta^n = f_\theta(X_{n-1}, \dots, X_{n-D_m})$ . The Bayes predictor which minimizes  $R(\theta)$  over the set of all predictors is clearly the inaccessible function  $f_{\theta^*}$ . Let then introduce the excess loss of the predictor  $f_\theta$  (with respect to  $f_{\theta^*}$ )

$$\ell(\theta, \theta^*) := R(\theta) - R(\theta^*) = \mathbb{E}[(f_\theta^{n+1} - f_{\theta^*}^{n+1})^2] \geq 0.$$

Given a model  $m$ , we define its best predictor  $f_{\theta_m^*}$  by

$$\theta_m^* = \underset{\theta \in \Theta_m}{\operatorname{argmin}} R(\theta).$$

Its empirical version minimizing the least-squares contrast is

$$\hat{\theta}_m = \underset{\theta \in \Theta_m}{\operatorname{argmin}} \gamma_n(\theta) \quad \text{where} \quad \gamma_n(\theta) = \frac{1}{n} \sum_{t=1}^n (X_t - f_\theta^t)^2. \quad (5.2.2)$$

In this work, we will consider that the excess loss is measured on the design points, that is to say

$$\ell(\hat{\theta}, \theta^*) = \mathbb{E}[\|F_{\hat{\theta}} - F_{\theta^*}\|_n^2] \quad (5.2.3)$$

where  $F_\theta := (f_\theta^1, \dots, f_\theta^n)^\top$  and  $\|x\|_n^2 = \frac{1}{n} \sum_{t=1}^n x_t^2$ .

Given that all the models which can be considered must have finite dimensions for fixed  $n$ , making all  $S_m$  wrong models, it is classical to let the dimension of competitive models grow with the number of observations. This will help reduce the excess loss and provide a better approximation of  $f_{\theta^*}$ .

Let  $\mathcal{M}_n$  a countable collection of hierarchical model  $S_m$  and  $K_n$  is the dimension of the largest model in  $\mathcal{M}_n$  satisfying  $|\mathcal{M}_n| \leq K_n < n$ . We follow the classical approach of model selection which consists in minimizing the penalized LSE. Let  $\text{pen}: \mathcal{M}_n \rightarrow \mathbb{R}^+$  be a penalty function, possibly data-dependent, and define

$$\hat{m} = \underset{m \in \mathcal{M}_n}{\operatorname{argmin}} \{C(m)\} \quad \text{with} \quad C(m) := \gamma_n(\hat{\theta}_m) + \text{pen}(S_m). \quad (5.2.4)$$

Thus, the best possible choice over  $\mathcal{M}_n$  is  $m^*$  the so-called *oracle* defined as

$$m^* \in \arg \inf_{m \in \mathcal{M}_n} \ell(\hat{\theta}_m, \theta^*). \quad (5.2.5)$$

The oracle  $m^*$  is unachievable since it depends on  $\theta^*$  and the distribution  $P_{(X_1, \dots, X_n)}$  that are unknowns. However, we hope to select a model  $\hat{m}$  so that  $\ell(\hat{\theta}_{\hat{m}}, \theta^*)$  is closest to  $\ell(\hat{\theta}_{m^*}, \theta^*)$ .

The goal of this paper is to propose a data driven penalty in order to obtain an oracle inequality

$$\ell(\hat{\theta}_{\hat{m}}, \theta^*) \leq C_1 \inf_{m \in \mathcal{M}_n} \{\ell(\hat{\theta}_m, \theta^*)\} + \frac{C_2}{n} \quad (5.2.6)$$

with the leading constant  $C_1$  close to one and  $C_2 > 0$ . This goal could rather be to show that the excess risk of the selected estimator  $\hat{\theta}_{\hat{m}}$  realizes the best bias-variance trade-off, which would make our penalty an ideal choice in terms of excess risk.

$$\ell(\hat{\theta}_{\hat{m}}, \theta^*) \leq C'_1 \inf_{m \in \mathcal{M}_n} \left\{ \ell(\theta_m^*, \theta^*) + \text{pen}(S_m) \right\} + \frac{C'_2}{n} \quad (5.2.7)$$

with the leading constant  $C'_1 = 1 + \delta$  with  $\delta > 0$  (and close to 0) and  $C'_2 > 0$ .

That is to say that the selected model  $\hat{m}$  will be large enough to reduce its bias, but not too large to avoid high variance.

## 5.2.2 Notations

We will use the following norms:

- $\|\cdot\|$  denotes the usual Euclidean norm on  $\mathbb{R}^\nu$ , with  $\nu \geq 1$ ;
- $\|A\|_{\text{op}}$  is the operator norm of  $A$  as the square root of the largest eigenvalue of  $A^\top A$ . If  $A$  is symmetric, then  $\|A\|_{\text{op}}$  is the largest (in absolute value) eigenvalue of  $A$ .
- if  $X$  is a  $\mathbb{R}^\nu$ -random variable and  $r \geq 1$ , we set  $\|X\|_r = (\mathbb{E}[\|X\|^r])^{1/r} \in [0, \infty]$ .

## 5.2.3 Preliminary Results

As we are in dependence setting, we are going to leverage the  $\tau$ -mixing property of  $(X_t)_{t \in \mathbb{Z}}$  in order to obtain some exponential Inequalities. The  $\tau$ -mixing coefficients are a measure of the dependence of the process and has been introduced by [Dedecker and Prieur \(2005\)](#). This will help us build 'independents' random vectors and apply classical exponential Inequalities. Let then introduce some notations.

Let  $(\Omega, \mathcal{C}, \mathbb{P})$  be a probability space,  $\mathcal{M}$  a  $\sigma$ -subalgebra of  $\mathcal{C}$  and  $Z$  a random variable with values in a Banach space  $(E, \|\cdot\|_E)$ . Assume that  $\mathbb{E}|Z| < \infty$  and define

$$\tau^{(p)}(\mathcal{M}, Z) = \left\| \sup_{f \in \Lambda(E)} \left\{ \left| \int f(x) \mathbb{P}_{Z|\mathcal{M}}(dx) - \int f(x) \mathbb{P}_Z(dx) \right| \right\} \right\|_p$$

where  $\Lambda(E)$  is the set of 1-Lipschitz function, i.e. the functions  $f$  from  $(E, \|\cdot\|_E)$  to  $\mathbb{R}$  such that  $|f(x) - f(y)| \leq \|x - y\|_E$ .

Using the definition of  $\tau$ , we will measure the dependence of the strictly stationary sequence  $(Z_t)_{t \in \mathbb{Z}}$  thanks to the coefficients defined as follows. For any  $s \geq 0$ , let introduce the norm  $\|x - y\|_{\mathbb{R}^k} = (|x_1 - y_1| + \dots + |x_k - y_k|)$  and setting  $\mathcal{M}_i = \sigma(Z_t, t \leq i)$  and if  $\mathbb{E}(|Z_1|) < \infty$ , let

$$\tau_{Z, \infty}^{(p)}(s) = \sup_{l > 0} \left\{ \max_{1 \leq k \leq l} \frac{1}{k} \sup \left\{ \tau^{(p)}(\mathcal{M}_i, (Z_{i_1}, \dots, Z_{i_k})) \mid i + s \leq i_1 < \dots < i_k \right\} \right\}.$$



Finally, the time series  $(Z_t)_{t \in \mathbb{Z}}$  is  $\tau_{Z,\infty}^{(p)}$ -weakly dependent when its coefficients  $\tau_{Z,\infty}^{(p)}$  tend to 0 as  $s$  tends to infinity.

The next Proposition that is a consequence of Theorem 3.1 in [Doukhan and Wintenberger \(2008\)](#) gives a link between the  $\tau$ -mixing coefficients of the process  $(X_t)_{t \in \mathbb{Z}}$  and the coefficients  $\theta_i^*$  of the model (5.1.3).

**Proposition 5.1.** *Assume **A1** holds and if  $|\theta_t^*| = O(t^{-\gamma})$  with  $\gamma > 1$ , there exists a  $\tau$ -weakly dependent stationary solution of (5.1.1) and a constant  $C_\tau > 0$  such that for  $r > 0$*

$$\tau_{X,\infty}^{(2)}(r) \leq C_\tau \left( \frac{\log r}{r} \right)^{\gamma-1} \quad (5.2.8)$$

*Proof.* With  $G(x, \xi_0) = \sigma \xi_0 + f_{\theta^*}(x)$  for any  $x \in \mathbb{R}^\infty$ , it holds

$$\|G(x, \xi_0) - G(y, \xi_0)\|_2 = |f_{\theta^*}(x) - f_{\theta^*}(y)| \leq \sum_{i=1}^{\infty} |\theta_i^*| |x_i - y_i|.$$

Therefore (5.2.8) is a straightforward application of Theorem 3.1 in [Doukhan and Wintenberger \(2008\)](#).  $\square$

As we are going to need independence for block of random variables, let denote for  $t = 1, \dots, n$  the random vector  $\vec{X}_t := (X_{t-1}, \dots, X_{t-K_n})^\top$ . One can see that the process  $(\vec{X}_t)_{t \in \mathbb{Z}}$  is also mixing with  $\tau_{\vec{X},\infty}^{(1)}$  upper bounded by  $K_n \tau_{X,\infty}^{(1)}$  (see Lemma 5.1).

Now, we construct random variables approximating  $\vec{X}_t$ 's enjoying the independence by block property. Let  $s_n, q_n$  two integers such that  $n = 2 s_n q_n$ . We are going to build  $2 s_n$  blocks of length  $q_n$  so that the even index blocks are independent and so the odd index blocks.

For  $k = 0, \dots, s_n - 1$  let denote by

$$A_k = (\vec{X}_{2kq_n+1}, \dots, \vec{X}_{(2k+1)q_n}) \quad \text{and} \quad B_k = (\vec{X}_{(2k+1)q_n+1}, \dots, \vec{X}_{(2k+2)q_n}).$$

We recall a result of [Lerasle et al. \(2011\)](#) which is a consequence of the coupling in [Dedecker and Priour \(2005\)](#).

**Proposition 5.2.** *Let  $(X_t)_{t \in \mathbb{Z}}$  be the stationary mixing process process obtained in Proposition 5.1. Let also  $s_n, q_n, A_k, B_k$  defined as above for  $k = 0, \dots, s_n - 1$ . There exist random vectors  $A_k^* = (\vec{X}_{2kq_n+1}^*, \dots, \vec{X}_{(2k+1)q_n}^*)$ ,  $B_k^* = (\vec{X}_{(2k+1)q_n+1}^*, \dots, \vec{X}_{(2k+2)q_n}^*)$  such that:*

1. *For  $k = 0, \dots, s_n - 1$ ,  $A_k^*$  has the same law as  $A_k$ , also  $B_k^*$  and  $B_k$ .*
2. *The random vectors  $(A_k^*)_{0 \leq k \leq s_n - 1}$  are independent and so are the vectors  $(B_k^*)_{0 \leq k \leq s_n - 1}$ .*
- 3.

$$\|A_k - A_k^*\|_1 \leq q_n K_n \tau_{X,\infty}^{(1)}(q_n)$$

$$\text{and} \quad \|B_k - B_k^*\|_1 \leq q_n K_n \tau_{X,\infty}^{(1)}(q_n).$$

To prove the oracle inequality, we will assume some constraints on the observations.

**A2**  $X_t$  is sub-Gaussian with variance proxy  $\sigma_0^2 > 0$  i.e.

$$\mathbb{E}[e^{\lambda X_t}] \leq e^{\lambda^2 \sigma_0^2 / 2} \quad \text{for any } \lambda > 0.$$

Condition **A2** implies that the vector  $Z_t^m = (X_{t-1}, \dots, X_{t-D_m})^\top$  which will be prominent in the proofs, is sub-Gaussian with variance proxy  $D_m \sigma_0^2$ . Indeed for any  $v \in \mathbb{R}^{D_m}$  such that  $\|v\| = 1$ ,

$$\begin{aligned} \mathbb{E}[\exp(\lambda v^\top Z_t^m)] &= \mathbb{E}\left[\prod_{i=1}^{D_m} \exp(\lambda v_i (X_{t-i}))\right] \\ &\leq \prod_{i=1}^{D_m} \left\| \exp(\lambda v_i (X_{t-i})) \right\|_{D_m} \\ &= \prod_{i=1}^{D_m} \exp(\lambda^2 D_m \sigma_0^2 v_i^2 / 2) \\ &= e^{\frac{\lambda^2}{2} D_m \sigma_0^2}, \end{aligned}$$

where the Inequality follows from Hölder's Inequality.

The following assumption provides a sufficient condition to ensure the invertibility of both  $\hat{\Sigma}_m := \mathbf{M}_m^\top \mathbf{M}_m$  and  $\Sigma_m := \mathbb{E}[\hat{\Sigma}_m]$  where  $\mathbf{M}_m = [X_{i-1}, \dots, X_{i-D_m}]_{i=1}^n$ .

**A3:** For any  $f_\theta \in S_m$ ,  $\langle \alpha, \partial_\theta f_\theta \rangle = 0$  a.s.  $\implies \alpha = 0$

This condition means that the columns of the matrix  $\mathbf{M}_m$  are linearly independents.

We will also need to bound eigenvalues of the matrices  $\Sigma_m$  for any  $m \in \mathcal{M}_n$ . To do that, we will leverage the relation between the spectral density of the process and these eigenvalues. Let us denote by  $r$ , the covariance function  $r(h) := \mathbb{E}[X_t X_{t+h}]$  for any integer  $h$ . Let also introduce the function  $g : [-\pi, \pi] \rightarrow \mathbb{C}$  such that for any  $\lambda$ ,

$$g(\lambda) = \frac{1}{2\pi} \sum_{h \in \mathbb{Z}} r(h) e^{-ih\lambda},$$

which exists under **A1** with  $|\theta_t^*| = O(t^{-\gamma})$  where  $\gamma \geq 1$ . Therefore,  $r$  is the inverse transform of  $g$  and  $r(h) = \int_{-\pi}^{\pi} e^{ih\lambda} g(\lambda) d\lambda$  for any  $h \in \mathbb{Z}$ . We will assume that

**A4:** There exists a constant  $a > 0$  such that  $\inf_{-\pi \leq \lambda < \pi} g(\lambda) \geq a$ .

This is a very weak assumption, and we are going to give the value of  $a$  for AR( $p$ ) process with  $p \in \mathbb{N}^*$ . Let denote  $\theta^*(z) = 1 - \sum_{j=1}^p \theta_j^* z^j$ , it is well known for such process that

$$g(\lambda) = \frac{\sigma^2}{2\pi |\theta^*(e^{-i\lambda})|^2}.$$

For instance for  $p$  equal to one, and  $X_t = \theta_1^* X_{t-1} + \sigma \xi_t$  with  $|\theta_1^*| < 1$ , it follows

$$\begin{aligned} g(\lambda) &= \frac{\sigma^2}{2\pi |1 - \theta_1^* e^{-i\lambda}|^2} \\ &= \frac{\sigma^2}{2\pi (1 - 2\theta_1^* \cos(\lambda) + (\theta_1^*)^2)}, \end{aligned}$$

and then it is simple to see that

$$a := \frac{\sigma^2}{2\pi (1 + |\theta_1^*|)^2} \leq g(\lambda) \leq \frac{\sigma^2}{2\pi (1 - |\theta_1^*|)^2}.$$

For  $p \geq 1$  and  $X_t = \sum_{j=1}^p \theta_j^* X_{t-j} + \sigma \xi_t$  satisfying  $\sum_{j=1}^p \theta_j^* < 1$  and  $\theta_j^* \geq 0$ , we have

$$\begin{aligned} g(\lambda) &= \frac{\sigma^2}{2\pi |1 - \sum_{j=1}^p \theta_j^* e^{-ij\lambda}|^2} \\ &= \sigma^2 (2\pi)^{-1} \left( 1 + \sum_{j=1}^p (\theta_j^*)^2 - 2 \sum_{j=1}^p \theta_j^* \cos(j\lambda) + 2 \sum_{k=1}^{p-1} \theta_k^* \left\{ \sum_{j=k+1}^p \theta_j^* \cos((j-k)\lambda) \right\} \right)^{-1}. \end{aligned}$$

Thus, using  $-1 \leq \cos(x) \leq 1$  for any real  $x$ , it follows for every  $\lambda$

$$\begin{aligned} \sigma^2 (2\pi)^{-1} \left( 1 + \sum_{j=1}^p (\theta_j^*)^2 + 2 \sum_{j=1}^p \theta_j^* + 2 \sum_{k=1}^{p-1} \theta_k^* \left\{ \sum_{j=k+1}^p \theta_j^* \right\} \right)^{-1} &\leq g(\lambda) \\ &\leq \sigma^2 (2\pi)^{-1} \left( 1 + \sum_{j=1}^p (\theta_j^*)^2 - 2 \sum_{j=1}^p \theta_j^* - 2 \sum_{k=1}^{p-1} \theta_k^* \left\{ \sum_{j=k+1}^p \theta_j^* \right\} \right)^{-1}. \end{aligned}$$

For such AR( $p$ ) process, one can take the constant  $a$  in **A4** to be equal to

$$a = \sigma^2 (2\pi)^{-1} \left( 1 + \sum_{j=1}^p (\theta_j^*)^2 + 2 \sum_{j=1}^p \theta_j^* + 2 \sum_{k=1}^{p-1} \theta_k^* \left\{ \sum_{j=k+1}^p \theta_j^* \right\} \right)^{-1}.$$

We can now state an important intermediate result which provides uniform lower and upper bound on the spectral norm of the matrices  $\Sigma_m$ .

**Proposition 5.3.** *Under **A1** with  $|\theta_t^*| = O(t^{-\gamma})$  where  $\gamma \geq 2$ , we have for any  $m \in \mathcal{M}_n$*

$$\|\Sigma_m\|_{\text{op}} \leq \pi^{-1} \sum_{i=0}^{\infty} |\mathbb{E}[X_0 X_i]| < \infty. \quad (5.2.9)$$

Moreover and under **A3-A4**, it holds

$$\|\Sigma_m^{-1}\|_{\text{op}} \leq 1/a. \quad (5.2.10)$$

Let us introduce extra important notations. Let denote by  $\mu$  the law of the vector  $\vec{X}_t$  and

$$\Omega_n = \left\{ \omega : \left| \frac{\|F_\theta\|_n^2}{\|F_\theta\|_\mu^2} - 1 \right| \leq \frac{1}{2}, \quad \forall F_\theta \in \bigcup_{m, m' \in \mathcal{M}_n} (S_m + S_{m'}) \right\}$$

where  $\|F_\theta\|_\mu^2 := \frac{1}{n} \mathbb{E} \left[ \sum_{t=1}^n (f_\theta^t)^2 \right] = \int (f_\theta^1)^2 d\mu$ . It is common to consider the set  $\Omega_n$  which makes a link between the empirical norm  $\|\cdot\|_n$  and the  $\mathbb{L}_2$  norm (see for instance Baraud et al. (2001b), Hsu et al. (2011), van de Geer (2002), Comte and Genon-Catalot (2020) among others). We will see that in our framework,  $\Omega_n$  holds with large probability. In all of this work, we assume that  $q_n$  was chosen to verify

$$\mathbf{A5} : \quad \left( \frac{\log q_n}{q_n} \right)^{\gamma-1} \leq \frac{A}{n}, \quad (5.2.11)$$

for some constant  $A$  and  $\gamma > 1$ . Also we choose the integer  $s_n$  such that

$$\mathbf{A6} : \quad \frac{s_n}{2} \min \left\{ \left( \frac{1}{2^7 \sigma_0^2 K_n} \right)^2, \frac{1}{2^8 \sigma_0^2 K_n} \right\} \geq 3 \log n, \quad (5.2.12)$$

This means that  $s_n$  is of the form  $s_n = C \log n$  where  $C \geq 6 \max \left\{ (2^7 \sigma_0^2 K_n)^2, 2^8 \sigma_0^2 K_n \right\}$ .

**Proposition 5.4.** *Under assumptions **A1**, **A6** and if  $|\theta_t^*| = O(t^{-\gamma})$  with  $\gamma \geq 8$ , there exists a constant  $C$  such that*

$$\mathbb{P}(\Omega_n^c) \leq \frac{C}{n^3}. \quad (5.2.13)$$

### 5.3 Bias-Variance Result

We are now able to state the main result of the paper.

**Theorem 5.1.** *Let consider observations  $(X_1, \dots, X_n)$  arising from a solution of the process (5.1.1) satisfying **A1** with  $|\theta_t^*| = O(t^{-\gamma})$  where  $\gamma \geq 8$  and also verifying **A2** and **A4**. Let  $\mathcal{M}_n$  be some countable family of AR models satisfying **A3** and **A5-A6**. For any constant  $x > 0$ , let a penalty function  $\text{pen} : \mathcal{M}_n \rightarrow \mathbb{R}^+$  such that*

$$\text{pen}(S_m) \geq 8x^3 \sigma^2 \frac{D_m}{n}. \quad (5.3.1)$$

*Then, the LSE  $\hat{\theta}_{\hat{m}}$  with  $\hat{m}$  given in (5.2.4), satisfies*

$$\mathbb{E} \left[ \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta^*}\|_n^2 \mathbb{I}_{\Omega_n} \right] \leq C_1(x) \inf_{m \in \mathcal{M}_n} \left\{ \mathbb{E} \left[ \|F_{\theta^*} - F_{\theta_m^*}\|_n^2 \right] + 2 \text{pen}(S_m) \right\} + \frac{x(x+2)}{x-2} \frac{C_2}{n} \quad (5.3.2)$$

where  $C_1(x) = \left( \frac{x+2}{x-2} \right)^2 > 1$  and  $C_2 > 0$ .

As we can see, this result is almost similar to that of Baraud et al. (2001b) obtained in non parametric framework. However, their result is only valid if we want to estimate the function  $F_{\theta^*}$  on some compact set. This restriction is lifted in our parametric framework.

### 5.4 PROOFS

#### 5.4.1 Proof of Theorem 5.1

*Proof.* We follow the scheme of the proof of Baraud et al. (2001b). Let fix  $m \in \mathcal{M}_n$ . From the definition (5.2.4), we have

$$\gamma_n(\hat{\theta}_{\hat{m}}) + \text{pen}(S_{\hat{m}}) \leq \gamma_n(\hat{\theta}_m) + \text{pen}(S_m) \leq \gamma_n(\theta_m^*) + \text{pen}(S_m). \quad (5.4.1)$$

Since,

$$\gamma_n(\hat{\theta}_m) = \frac{1}{n} \sum_{t=1}^n (X_t - f_{\hat{\theta}_m}^t)^2 = \frac{1}{n} \sum_{t=1}^n (f_{\theta^*}^t - f_{\hat{\theta}_m}^t)^2 + \frac{\sigma^2}{n} \sum_{t=1}^n \xi_t^2 - \frac{2\sigma}{n} \sum_{t=1}^n \xi_t (f_{\hat{\theta}_m}^t - f_{\theta^*}^t),$$

(5.4.1) yields to

$$\begin{aligned} \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta^*}\|_n^2 &\leq \|F_{\theta_m^*} - F_{\theta^*}\|_n^2 \\ &\quad + \frac{2\sigma}{n} \sum_{t=1}^n \xi_t (f_{\hat{\theta}_{\hat{m}}}^t - f_{\theta_m^*}^t) + \text{pen}(S_m) - \text{pen}(S_{\hat{m}}). \end{aligned} \quad (5.4.2)$$

The difficult part of this proof is to handle the inner product  $\frac{2}{n} \sum_{t=1}^n \sigma \xi_t (f_{\hat{\theta}_{\hat{m}}}^t - f_{\theta_m^*}^t)$ , which can be rewritten as

$$\begin{aligned} \frac{2}{n} \sum_{t=1}^n \sigma \xi_t (f_{\hat{\theta}_{\hat{m}}}^t - f_{\theta_m^*}^t) &= \frac{2}{n} \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta_m^*}\|_{\mu} \sum_{t=1}^n \sigma \xi_t \frac{(f_{\hat{\theta}_{\hat{m}}}^t - f_{\theta_m^*}^t)}{\|F_{\hat{\theta}_{\hat{m}}} - F_{\theta_m^*}\|_{\mu}} \\ &\leq \frac{2}{n} \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta_m^*}\|_{\mu} \sup_{g \in B(\hat{m}, \mu)} \sum_{t=1}^n \sigma \xi_t g_{\theta}^t \\ &\leq x^{-1} \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta_m^*}\|_{\mu}^2 + n^{-2} x \left( \sup_{g \in B(\hat{m}, \mu)} \sum_{t=1}^n \sigma \xi_t g_{\theta}^t \right)^2 \end{aligned}$$

since  $2ab \leq x^{-1}a^2 + xb^2$  for any  $x > 0$  and where

$$B(m', \mu) = \left\{ F_{\theta} \in S_m + S_{m'} : \|F_{\theta}\|_{\mu}^2 \leq 1 \right\}.$$

Moreover, on the set  $\Omega_n$ , it holds

$$\begin{aligned} \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta_m^*}\|_{\mu}^2 &\leq 2 \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta_m^*}\|_n^2 \\ &\leq 2 \left( \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta^*}\|_n + \|F_{\theta^*} - F_{\theta_m^*}\|_n \right)^2 \\ &\leq 2(1+y) \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta^*}\|_n^2 + 2(1+y^{-1}) \|F_{\theta^*} - F_{\theta_m^*}\|_n^2 \end{aligned}$$

for some  $y > 0$ . As a result,

$$\begin{aligned} \frac{2}{n} \sum_{t=1}^n \sigma \xi_t (f_{\hat{\theta}_{\hat{m}}}^t - f_{\theta_m^*}^t) &\leq 2 \frac{(1+y)}{x} \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta^*}\|_n^2 + 2 \frac{(1+y^{-1})}{x} \|F_{\theta^*} - f_{\theta_m^*}\|_n^2 \\ &\quad + \frac{x}{n^2} \left( \sup_{g \in B(\hat{m}, \mu)} \sum_{t=1}^n \sigma \xi_t g_{\theta}^t \right)^2. \end{aligned}$$

Therefore, from (5.4.2), it holds on  $\Omega_n$

$$\begin{aligned} \left( 1 - 2 \frac{(1+y)}{x} \right) \|F_{\hat{\theta}_{\hat{m}}} - F_{\theta^*}\|_n^2 &\leq \left( 1 + 2 \frac{(1+y^{-1})}{x} \right) \|F_{\theta^*} - F_{\theta_m^*}\|_n^2 \\ &\quad + \text{pen}(S_m) - \text{pen}(S_{\hat{m}}) + \frac{x\sigma^2}{n^2} \left( \sup_{g \in B(\hat{m}, \mu)} \sum_{t=1}^n \xi_t g_{\theta}^t \right)^2 \\ &\leq \left( 1 + 2 \frac{(1+y^{-1})}{x} \right) \|F_{\theta^*} - F_{\theta_m^*}\|_n^2 + \text{pen}(S_m) - \text{pen}(S_{\hat{m}}) \\ &\quad + 8x^3\sigma^2 \frac{D_m + D_{\hat{m}}}{n} + \frac{x\sigma^2}{n^2} \left[ \left( \sup_{g \in B(\hat{m}, \mu)} \sum_{t=1}^n \xi_t g_{\theta}^t \right)^2 - 8nx^2D(S_{\hat{m}}) \right]_+ \end{aligned}$$

where  $D(S_{\widehat{m}}) = \dim(S_m + S_{\widehat{m}}) \leq D_m + D_{\widehat{m}}$ . Hence using the condition on the penalty (5.3.1),

$$\left(1 - 2 \frac{(1+y)}{x}\right) \|F_{\widehat{\theta}_{\widehat{m}}} - F_{\theta^*}\|_n^2 \leq \left(1 + 2 \frac{(1+y^{-1})}{x}\right) \|F_{\theta^*} - F_{\theta_m^*}\|_n^2 + 2 \text{pen}(S_m) + x \sigma^2 V_{\widehat{m}} \quad (5.4.3)$$

with

$$V_{m'} = \left[ \left( \sup_{g_\theta \in B(m', \mu)} \nu_n(g_\theta) \right)^2 - 8 \frac{x^2}{n} D(S_{m'}) \right]_+,$$

where  $\nu_n(g_\theta) := n^{-1} \sum_{t=1}^n \xi_t g_\theta^t$ .

The proof will be established after controlling the expectation of  $V_{\widehat{m}}$  which involves the supremum of an empirical process.

Now we leverage the mixing property in order to apply Talagrand's Inequality (Theorem 5.2) to tackle  $\mathbb{E}[V_{\widehat{m}}]$ .

We have

$$V_{\widehat{m}} = \left[ \sup_{g_\theta \in B(\widehat{m}, \mu)} (\nu_n(g_\theta))^2 - 8 \frac{x^2}{n} D(S_{\widehat{m}}) \right]_+ \leq 2 \sup_{g_\theta \in B(\widehat{m}, \mu)} (\nu_n(g_\theta) - \nu_n^*(g_\theta))^2 + V_{\widehat{m}}^* \quad (5.4.4)$$

where

$$V_{\widehat{m}}^* = \left[ 2 \sup_{g_\theta \in B(\widehat{m}, \mu)} (\nu_n^*(g_\theta))^2 - 8 \frac{x^2}{n} D(S_{\widehat{m}}) \right]_+ \quad \text{and} \quad \nu_n^*(g_\theta) = \frac{1}{n} \sum_{t=1}^n \xi_t g_\theta(\vec{X}_t^*).$$

1/ **Control of  $\mathbb{E} \left[ \sup_{g_\theta \in B(\widehat{m}, \mu)} (\nu_n(g_\theta) - \nu_n^*(g_\theta))^2 \right]$ .** Let  $m' \in \mathcal{M}_n$  and  $g_\theta \in B(m', \mu)$ . Since the parameter set are compacts and  $\theta \mapsto g_\theta$  is continuous, there exists  $\theta_0 \in \Theta_m \cup \Theta_{m'}$  such that

$$\sup_{g_\theta \in B(m', \mu)} (\nu_n(g_\theta) - \nu_n^*(g_\theta))^2 = \frac{1}{n^2} \left( \sum_{t=1}^n \xi_t (g_{\theta_0}(\vec{X}_t) - g_{\theta_0}(\vec{X}_t^*)) \right)^2.$$

As  $\xi_t$  and  $\mathcal{F}_t$  are independents, it follows that

$$\begin{aligned} \mathbb{E} \left[ \sup_{g_\theta \in B(m', \mu)} (\nu_n(g_\theta) - \nu_n^*(g_\theta))^2 \right] &= \frac{1}{n^2} \sum_{t=1}^n \mathbb{E} \left[ \xi_t^2 (g_{\theta_0}(\vec{X}_t) - g_{\theta_0}(\vec{X}_t^*))^2 \right] \\ &= \frac{1}{n} \mathbb{E} \left[ (g_{\theta_0}(\vec{X}_0) - g_{\theta_0}(\vec{X}_0^*))^2 \right] \end{aligned}$$

since  $\mathbb{E}[\xi_0^2] = 1$ . In addition,

$$\begin{aligned} \mathbb{E} \left[ (g_{\theta_0}(\vec{X}_0) - g_{\theta_0}(\vec{X}_0^*))^2 \right] &= \sum_{i=1}^{D(S_{m'})} \sum_{j=1}^{D(S_{m'})} \theta_{0,i} \theta_{0,j} \mathbb{E}[(X_{-i} - X_{-i}^*)(X_{-j} - X_{-j}^*)] \\ &\leq \left( \sum_{i=1}^{D(S_{m'})} \theta_{0,i} \|X_{-i} - X_{-i}^*\|_2 \right)^2 \end{aligned}$$

using Cauchy-Schwarz Inequality. It then follows as  $\sum_{i=1}^{D(S_{m'})} |\theta_{0,i}| < 1$

$$\begin{aligned} \mathbb{E} \left[ \sup_{g_\theta \in B(m', \mu)} \left( \nu_n(g_\theta) - \nu_n^*(g_\theta) \right)^2 \right] &\leq \frac{1}{n} (\tau^{(2)}(q_n))^2 \\ &\leq \frac{C_\tau^2}{n} \left( \frac{\log q_n}{q_n} \right)^{2\gamma-2} \end{aligned}$$

where the last inequality follows from Proposition 5.1. Thus,

$$\begin{aligned} \mathbb{E} \left[ \sup_{g_\theta \in B(\hat{m}, \mu)} \left( \nu_n(g_\theta) - \nu_n^*(g_\theta) \right)^2 \right] &\leq \sum_{m' \in \mathcal{M}_n} \mathbb{E} \left[ \sup_{g_\theta \in B(m', \mu)} \left( \nu_n(g_\theta) - \nu_n^*(g_\theta) \right)^2 \right] \\ &\leq K_n \frac{C_\tau^2}{n} \left( \frac{\log q_n}{q_n} \right)^{2\gamma-2} \\ &\leq \frac{A^2 C_\tau^2}{n^2}, \end{aligned}$$

using Assumption **A5** and since  $K_n \leq n$ .

2/ **Control of  $\mathbb{E}[V_{\hat{m}}^*]$ .**

First, let us rewrite  $\nu_n^*(g_\theta)$  for  $g_\theta \in B(m', \mu)$ . Setting  $\vec{X}_t = (X_{t-1}^*, \dots, X_{t-D(S_{m'})}^*)^\top$ , we have

$$\begin{aligned} \nu_n^*(g_\theta) &= \frac{1}{n} \sum_{t=1}^n \xi_t g_\theta(\vec{X}_t^*) \\ &= \frac{1}{2 s_n q_n} \sum_{k=0}^{s_n-1} \left( \sum_{i=1}^{q_n} \xi_{2kq_n+i} g_\theta(\vec{X}_{2kq_n+i}^*) + \sum_{i=1}^{q_n} \xi_{(2k+1)q_n+i} g_\theta(\vec{X}_{(2k+1)q_n+i}^*) \right) \\ &= \nu_{n,1}^*(g_\theta) + \nu_{n,2}^*(g_\theta) \end{aligned}$$

with

$$\nu_{n,1}^*(g_\theta) = \frac{1}{s_n} \sum_{k=0}^{s_n-1} \nu_{n,1,k}^*(g_\theta) \quad \text{and} \quad \nu_{n,2}^*(g_\theta) = \frac{1}{s_n} \sum_{k=0}^{s_n-1} \nu_{n,2,k}^*(g_\theta)$$

where

$$\nu_{n,1,k}^*(g_\theta) = \frac{1}{2 q_n} \sum_{i=1}^{q_n} \xi_{2kq_n+i} g_\theta(\vec{X}_{2kq_n+i}^*) \quad \text{and} \quad \nu_{n,2,k}^*(g_\theta) = \frac{1}{2 q_n} \sum_{i=1}^{q_n} \xi_{(2k+1)q_n+i} g_\theta(\vec{X}_{(2k+1)q_n+i}^*)$$

Now let remark that  $\nu_{n,1}^*(g_\theta)$  and  $\nu_{n,2}^*(g_\theta)$  are both sum of  $s_n$  independent random variables by virtue of Proposition 5.2. Hence,

$$\begin{aligned} V_{\hat{m}}^* &\leq \left( \sup_{g_\theta \in B(\hat{m}, \mu)} 4 (\nu_{n,1}^*(g_\theta))^2 - 4 x^2 n^{-1} D(S_{\hat{m}}) \right)_+ \\ &\quad + \left( \sup_{g_\theta \in B(\hat{m}, \mu)} 4 (\nu_{n,2}^*(g_\theta))^2 - 4 x^2 n^{-1} D(S_{\hat{m}}) \right)_+. \end{aligned}$$

As a consequence it is sufficient to study  $\mathbb{E}_1^* := \mathbb{E} \left( \sup_{g \in B(\hat{m}, \mu)} 4 (\nu_{n,1}^*(g_\theta))^2 - 4 x^2 n^{-1} D(S_{\hat{m}}) \right)_+$

and the bound for  $\mathbb{E} \left( \sup_{g_\theta \in B(\hat{m}, \mu)} 4 (\nu_{n,2}^*(g_\theta))^2 - 4 x^2 n^{-1} D(S_{\hat{m}}) \right)_+$  will follow by using analogous arguments.

### Bounding $\mathbb{E}_1^*$

Since the noise  $(\xi_t)$  is not bounded, the process  $\nu_{n,1}^*$  is not bounded either. Let's use the technique used in [Comte and Genon-Catalot \(2020\)](#) to overcome this difficulty. Therefore, we decompose  $\xi_t$  as

$$\xi_t = \eta_t + \epsilon_t, \quad \eta_t = \xi_t \mathbb{I}_{|\xi_t| \leq k_n},$$

where  $k_n$  is a deterministic sequence or a constant to be chosen later. We then have

$$\nu_{n,1}^*(g_\theta) = v_{n,1}^*(g_\theta) + v_{n,2}^*(g_\theta), \quad \text{where}$$

$$v_{n,1}^*(g_\theta) = \frac{1}{s_n} \sum_{k=0}^{s_n-1} v_{n,1,k}^*(g_\theta) \quad \text{with} \quad v_{n,1,k}^*(g_\theta) = \frac{1}{2q_n} \sum_{i=1}^{q_n} \eta_{2kq_n+i} g_\theta(\vec{X}_{2kq_n+i}^*) \quad \text{and}$$

$$v_{n,2}^*(g_\theta) = \frac{1}{s_n} \sum_{k=0}^{s_n-1} v_{n,2,k}^*(g_\theta) \quad \text{with} \quad v_{n,2,k}^*(g_\theta) = \frac{1}{2q_n} \sum_{i=1}^{q_n} \epsilon_{2kq_n+i} g_\theta(\vec{X}_{2kq_n+i}^*).$$

Thus,

$$\begin{aligned} \mathbb{E}_1^* &\leq 8 \mathbb{E} \left[ \left( \sup_{g_\theta \in B(\hat{m}, \mu)} (v_{n,1}^*(g_\theta))^2 - 0.5 x^2 n^{-1} D(S_{\hat{m}}) \right)_+ \right] + 2 \mathbb{E} \left[ \sup_{g_\theta \in B(\hat{m}, \mu)} (v_{n,2}^*(g_\theta))^2 \right] \\ &\leq 8 \sum_{m' \in \mathcal{M}_n} \mathbb{E} \left[ \left( \sup_{g_\theta \in B(m', \mu)} (v_{n,1}^*(g_\theta))^2 - 0.5 x^2 n^{-1} D(S_{m'}) \right)_+ \right] \end{aligned} \quad (5.4.5)$$

$$+ 2 \mathbb{E} \left[ \sup_{g_\theta \in B(\hat{m}, \mu)} (v_{n,2}^*(g_\theta))^2 \right]. \quad (5.4.6)$$

We start by bounding the term in (5.4.5). Let  $m' \in \mathcal{M}_n$ . In order to apply Theorem 5.2, one has to find  $M, H$  and  $v$  such that

$$\begin{aligned} \sup_{g_\theta \in B(m', \mu)} |v_{n,1,k}^*(g_\theta)| &\leq M, \quad \mathbb{E} \left[ \sup_{g_\theta \in B(m', \mu)} |v_{n,1}(g_\theta)|^2 \right] \leq H^2, \\ &\text{and} \quad \sup_{g \in B(m', \mu)} \text{Var} (v_{n,1,k}^*(g_\theta)) \leq v. \end{aligned}$$

• Since the noise is bounded here and from the assumption **A1**, the process  $(X_t)$  is also bounded. Indeed, under **A1**, there exists  $(\phi_i^*)$  such that

$$X_t = \sum_{i=0}^{\infty} \phi_i^* \xi_{t-i} \quad \text{with} \quad \sum_{i=0}^{\infty} |\phi_i^*| < +\infty.$$

Therefore  $|X_t| \leq \Phi_0 k_n$  with  $\Phi_0 := \sum_{i=0}^{\infty} |\phi_i^*|$ . Moreover, for any  $g_\theta \in B(m', \mu)$ , we have

$$|g_\theta(\vec{X}_t)| = \left| \sum_{i=1}^{D(S_{m'})} \theta_i X_{t-i} \right| \leq \Phi_0 k_n \sum_{i=1}^{D(S_{m'})} |\theta_i| < \Phi_0 k_n.$$

As a result, we have

$$\begin{aligned} \sup_{g_\theta \in B(m', \mu)} |v_{n,1,k}^*(g_\theta)| &\leq \frac{1}{2q_n} \sup_{g_\theta \in B(m', \mu)} \sum_{i=1}^{q_n} |\eta_{2kq_n+i} g_\theta(\vec{X}_{2kq_n+i}^*)| \\ &\leq \frac{\Phi_0 k_n^2}{2} := M \end{aligned}$$



- Next, since the parameter set are compacts, there exists  $\theta_0 \in \Theta_m \cup \Theta_{m'}$  such that

$$\sup_{g_\theta \in B(m', \mu)} |v_{n,1}^*(g_\theta)|^2 = |v_{n,1}^*(g_{\theta_0})|^2.$$

Moreover,

$$\begin{aligned} \mathbb{E}[|v_{n,1}^*(g_{\theta_0})|^2] &= \frac{1}{s_n} \mathbb{E}[|v_{n,1,0}^*(g_{\theta_0})|^2] \\ &= \frac{1}{4 s_n q_n^2} \sum_{i,j=1}^{q_n} \mathbb{E}[\eta_i g_\theta(\vec{X}_i^*) \eta_j g_\theta(\vec{X}_j^*)] \\ &= \frac{1}{4 s_n q_n^2} \sum_{i=1}^{q_n} \mathbb{E}[(\eta_i g_\theta(\vec{X}_i^*))^2] \\ &\leq \frac{\Phi_0^2 k_n^2}{2 n} \leq \frac{\Phi_0^2 k_n^2}{2 n} D(S_{m'}) := H^2 \end{aligned}$$

since  $D(S_{m'}) \geq 1$ .

- Lastly, as  $\text{Var}[X] \leq E[X^2]$ , it follows from the previous series of equations

$$\text{Var}(v_{n,1,0}^*(g_\theta)) \leq \mathbb{E}[|v_{n,1,0}^*(g_{\theta_0})|^2] \leq \frac{\Phi_0^2 k_n^2}{4 q_n} := v.$$

As a consequence from Theorem 5.2 and taking  $\alpha = \frac{1}{2}(\frac{x^2}{2\Phi_0^2 k_n^2} - 1) > 0$ , we have

$$\begin{aligned} &\mathbb{E}\left[\left(\sup_{g_\theta \in B(m', \mu)} (v_{n,1}^*(g_\theta))^2 - 0.5 x^2 n^{-1} D(S_{m'})\right)_+\right] \\ &\leq \frac{2}{K} \left( \frac{\Phi_0^2 k_n^2}{4 q_n} e^{-K q_n D(S_{m'}) (\frac{x^2}{2\Phi_0^2 k_n^2} - 1)} + \frac{49 \Phi_0^2 k_n^4}{4 n^2 K C^2(\alpha)} e^{-2\sqrt{2} K C(\alpha) \frac{\sqrt{n} \sqrt{D(S_{m'})}}{k_n}} \right). \end{aligned}$$

Hence there exists a constant  $K'$  such that

$$\sum_{m' \in \mathcal{M}_n} \mathbb{E}\left[\left(\sup_{g_\theta \in B(m', \mu)} (v_{n,1}^*(g_\theta))^2 - 0.5 x^2 n^{-1} D(S_{m'})\right)_+\right] \leq \frac{K'}{n}. \quad (5.4.7)$$

- Now, let us upper bound the term in (5.4.6). For any  $m' \in \mathcal{M}_n$  and any  $g_\theta \in B(m', \mu)$ , we have

$$g_\theta(\vec{X}_t) = \sum_{i=1}^{D(S_{m'})} \theta_i X_{t-i} \leq \sup_{t-D(S_{m'}) \leq i < t} |X_i| \left( \sum_{i=1}^{D(S_{m'})} |\theta_i| \right) < \sup_{t-D(S_{m'}) \leq i < t} |X_i|.$$

Therefore,

$$\begin{aligned} v_{n,2,k}^*(g_\theta) &= \frac{1}{2 q_n} \sum_{i=1}^{q_n} \xi_{2kq_n+i} g_\theta(\vec{X}_{2kq_n+i}^*) \\ &< \frac{1}{2 q_n} \sum_{i=1}^{q_n} |\xi_{2kq_n+i}| \sup_{2kq_n+i-K_n \leq t < 2kq_n+i} |X_t^*| := Y_k^*, \end{aligned} \quad (5.4.8)$$

so that

$$\sup_{g \in B(\hat{m}, \mu)} v_{n,2}^*(g\theta) < \frac{1}{s_n} \sum_{k=0}^{s_n-1} Y_k^*.$$

Let us notice that  $(Y_k^*)_k$  is a family of independent random variables as  $(v_{n,1,k}^*(g))_k$ . Thus, it follows

$$\begin{aligned} \mathbb{E} \left[ \sup_{g\theta \in B(\hat{m}, \mu)} |v_{n,2}^*(g\theta)|^2 \right] &< \frac{1}{s_n^2} \sum_{i,j=0}^{s_n-1} \mathbb{E}[Y_i^* Y_j^*] \\ &< \frac{1}{s_n^2} \sum_{i=0}^{s_n-1} \mathbb{E}[Y_i^{*2}] \\ &= \frac{1}{s_n} \mathbb{E}[Y_0^{*2}]. \end{aligned}$$

Moreover,

$$\begin{aligned} \mathbb{E}[Y_0^{*2}] &= \frac{1}{4q_n^2} \sum_{i,j=1}^{q_n} \mathbb{E} \left[ |\xi_i| \sup_{i-K_n \leq t < i} |X_t^*| |\xi_j| \sup_{j-K_n \leq t < j} |X_t^*| \right] \\ &= \frac{1}{4q_n^2} \sum_{i=1}^{q_n} \mathbb{E} \left[ \left( \sup_{i-K_n \leq t < i} |X_t^*| \right)^2 \right] \\ &= \frac{\mu_2}{4q_n}, \end{aligned} \tag{5.4.9}$$

where  $\mu_2 = \mathbb{E}[X_t^2] < \infty$ . It follows

$$\mathbb{E} \left[ \sup_{g\theta \in B(\hat{m}, \mu)} |\nu_{n,1}^*(g\theta)|^2 \right] < \frac{\mu_2}{4s_n q_n} = \frac{\mu_2}{2n}. \tag{5.4.10}$$

Inequality (5.4.7) along with (5.4.10) yields to

$$\mathbb{E}_1^* \leq \frac{8K'}{n} + \frac{\mu_2}{n}.$$

We conclude that there exists  $K > 0$

$$\mathbb{E}[V_{\hat{m}}] \leq \frac{K}{n}. \tag{5.4.11}$$

Returning to (5.4.3), and taking expectation on both sides, it then follows

$$\left(1 - 2 \frac{(1+y)}{x}\right) \mathbb{E} \left[ \|f_{\hat{\theta}_{\hat{m}}}^t - f_{\theta^*}^t\|_n^2 \right] \leq \left(1 + 2 \frac{(1+y^{-1})}{x}\right) \mathbb{E} \left[ \|f_{\theta^*}^t - f_{\theta_m^*}^t\|_n^2 \right] + 2 \text{pen}(S_m) + x \frac{K}{n}. \tag{5.4.12}$$

For  $y = \frac{x-2}{x+2} > 0$ , so that  $1+y = \frac{2x}{x+2}$  and  $1+y^{-1} = \frac{2x}{x-2}$ , we obtain

$$\mathbb{E} \left[ \|f_{\hat{\theta}_{\hat{m}}}^t - f_{\theta^*}^t\|_n^2 \right] \leq C(x) \left( \mathbb{E} \left[ \|f_{\theta^*}^t - f_{\theta_m^*}^t\|_n^2 \right] + 2 \text{pen}(S_m) \right) + \frac{x(x+2)}{x-2} \frac{K}{n}$$

with  $C(x) = \frac{(x+2)^2}{(x-2)^2} > 1$ . □

### 5.4.2 Proof of Proposition 5.4

Since the collection  $\mathcal{M}_n$  is hierarchical, we have

$$\begin{aligned}\mathbb{P}(\Omega_n^c) &\leq \sum_{m \in \mathcal{M}_n} \mathbb{P}\left(\exists F_\theta \in S_m : \left| \frac{\|f_\theta\|_n^2}{\|f_\theta\|_\mu^2} - 1 \right| > \frac{1}{2}\right) \\ &\leq \sum_{m \in \mathcal{M}_n} \mathbb{P}(\Omega_m^c)\end{aligned}$$

where

$$\Omega_m = \left\{ \left| \frac{\|F_\theta\|_n^2}{\|F_\theta\|_\mu^2} - 1 \right| \leq \frac{1}{2} \quad \forall F_\theta \in S_m \right\}.$$

Let  $m \in \mathcal{M}_n$ . We have

$$\mathbb{P}(\Omega_m^c) \leq \mathbb{P}\left(\sup_{F_\theta \in S_m} \left| \frac{\|F_\theta\|_n^2}{\|F_\theta\|_\mu^2} - 1 \right| > \frac{1}{2}\right).$$

Moreover,

$$\sup_{F_\theta \in S_m, \|F_\theta\|_\mu^2=1} \left| \frac{\|F_\theta\|_n^2}{\|F_\theta\|_\mu^2} - 1 \right| > \frac{1}{2} \iff \sup_{F_\theta \in S_m, \|F_\theta\|_\mu^2=1} |\nu_n(F_\theta^2)| > \frac{1}{2}$$

with  $\nu_n(F_\theta^2) = n^{-1} \sum_{t=1}^n \left( (f_\theta^t)^2 - \mathbb{E}[(f_\theta^1)^2] \right)$ . Hence,

$$\mathbb{P}\left(\sup_{f \in S_m} \left| \frac{\|F_\theta\|_n^2}{\|F_\theta\|_\mu^2} - 1 \right| > \frac{1}{2}\right) \leq \mathbb{P}\left(\sup_{F_\theta \in S_m} |\nu_n(F_\theta^2)| > \frac{1}{2}\right).$$

For any  $F_\theta \in S_m$ , using the linearity we can write

$$(f_\theta(\vec{X}_t))^2 = \left( \sum_{i=1}^{D_m} \theta_i X_{t-i} \right)^2 = \sum_{i,j=1}^{D_m} \theta_i \theta_j X_{t-i} X_{t-j} = \theta^\top \widehat{\Sigma}_{m,t} \theta$$

where  $\theta = (\theta_1, \dots, \theta_{D_m})^\top$ ,  $\widehat{\Sigma}_{m,t} = Z_t^m (Z_t^m)^\top$  with  $Z_t^m = (X_{t-1}, \dots, X_{t-D_m})^\top$ . So that with  $\Sigma_m = \mathbb{E}[\widehat{\Sigma}_{m,t}]$ , it follows

$$\nu_n^*(F_\theta^2) = \frac{1}{n} \sum_{t=1}^n \theta^\top (\widehat{\Sigma}_t - \Sigma) \theta = \theta^\top (\widehat{\Sigma}_m - \Sigma_m) \theta,$$

where  $\widehat{\Sigma} = n^{-1} \sum_{t=1}^n \widehat{\Sigma}_t$ . As a result,

$$\sup_{F_\theta \in S_m, \|F_\theta\|_\mu^2=1} |\nu_n^*(F_\theta^2)| \leq \|\widehat{\Sigma}_m - \Sigma_m\|_{\text{op}}.$$

Indeed,

$$\begin{aligned}\sup_{F_\theta \in S_m} |\nu_n^*(F_\theta^2)| &= \sup_{\theta: \sum |\theta_i| < 1} \theta^\top (\widehat{\Sigma}_m - \Sigma_m) \theta = \sup_{\theta: \sum |\theta_i| < 1} \|\theta\|^2 \frac{\theta^\top (\widehat{\Sigma}_m - \Sigma_m) \theta}{\|\theta\|^2} \\ &\leq \sup_{\theta: \|\theta\|^2 \leq 1} \frac{\theta^\top (\widehat{\Sigma}_m - \Sigma_m) \theta}{\|\theta\|^2} = \|\widehat{\Sigma}_m - \Sigma_m\|_{\text{op}}\end{aligned}$$

since  $-1 < \theta_i < 1$  ensures that  $\|\theta\|^2 \leq \sum |\theta_i|$ . Hence,

$$\begin{aligned} \mathbb{P}\left(\sup_{F_\theta \in S_m, \|F_\theta\|_\mu^2=1} |\nu_n(F_\theta^2)| > \frac{1}{2}\right) &\leq \mathbb{P}\left(\|\hat{\Sigma}_m - \Sigma_m\|_{\text{op}} > \frac{1}{2}\right) \\ &\leq \mathbb{P}\left(\|\hat{\Sigma}_m - \hat{\Sigma}_m^*\|_{\text{op}} > \frac{1}{4}\right) + \mathbb{P}\left(\|\hat{\Sigma}_m^* - \Sigma_m\|_{\text{op}} > \frac{1}{4}\right) \\ &=: \mathbb{P}_1 + \mathbb{P}_2. \end{aligned}$$

Using Lemma 5.3 with  $u = 1/4$  and by virtue of **A6**, it follows

$$\begin{aligned} \mathbb{P}_2 &\leq 2 \exp\left(-3 \log n\right) \\ &\leq \frac{2}{n^3}. \end{aligned}$$

Now let bound  $\mathbb{P}_1$ . We know that for a  $D_m \times D_m$  matrix  $A$

$$\|A\|_{\text{op}} \leq \|A\|_\infty := \max_{1 \leq i \leq D_m} \sum_{j=1}^{D_m} |A_{ij}|$$

Thus, from Markov's Inequality,

$$\begin{aligned} \mathbb{P}_1 &\leq 4 \mathbb{E}\left[\|\hat{\Sigma}_m - \hat{\Sigma}_m^*\|_{\text{op}}\right] \\ &\leq 4 \mathbb{E}\left[\max_{1 \leq i \leq D_m} \sum_{j=1}^{D_m} |(\hat{\Sigma}_m - \hat{\Sigma}_m^*)_{i,j}|\right] \\ &\leq 4 \sum_{j=1}^{D_m} \mathbb{E}[|(\hat{\Sigma}_m - \hat{\Sigma}_m^*)_{i_0,j}|] \\ &\leq 4 \sum_{j=1}^{D_m} \mathbb{E}[|X_{t-i_0} X_{t-j} - X_{t-i_0}^* X_{t-j}^*|]. \end{aligned}$$

Moreover,  $|X_{t-i} X_{t-j} - X_{t-i}^* X_{t-j}^*| \leq |X_{t-i}| |X_{t-j} - X_{t-j}^*| + |X_{t-j}^*| |X_{t-i} - X_{t-i}^*|$  so that with Cauchy-Schwartz's Inequality,

$$\begin{aligned} \mathbb{E}[|X_{t-i} X_{t-j} - X_{t-i}^* X_{t-j}^*|] &\leq 2 \|X_0\|_2 \|X_{t-1} - X_{t-1}^*\|_2 \\ &\leq 2 \|X_0\|_2 \tau^{(2)}(q_n). \end{aligned}$$

Hence using Proposition 5.1, it follows

$$\begin{aligned} \mathbb{P}_1 &\leq 8 \|X_0\|_2 D_m \tau^{(2)}(q_n) \\ &\leq 8 \|X_0\|_2 D_m C_\tau \left(\frac{\log q_n}{q_n}\right)^{\gamma-1}. \end{aligned}$$

Moreover, since  $\gamma \geq 8$  and from assumption **A5**, one can find some constant  $A'$  such that

$$\left(\frac{\log q_n}{q_n}\right)^{\gamma-1} \leq \frac{A'}{n^4}.$$

As a result, with  $c_0 := 8 \|X_0\|_2 C_\tau A'$ , it holds

$$\mathbb{P}_1 \leq \frac{c_0}{n^3}.$$

As a consequence,

$$\mathbb{P}(\Omega_n^c) \leq \frac{2 + c_0}{n^3}.$$

□

### 5.4.3 Proof of Proposition 5.3

*Proof.* The proof of the will be based on the relation between the spectral density function and the maximum eigenvalues of the variance covariance matrix.

Denote by  $u \in \mathbb{R}^{D_m}$  the normalized eigenvector associated to the largest eigenvalue  $\lambda_{\max}(\Sigma_m)$ . Hence,

$$\begin{aligned} \lambda_{\max}(\Sigma_m) &= u^\top \Sigma_m u = \sum_{j,k=1}^{D_m} u_j r(j-k) u_k = \int_{-\pi}^{\pi} g(\lambda) \sum_{j,k=1}^{D_m} u_j e^{i(j-k)\lambda} u_k d\lambda \\ &= \int_{-\pi}^{\pi} g(\lambda) \left| \sum_{j=1}^{D_m} u_j e^{ij\lambda} \right|^2 d\lambda \leq \sup_{-\pi \leq \lambda < \pi} g(\lambda) \int_{-\pi}^{\pi} \left| \sum_{j=1}^{D_m} u_j e^{ij\lambda} \right|^2 d\lambda \\ &\leq \sup_{-\pi \leq \lambda < \pi} g(\lambda), \end{aligned}$$

since, using Parseval identity,  $\int_{-\pi}^{\pi} \left| \sum_{j=1}^{D_m} u_j e^{ij\lambda} \right|^2 d\lambda = \sum_{j=1}^{D_m} u_j^2 = 1$ .

But, from Lemma 5.2 and since  $\gamma \geq 2$ , it follows

$$\begin{aligned} \left| \sup_{-\pi \leq \lambda < \pi} g(\lambda) \right| &\leq \frac{1}{2\pi} \sum_{h \in \mathbb{Z}} |r(h)| \\ &\leq \frac{C}{\pi} \sum_{h=0}^{+\infty} \frac{1}{(h+1)^\gamma} < \infty. \end{aligned}$$

Given that  $\Sigma_m$  is symmetric, it follows

$$\|\Sigma_m\|_{\text{op}} = \lambda_{\max}(\Sigma_m) \leq \frac{C}{\pi} \sum_{h=0}^{+\infty} \frac{1}{(h+1)^\gamma},$$

which concludes the proof of (5.2.9).

Now we end by the proof of (5.2.10). Reasoning as above, and by virtue of **A4**, one can show that

$$\lambda_{\min}(\Sigma_m) \geq \inf_{-\pi \leq \lambda < \pi} g(\lambda) \geq a$$

which yields to

$$\|\Sigma_m^{-1}\|_{\text{op}} = \frac{1}{\lambda_{\min}(\Sigma_m)} \leq \frac{1}{a},$$

so that (5.2.10) is established.

□

#### 5.4.4 Technical Lemmas

**Lemma 5.1.** Assume **A1** holds and  $(X_t)$  the mixing stationary solution of (5.1.1). Then, the process  $(\vec{X}_t)$  is mixing and

$$\tau_{\vec{X},\infty}^{(1)}(r) \leq K_n \tau_{X,\infty}^{(1)}(r-1). \quad (5.4.13)$$

*Proof.* Let set by  $\mathcal{M}_{\vec{X}}^i = \sigma(\vec{X}_t, t \leq i)$  and  $\mathcal{M}_X^i = \sigma(X_t, t \leq i)$  for an integer  $i$ . One would like to bound  $\tau(\mathcal{M}_{\vec{X}}^i, (\vec{X}_{j_1}, \dots, \vec{X}_{j_k}))$  for  $j_k > \dots > j_1 \geq i + r$ .

Let assume that the universe  $\Omega$  is rich enough so that, one can find  $\vec{X}_{j_l}^* = (X_{j_l-1}^*, \dots, X_{j_l-K_n}^*)^\top$  with  $l = 1, \dots, k$  verifying

1.  $(\vec{X}_{j_1}^*, \dots, \vec{X}_{j_k}^*)$  is distributed as  $(\vec{X}_{j_1}, \dots, \vec{X}_{j_k})$  and independent of  $\mathcal{M}_{\vec{X}}^i$ ;
2.  $(X_{j_1-1}^*, \dots, X_{j_k-1}^*)^\top$  is distributed as  $(X_{j_1-1}, \dots, X_{j_k-1})^\top$  and independent of  $\mathcal{M}_X^i$ .

As a result,

$$\begin{aligned} \tau(\mathcal{M}_{\vec{X}}^i, (\vec{X}_{j_1}, \dots, \vec{X}_{j_k})) &\leq \sum_{l=1}^k \|\vec{X}_{j_l} - \vec{X}_{j_l}^*\|_1 = \sum_{l=1}^k \sum_{t=1}^{K_n} \mathbb{E}[|X_{j_l-t} - X_{j_l-t}^*|] \\ &\leq K_n \sum_{l=1}^k \mathbb{E}[|X_{j_l-1} - X_{j_l-1}^*|] \\ &= K_n \left\| (X_{j_1-1}, \dots, X_{j_k-1})^\top - (X_{j_1-1}^*, \dots, X_{j_k-1}^*)^\top \right\|_1 \\ &= K_n \tau(\mathcal{M}_X^i, (X_{j_1-1}, \dots, X_{j_k-1})). \end{aligned}$$

This fact along with the definition of  $\tau_{\vec{X},\infty}^{(1)}(r)$  leads to (5.4.13). □

**Lemma 5.2.** Under **A1** with  $|\theta_t^*| = O(t^{-\gamma})$  where  $\gamma > 1$ , we have

$$r(h) = \mathbb{E}[X_0 X_h] = O((h+1)^{-\gamma})$$

*Proof.* By virtue of **A1**, the process  $(X_t)_t$  is causal; that is there exists  $(\phi_i)_{i \in \mathbb{N}}$  such that  $X_t = \sum_{i=0}^{+\infty} \phi_i \xi_{t-i}$  with  $\sum_{i=0}^{+\infty} |\phi_i| < \infty$ . The sequence  $(\phi_i)_{i \in \mathbb{N}}$  is given by the relation  $\phi(z) = \sum_{i=0}^{+\infty} \phi_i z^i = \frac{1}{\theta(z)}$  with  $\theta(z) = 1 - \sum_{i=0}^{+\infty} \theta_i^* z^i$ . Equating coefficients of  $z_j, j = 0, 1, \dots$ , we find that  $\phi_0 = 1$  and for  $i \geq 1$

$$\phi_i = \sum_{j=1}^i \theta_j^* \phi_{i-j}.$$

This fact allows us to deduce that the sequences  $(\phi_i)_{i \in \mathbb{N}}$  and  $(\theta_i^*)_{i \in \mathbb{N}}$  decay at the same rate. Therefore, since  $|\theta_t^*| = O((t+1)^{-\gamma})$ , there exists  $h_0 \in \mathbb{Z}$  such that for any  $h \geq h_0$ , it holds  $|\phi_t| \leq C(t+1)^{-\gamma}$  for some constant  $C > 0$ . Thus,

$$\begin{aligned} r(h) &= \sum_{j=0}^{\infty} \phi_j \phi_{j+h} \\ &\leq C^2 \sum_{j=0}^{\infty} \frac{1}{(j+1)^\gamma} \frac{1}{(j+h+1)^\gamma} \\ &\leq C^2 (h+1)^{-\gamma} \sum_{j=0}^{\infty} \frac{1}{(j+1)^\gamma} \leq C^2 \frac{\pi^2}{6} (h+1)^{-\gamma}, \end{aligned}$$

where the last inequality follows from the fact that  $\gamma \geq 2$  and that established the Lemma.  $\square$

**Lemma 5.3.** *Under assumptions **A2**, it holds for any model  $m \in \mathcal{M}_n$ , and for all  $u > 0$*

$$\mathbb{P}\left(\|\hat{\Sigma}_m^* - \Sigma_m\|_{op} \geq u\right) \leq 2 \exp\left\{-\frac{s_n}{2} \min\left\{\left(\frac{u}{16 D_m \sigma_0^2}\right)^2, \frac{u}{32 D_m \sigma_0^2}\right\}\right\}$$

*Proof.* One can write for a matrix  $A$

$$\|A\|_{op} = \max_{v: \|v\|=1} |v^\top A v| = |v_0^\top A v_0|.$$

Therefore one can find a vector  $v_0 \in \mathbb{R}^{D_m}$  with  $\|v_0\| = 1$  such that

$$\mathbb{P}\left(\|\hat{\Sigma}_m^* - \Sigma_m\|_{op} \geq u\right) = \mathbb{P}\left(|v_0^\top (\hat{\Sigma}_m^* - \Sigma_m) v_0| \geq u\right).$$

But,

$$\begin{aligned} v_0^\top (\hat{\Sigma}_m^* - \Sigma_m) v_0 &= \frac{1}{n} \sum_{t=1}^n (v_0^\top \hat{\Sigma}_{m,t}^* v_0 - v_0^\top \Sigma_m v_0) \\ &= \frac{1}{n} \sum_{t=1}^n (v_0^\top (Z_t^{*m}) (Z_t^{*m})^\top v_0 - v_0^\top \Sigma_m v_0) \\ &= \frac{1}{n} \sum_{t=1}^n (Y_t^2 - \mathbb{E}[Y_t^2]) \end{aligned}$$

with  $Y_t = v_0^\top Z_t^m = \sum_{i=1}^{D_m} v_0^i X_{t-i}^*$ . From **A2**,  $Y_t$  is  $\text{SG}(D_m \sigma_0^2)$ . Therefore,  $Y_t^2$  is  $\text{SE}(256 D_m^2 \sigma_0^4, 16 D_m \sigma_0^2)$  (where SE stands for Sub-Gaussian and SE for Sub-Exponential).

Moreover, we can write

$$\begin{aligned} v_0^\top (\hat{\Sigma}_m^* - \Sigma_m) v_0 &= \frac{1}{n} \sum_{t=1}^n (Y_t^2 - \mathbb{E}[Y_t^2]) \\ &= \frac{1}{s_n} \sum_{k=0}^{s_n-1} \left( \frac{1}{2q_n} \sum_{i=1}^{q_n} (Y_{2kq_n+i}^2 - \mathbb{E}[Y_1^2]) \right) + \frac{1}{s_n} \sum_{k=0}^{s_n-1} \left( \frac{1}{2q_n} \sum_{i=1}^{q_n} (Y_{(2k+1)q_n+i}^2 - \mathbb{E}[Y_1^2]) \right) \\ &= \mathbf{Y}_1 + \mathbf{Y}_2. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbf{Y}_1 &= \frac{1}{s_n} \sum_{k=0}^{s_n-1} \mathbf{Y}_{1,k} \quad \text{and} \quad \mathbf{Y}_2 = \frac{1}{s_n} \sum_{k=0}^{s_n-1} \mathbf{Y}_{2,k} \quad \text{with} \\ \mathbf{Y}_{1,k} &= \frac{1}{2q_n} \sum_{i=1}^{q_n} (Y_{2kq_n+i}^2 - \mathbb{E}[Y_1^2]) \quad \text{and} \quad \mathbf{Y}_{2,k} = \frac{1}{2q_n} \sum_{i=1}^{q_n} (Y_{(2k+1)q_n+i}^2 - \mathbb{E}[Y_1^2]). \end{aligned}$$

$\{\mathbf{Y}_{1,k}\}$  and  $\{\mathbf{Y}_{2,k}\}$  are independent random vectors by virtue of Proposition 5.2. Now, let us show that  $\mathbf{Y}_{i,k}$  are sub-exponentials. For  $\lambda$  such that  $|\lambda| < \frac{1}{16 D_m \sigma_0^2}$ , and denoting

$w_i = Y_{2kq_n+i}^2 - \mathbb{E}[Y_1^2]$ , we have

$$\begin{aligned}
\mathbb{E}[e^{\lambda \mathbf{Y}_{1,k}}] &= \mathbb{E}\left[\exp\left(\frac{1}{2q_n} \sum_{i=1}^{q_n} \lambda w_i\right)\right] \\
&= \mathbb{E}\left[\prod_{i=1}^{q_n} \exp\left(\frac{\lambda w_i}{2q_n}\right)\right] \\
&= \mathbb{E}\left[\prod_{i=1}^{q_n} \left(\exp\left(\frac{\lambda w_i}{2}\right)\right)^{1/q_n}\right] \\
&\leq \prod_{i=1}^{q_n} \left(\mathbb{E}\left[\exp\left(\frac{\lambda w_i}{2}\right)\right]\right)^{1/q_n} \\
&\leq e^{\frac{\lambda^2}{2} 64 D_m^2 \sigma_0^4},
\end{aligned}$$

where we have used Hölder's Inequality. Hence  $\mathbf{Y}_{1,k}$  is  $\text{SE}(64 D_m^2 \sigma_0^4, 16 D_m \sigma_0^2)$ . As a result, using exponential inequalities for SE random variables, it follows

$$\mathbb{P}(\mathbf{Y}_1 \geq u/2) \leq \exp\left\{-\frac{s_n}{2} \min\left\{\left(\frac{u}{16 D_m \sigma_0^2}\right)^2, \frac{u}{32 D_m \sigma_0^2}\right\}\right\}$$

so that

$$\mathbb{P}\left(|v_0^\top (\hat{\Sigma}_m^* - \Sigma_m) v_0| \geq u/2\right) \leq 2 \exp\left\{-\frac{s_n}{2} \min\left\{\left(\frac{u}{16 D_m \sigma_0^2}\right)^2, \frac{u}{32 D_m \sigma_0^2}\right\}\right\}.$$

□

**Lemma 5.4.** Assume **A3** holds, then  $\hat{\Sigma}_m$  is a.e. invertible. Also,  $\Sigma_m$  is invertible.

*Proof.* We can write  $\hat{\Sigma}_m = \mathbf{M}_m^\top \mathbf{M}_m$  with  $\mathbf{M}_m = [X_{i-1}, \dots, X_{i-D_m}]_{i=1}^n$ . By virtue of **A3**,  $\mathbf{M}_m$  is of full rank which implies the a.e. invertibility of  $\hat{\Sigma}_m$ .

Moreover,  $\Sigma_m = \mathbb{E}[\hat{\Sigma}_m] = \mathbb{E}[Z_0^m (Z_0^m)^\top]$  with  $Z_0^m = (X_{-1}, \dots, X_{-D_m})^\top$ . Let  $\mathbf{u} \in \mathbb{R}^{D_m}$ , it follows  $\mathbf{u}^\top \Sigma_m \mathbf{u} = \mathbb{E}[(Z_0^m)^\top \mathbf{u}]^2 \geq 0$ . Let show that whenever the equality holds ( $\mathbf{u}^\top \Sigma_m = 0$ ),  $\mathbf{u} = 0$ .

Since  $((Z_0^m)^\top \mathbf{u})^2 \geq 0$ , its expectation vanishes if and only if  $(Z_0^m)^\top \mathbf{u} = 0$  a.e. which yields to  $\mathbf{u} = 0$  by **A3**. Hence,  $\Sigma_m$  is positive definite and then invertible. □

## 5.5 THEORETICAL TOOLS

The next Theorem is a Talagrand's Inequality given in Klein et al. (2005).

**Theorem 5.2.** Let  $Y_1, \dots, Y_n$  be independent random variables and let  $\mathcal{F}$  be a countable class of uniformly bounded measurable functions. Then for all  $\alpha > 0$ ,

$$\mathbb{E}\left[\sup_{g \in \mathcal{F}} |\eta_n(g)|^2 - 2(1 + 2\alpha) H^2\right]_+ \leq \frac{2}{K} \left(\frac{v}{n} e^{-K\alpha \frac{nH^2}{v}} + \frac{49M^2}{4Kn^2C^2(\alpha)} e^{-\frac{2\sqrt{2}KC(\alpha)\sqrt{\alpha}}{7\sqrt{2}} \frac{nH}{M}}\right)$$

with  $\eta_n(g) = n^{-1} \sum_{t=1}^n (g(Y_t) - \mathbb{E}[g(Y_t)])$  for any  $g \in \mathcal{F}$ ;

$C(\alpha) = (\sqrt{1 + \alpha} - 1) \wedge 1$ ,  $K = 1/6$

$$\sup_{g \in \mathcal{F}} \|g\|_\infty \leq M, \quad \mathbb{E}\left[\sup_{g \in \mathcal{F}} |\eta_n(g)|\right] \leq H, \quad \sup_{g \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \text{Var}(g(Y_t)) \leq v.$$





# 6

## Conclusion Générale et Perspectives

Dans cette thèse, il était question pour nous de proposer et d'étudier sur les plans théorique et numérique des critères de sélection de modèles pour les séries affines causales.

Dans un premier temps, nous avons montré qu'il était indispensable de faire dépendre la pénalité du critère de sélection de la vitesse de décroissance des coefficients de Lipschitz du processus afin d'obtenir un critère consistant en probabilité. Nous avons également proposé un test d'adéquation du modèle sélectionné basé sur l'autocorrélation du carré des résidus du modèle. Les simulations numériques ont montré des résultats assez satisfaisants.

Deuxièmement, nous avons proposé une généralisation du critère de Hannan et Quinn à la classe des séries affines causales. Cette généralisation fournit une nouvelle constante multiplicative (dans la pénalité) connue pour les modèles classiques (ARMA, GARCH ou APARCH) et data-driven calibrable pour les modèles complexes tels que les ARMA-GARCH. Là également, quelques études de simulation ont attesté de la bonne qualité des critères obtenus.

Dans un troisième temps, nous nous sommes intéressés à la construction des critères asymptotiquement efficaces. Nous avons généralisé le critère AIC d'Akaike en s'appuyant sur la pénalité dite idéale. Le comportement asymptotique de cette pénalité idéale nous a suggéré un terme de pénalité qui vaut exactement  $2 D_m$  comme dans l'AIC pour des modèles assez simples, et pour des modèles complexes, nous avons donné une formule moins explicite mais qui peut être calibrée au moyen des données. Dans cette troisième partie, nous avons également à la suite de Schwartz, dérivé le critère BIC basé sur la maximisation de la probabilité a posteriori de choisir le vrai modèle.

Enfin, nous nous sommes restreints à l'étude non asymptotique d'un processus particulier de la classe des modèles affines causaux. Un estimateur des moindres carrés pénalisé est construit à partir d'un critère de sélection adaptatif et la sélection est opérée parmi une collection de modèles linéaires. Nous avons montré que l'estimateur final est presque aussi performant que le meilleur sur la collection considérée, *i.e.* qu'il réalise, à une constante près, le compromis biais-variance. La pénalité obtenue généralise celle de Mallows

et dépend d'une constante que l'on pourrait estimer avec des algorithmes de calibration adaptative.

Au sortir de cette thèse, il est important de souligner que nos travaux ouvrent plusieurs perspectives de recherche.

1/ A la suite du Chapitre 4, il serait intéressant d'obtenir l'inégalité d'oracle (5.2.6) cette fois ci en considérant que les observations proviennent du modèle  $AR(\infty)$  suivant:

$$X_t = \sum_{k=1}^{+\infty} \theta_k^* X_{t-k} + \sigma \xi_t \quad \text{pour tout } t \in \mathbb{Z}. \quad (6.0.1)$$

où  $(\xi_t)_{t \in \mathbb{Z}}$  est un bruit blanc **faible**. La tâche devient beaucoup plus difficile car l'hypothèse d'indépendance du bruit est fondamentale pour garantir l'existence d'une solution mélangeante, propriété nécessaire pour utiliser des inégalités exponentielles usuelles. L'on pourrait ainsi dans un premier temps, supposer que le processus  $(X_t)_{t \in \mathbb{Z}}$  est mélangeant. Cela pourrait permettre de traiter les GARCH puisque l'on peut toujours les écrire comme des ARMA à bruit blanc faible.

2/ Une extension logique du Chapitre 4 serait de considérer que les observations proviennent d'un modèle fonctionnel autoregressif

$$X_t = f^*(X_{t-1}, X_{t-2}, \dots) + \sigma \xi_t \quad \text{pour tout } t \in \mathbb{Z} \quad (6.0.2)$$

avec bruit blanc fort mais où la fonction  $f^*$  n'est plus nécessairement linéaire mais vérifie l'hypothèse de contraction

$$|f^*(x) - f^*(y)| \leq \sum_{i=1}^{\infty} \alpha_i |x_i - y_i|$$

pour tout  $x, y \in \mathbb{R}^{\infty}$  avec  $\sum_{i=1}^{\infty} \alpha_i < 1$ .

Un modèle candidat  $S_m$  à considérer pourrait être dans ce cas, un  $D_m$  sous-espace linéaire de  $\mathbb{L}^2(\mathbb{R}^{D_m})$ . Si  $\mathcal{M}_n$  désigne la famille de modèles candidats, l'on définit l'estimateur pour tout  $S_m \in \mathcal{M}_n$  comme

$$\hat{f}_m = \operatorname{argmin}_{f \in S_m} \gamma_n(f) \quad \text{avec} \quad \gamma_n(f) = \frac{1}{n} \sum_{t=1}^n (X_t - f(X_{t-1}, X_{t-2}, \dots, X_{t-D_m}))^2.$$

A travers une procédure de pénalisation, l'objectif serait de trouver la bonne fonction de pénalité pen de sorte qu'on ait

$$\ell(\hat{f}_{\hat{m}}, f^*) \leq C_1 \inf_{m \in \mathcal{M}_n} \{\ell(\hat{f}_m, f^*)\} + C_2 n^{-1} \quad (6.0.3)$$

avec

$$\hat{m} = \operatorname{argmin}_{m \in \mathcal{M}_n} \{C(m)\} \quad \text{avec} \quad C(m) := \gamma_n(\hat{\theta}_m) + \operatorname{pen}(S_m), \quad (6.0.4)$$

$\ell(\hat{f}_m, f^*) = \mathbb{E}[(X_t - \hat{f}_m(X_{t-1}, X_{t-2}, \dots, X_{t-D_m}))^2] - \sigma^2$  et la constante  $C_1$  proche de 1.

3/ De plus en plus, l'hypothèse de stationnarité est remise en question et l'intérêt pour des séries non stationnaires surtout localement stationnaires est croissant ces dernières années. [Bardet et al. \(2020a\)](#) ont prouvé des résultats généraux (consistance, TLC) pour des processus localement stationnaires. Il serait intéressant de s'appuyer sur ces résultats pour faire de la sélection de modèles sur ces processus.

## Bibliography

- H. Akaike. Fitting autoregressive models for prediction. *Annals of the institute of Statistical Mathematics*, 21(1):243–247, 1969.
- H. Akaike. Information theory and an extension of the maximum likelihood principle. *Proceedings of the 2nd international symposium on information, Akademiai Kiado, Budapest*, 1973.
- D. Allen. The relationship between variable selection and data augmentation and a method for prediction. *Technometrics*, 16(1):125–127, 1974.
- P. Alquier and O. Wintenberger. Model selection for weakly dependent time series forecasting. *Bernoulli*, 18(3):883–913, 2012.
- O. Arkoun, J.-Y. Brua, and S. Pergamenschikov. Sequential model selection method for nonparametric autoregression. *Sequential Anal.*, 38(4):437–460, 2019.
- S. Arlot. *Rééchantillonnage et Sélection de modèles*. PhD thesis, Université Paris Sud-Paris XI, 2007.
- S. Arlot and P. Massart. Data-driven calibration of penalties for least-squares regression. *Journal of Machine learning research*, 10:245–279, 2009.
- Y. Baraud, F. Comte, and G. Viennet. Model selection for (auto-) regression with dependent data. *ESAIM: Probability and Statistics*, 5:33–49, 2001a.
- Y. Baraud, F. Comte, G. Viennet, et al. Adaptive estimation in autoregression or-mixing regression via model selection. *The Annals of Statistics*, 29(3):839–875, 2001b.
- J.-M. Bardet and O. Wintenberger. Asymptotic normality of the quasi-maximum likelihood estimator for multidimensional causal processes. *The Annals of Statistics*, 37(5B):2730–2759, 2009.
- J.-M. Bardet, W. Kengne, and O. Wintenberger. Detecting multiple change-points in general causal time series using penalized quasi-likelihood. *Electronic journal of statistics*, 6:435–477, 2012.
- J.-M. Bardet, Y. Boularouk, and K. Djballah. Asymptotic behavior of the laplacian quasi-maximum likelihood estimator of affine causal processes. *Electronic journal of statistics*, 11(1):452–479, 2017.
- J.-M. Bardet, P. Doukhan, and O. Wintenberger. Contrast estimation of general locally stationary processes using coupling. *Preprint arXiv:2005.07397*, 2020a.

- J.-M. Bardet, K. Kamila, and W. Kengne. Consistent model selection criteria and goodness-of-fit test for common time series models. *Electronic Journal of Statistics*, 14(1):2009–2052, 2020b.
- J.-P. Baudry, C. Maugis, and B. Michel. Slope heuristics: overview and implementation. *Statistics and Computing*, 22(2):455–470, 2012.
- E. Beale, M. Kendall, and D. Mann. The discarding of variables in multivariate analysis. *Biometrika*, 54(3-4):357–366, 1967.
- I. Berkes, L. Horváth, and P. Kokoszka. GARCH processes: structure and estimation. *Bernoulli*, 9:201–227, 2003.
- L. Birgé and P. Massart. Gaussian model selection. *Journal of the European Mathematical Society*, 3(3):203–268, 2001.
- L. Birgé and P. Massart. Minimal penalties for gaussian model selection. *Probability theory and related fields*, 138(1-2):33–73, 2007a.
- L. Birgé and P. Massart. Minimal penalties for gaussian model selection. *Probability theory and related fields*, 138(1-2):33–73, 2007b.
- H. Bozdogan. Model selection and akaike’s information criterion (aic): The general theory and its analytical extensions. *Psychometrika*, 52(3):345–370, 1987.
- P. Brockwell and R. Davis. *Time Series: Theory and Methods: Theory and Methods*. Springer Science & Business Media, 1991.
- J. Cavanaugh. Unifying the derivations for the akaike and corrected akaike information criteria. *Statistics & Probability Letters*, 33(2):201–208, 1997.
- A. Celisse. *Model selection via cross-validation in density estimation, regression, and change-points detection*. PhD thesis, Université Paris Sud-Paris XI, 2008.
- A. Celisse et al. Optimal cross-validation in density estimation with the  $\ell^2$ -loss. *Annals of Statistics*, 42(5):1879–1910, 2014.
- C.-F. Chen. On asymptotic normality of limiting density functions with bayesian ompliations. *J. R. Statist. Soc. B*, 47:540–546, 1985.
- F. Comte and V. Genon-Catalot. Regression function estimation as a partly inverse problem. *Annals of the Institute of Statistical Mathematics*, 72(4):1023–1054, 2020.
- F. Comte, J. Dedecker, and M.-L. Taupin. Adaptive density deconvolution with dependent inputs. *Mathematical methods of Statistics*, 17(2):87, 2008.
- I. Csizsár and P. Shields. Consistency of the bic order estimator. *Electronic research announcements of the American mathematical society*, 5(17):123–127, 1999.
- I. Csizsár, P. C. Shields, et al. The consistency of the bic markov order estimator. *The Annals of Statistics*, 28(6):1601–1619, 2000.
- J. Dedecker and P. Doukhan. A new covariance inequality and applications. *Stochastic Processes and Applications*, 106:63–80, 2003.

- J. Dedecker and C. Prieur. Coupling for  $\tau$ -Dependent Sequences and Applications. *Journal of Theoretical Probability*, 17:861–885, 2004.
- J. Dedecker and C. Prieur. New dependence coefficients. examples and applications to statistics. *Probability Theory and Related Fields*, 132(2):203–236, 2005.
- J. Dedecker, P. Doukhan, G. Lang, J. R. León, S. Louhichi, and C. Prieur. *Weak dependence: With Examples and Applications*. Lecture Notes in Statistics **190**, Springer-Verlag, New York, 2007.
- J. Ding, V. Tarokh, and Y. Yang. Model selection techniques: An overview. *IEEE Signal Processing Magazine*, 35(6):16–34, 2018.
- Z. Ding, C. Granger, and R. Engle. A long memory property of stock market returns and a new model. *Journal of empirical finance*, 1(1):83–106, 1993.
- P. Doukhan and O. Wintenberger. Weakly dependent chains with infinite memory. *Stochastic Processes and their Applications*, 118(11):1997–2013, 2008.
- P. Duchesne and C. Francq. On diagnostic checking time series models with portmanteau test statistics based on generalized inverses and. In *COMPSTAT 2008*, pages 143–154. Springer, 2008.
- D. F. Findley and C.-Z. Wei. Aic, overfitting principles, and the boundedness of moments of inverse matrices for vector autotregressions and related models. *Journal of Multivariate Analysis*, 83(2):415–450, 2002.
- C. Francq and J.-M. Zakoïan. Maximum likelihood estimation of pure garch and arma-garch processes. *Bernoulli*, 10:605–637, 2004.
- C. Francq and J.-M. Zakoïan. *GARCH models: structure, statistical inference and financial applications*. John Wiley & Sons, 2010.
- J. Gao and H. Tong. Semiparametric non-linear time series model selection. *Journal of the Royal Statistical Society: Series B*, 66(2):321–336, 2004.
- A. Garivier. Consistency of the unlimited bic context tree estimator. *IEEE Transactions on Information theory*, 52(10):4630–4635, 2006.
- A. Goldenshluger and A. Zeevi. Nonasymptotic bounds for autoregressive time series modeling. *Annals of statistics*, pages 417–444, 2001.
- U. S. government. Particulate matter (pm) pollution. <http://www.epa.gov/pm-pollution/particulate-matter-pm-basics>, 2017.
- E. Hannan. The estimation of the order of an arma process. *The Annals of Statistics*, 8(5):1071–1081, 1980.
- E. J. Hannan and M. Deistler. *The statistical theory of linear systems*. SIAM, 2012.
- E. J. Hannan and B. G. Quinn. The determination of the order of an autoregression. *Journal of the Royal Statistical Society: Series B (Methodological)*, 41(2):190–195, 1979.
- R. R. Hocking and R. Leslie. Selection of the best subset in regression analysis. *Technometrics*, 9(4):531–540, 1967.

- A. Hoerl and R. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- D. Hsu, S. M. Kakade, and T. Zhang. An analysis of random design linear regression. *arXiv preprint arXiv:1106.2363*, 2011.
- H.-L. Hsu, C.-K. Ing, and H. Tong. On model selection from a finite family of possibly misspecified time series models. *The Annals of Statistics*, 47(2):1061–1087, 2019a.
- H.-L. Hsu, C.-K. Ing, H. Tong, et al. On model selection from a finite family of possibly misspecified time series models. *The Annals of Statistics*, 47(2):1061–1087, 2019b.
- C. Hurvich and C.-L. Tsai. Regression and time series model selection in small samples. *Biometrika*, 76(2):297–307, 1989.
- C.-K. Ing. Accumulated prediction errors, information criteria and optimal forecasting for autoregressive time series. *The Annals of Statistics*, 35(3):1238–1277, 2007.
- C.-K. Ing and C.-Z. Wei. On same-realization prediction in an infinite-order autoregressive process. *Journal of Multivariate Analysis*, 85(1):130–155, 2003.
- C.-K. Ing and C.-Z. Wei. Order selection for same-realization predictions in autoregressive processes. *The Annals of Statistics*, 33(5):2423–2474, 2005.
- C.-K. Ing, C.-Z. Wei, et al. Order selection for same-realization predictions in autoregressive processes. *The Annals of Statistics*, 33(5):2423–2474, 2005.
- C.-K. Ing, C.-Y. Sin, and S.-H. Yu. Model selection for integrated autoregressive processes of infinite order. *Journal of Multivariate Analysis*, 106:57–71, 2012.
- T. Jeantheau. Strong consistency of estimators for multivariate arch models. *Econometric Theory*, 14(1):70–86, 1998.
- G. Kapetanios. Model selection in threshold models. *Journal of Time Series Analysis*, 22(6):733–754, 2001.
- R. L. Kashyap. Optimal choice of ar and ma parts in autoregressive moving average models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (2):99–104, 1982.
- R. Kass, L. Tierney, and J. Kadane. The validity of posterior expansions based on laplace’s method. *Essays in Honor of George Barnard*, eds. S. Geisser, J. S. Hodges, S. J. Press, and A. Zellner, Amsterdam: North-Holland, pages 473–188, 1990.
- W. Kengne. Strongly consistent model selection for general causal time series. *Statistics & Probability Letters*, page 109000, 2020.
- T. Klein, E. Rio, et al. Concentration around the mean for maxima of empirical processes. *The Annals of Probability*, 33(3):1060–1077, 2005.
- A. Kock. Consistent and conservative model selection with the adaptive lasso in stationary and nonstationary autoregressions. *Econometric Theory*, 32(1):243–259, 2016.
- E. Kounias and T. Weng. An inequality and almost sure convergence. *The Annals of Mathematical Statistics*, 40(3):1091–1093, 1969.

- É. Lebarbier. Detecting multiple change-points in the mean of gaussian process by model selection. *Signal processing*, 85(4):717–736, 2005.
- E. Lebarbier and T. Mary-Huard. *Le critère BIC: fondements théoriques et interprétation*. PhD thesis, INRIA, 2004.
- M. Lerasle. Optimal model selection for density estimation of stationary data under various mixing conditions. *The Annals of Statistics*, 39(4):1852–1877, 2011.
- M. Lerasle et al. Optimal model selection for density estimation of stationary data under various mixing conditions. *The Annals of Statistics*, 39(4):1852–1877, 2011.
- G. Li and W. Li. Least absolute deviation estimation for fractionally integrated autoregressive moving average time series models with conditional heteroscedasticity. *Biometrika*, 95(2):399–414, 2008.
- K.-C. Li. Asymptotic optimality for  $c_p, c_l$ , cross-validation and generalized cross-validation: Discrete index set. *The Annals of Statistics*, 15(3):958–975, 1987.
- W. Li. On the asymptotic standard errors of residual autocorrelations in nonlinear time series modelling. *Biometrika*, 79(2):435–437, 1992.
- W. Li and T. Mak. On the squared residual autocorrelations in non-linear time series with conditional heteroskedasticity. *Journal of Time Series Analysis*, 15(6):627–636, 1994.
- S. Ling and W.-K. Li. Diagnostic checking of nonlinear multivariate time series with multivariate arch errors. *Journal of Time Series Analysis*, 18(5):447–464, 1997.
- S. Ling and M. McAleer. Asymptotic theory for a vector arma-garch model. *Econometric theory*, 19(2):280–310, 2003.
- J. Lv and J. Liu. Model selection principles in misspecified models. *Journal of the Royal Statistical Society: Series B*, 76(1):141–167, 2014.
- C. Mallows. Some comments on  $cp$ . *Technometrics*, 15(4):661–675, 1973.
- P. Massart. *Concentration inequalities and model selection*. Springer, 2007.
- A. McQuarrie and C. Tsai. *Regression and Time Series Model Selection*. World Scientific Pub Co Inc, 1998.
- F. Papangelou. On a distributional bound arising in autoregressive model fitting. *Journal of applied probability*, 31(2):401–408, 1994.
- C. Rao, Y. Wu, S. Konishi, and R. Mukerjee. On model selection. *Lecture Notes-Monograph Series*, pages 1–64, 2001.
- Y. Ren and X. Zhang. Subset selection for vector autoregressive processes via adaptive lasso. *Statistics & probability letters*, 80(23-24):1705–1712, 2010.
- G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6:461–464, 1978.
- S. L. Sclove. Application of model-selection criteria to some problems in multivariate analysis. *Psychometrika*, 52(3):333–343, 1987.



- Q. Shao and L. Yang. Oracally efficient estimation and consistent model selection for auto-regressive moving average time series with trend. *Journal of the Royal Statistical Society: Series B*, 79(2):507–524, 2017.
- P. Shi and C.-L. Tsai. Regression model selection-a residual likelihood approach. *Journal of the Royal Statistical Society: Series B*, 64(2):237–252, 2002.
- R. Shibata. Asymptotically efficient selection of the order of the model for estimating parameters of a linear process. *The Annals of Statistics*, pages 147–164, 1980.
- R. Shibata. Consistency of model selection and parameter estimation. *Journal of Applied Probability*, pages 127–141, 1986.
- C.-Y. Sin and H. White. Information criteria for selecting possibly misspecified parametric models. *Journal of Econometrics*, 71(1-2):207–225, 1996.
- M. Stone. Cross-validatory choice and assessment of statistical predictions. *Journal of the royal statistical society. Series B*, pages 111–147, 1974.
- D. Straumann and T. Mikosch. Quasi-maximum-likelihood estimation in conditionally heteroscedastic time series: A stochastic recurrence equations approach. *The Annals of Statistics*, 34(5):2449–2495, 2006.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B*, pages 267–288, 1996.
- R. Tsay. Order selection in nonstationary autoregressive models. *The Annals of Statistics*, 12(4):1425–1433, 1984.
- Y. Tse and X. Zuo. Testing for conditional heteroscedasticity: Some monte carlo results. *Journal of Statistical Computation and Simulation*, 58(3):237–253, 1997.
- S. A. van de Geer. On hoeffding’s inequality for dependent random variables. In *Empirical process techniques for dependent data*, pages 161–169. Springer, 2002.
- H. White. Maximum likelihood estimation of misspecified models. *Econometrica*, pages 1–25, 1982.
- H. Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101(476):1418–1429, 2006.
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B*, 67(2):301–320, 2005.