

Saskia Rinke and Philipp Sibbertsen\*

# Information criteria for nonlinear time series models

DOI 10.1515/snde-2015-0026

**Abstract:** In this paper the performance of different information criteria for simultaneous model class and lag order selection is evaluated using simulation studies. We focus on the ability of the criteria to distinguish linear and nonlinear models. In the simulation studies, we consider three different versions of the commonly known criteria AIC, SIC and AICc. In addition, we also assess the performance of WIC and evaluate the impact of the error term variance estimator. Our results confirm the findings of different authors that AIC and AICc favor nonlinear over linear models, whereas weighted versions of WIC and all versions of SIC are able to successfully distinguish linear and nonlinear models. However, the discrimination between different nonlinear model classes is more difficult. Nevertheless, the lag order selection is reliable. In general, information criteria involving the unbiased error term variance estimator overfit less and should be preferred to using the usual ML estimator of the error term variance.

**Keywords:** information criteria; Monte Carlo; nonlinear time series; threshold models.

**JEL Numbers:** C15; C22

## 1 Introduction

In time series analysis the identification of a model that is able to appropriately describe special features of a given data set, like cyclical behavior or persistence of shocks in the series, is crucial. This is due to the fact that fitting a misspecified model to the data will lead to biased estimates and all further inference based on previous results, e.g. forecasting, will be misleading.

In order to identify the best fitting model there are two different strands of procedures in the literature, namely hypothesis testing and model selection using information criteria. In the context of linear time series models it is common practice to determine the lag order using information criteria. However, when nonlinear models are considered, the testing approach is preferred. So, instead of calculating information criteria for different models, linearity tests are applied (cf. Luukkonen, Saikkonen, and Teräsvirta 1988a,b; Tong 1990). In a first step a linear AR process is fitted to the data which lag length is determined using information criteria. Afterwards, this specification is tested against a nonlinear alternative (cf. Luukkonen, Saikkonen, and Teräsvirta 1988a). Pitarakis 2006 shows that the lag order selection can seriously influence the power properties of linearity tests since the linear model will be misspecified if the true data generating process (DGP) is actually nonlinear. Moreover, Luukkonen, Saikkonen, and Teräsvirta (1988b) show that linearity tests designed to detect a certain kind of nonlinearity may also have power against other nonlinear models. Hence, the rejection of the null of a linear model may not tell which nonlinear model should be used to model the data.

---

\*Corresponding author: Philipp Sibbertsen, Institute of Statistics, Leibniz University Hannover, School of Economics and Management, Königsworther Platz 1, D-30167 Hannover, Germany, Tel.: +49-511-762-3783, Fax: +49-511-762-3923, e-mail: sibbertsen@statistik.uni-hannover.de

Saskia Rinke: Institute of Statistics, Leibniz University Hannover, School of Economics and Management, Königsworther Platz 1, D-30167 Hannover, Germany

Despite these drawbacks of testing, there are probably two reasons why the testing approach is preferred to model selection for nonlinear models. Firstly, it may not be clear how to calculate the value of an information criterion if a multiple regime model is fitted to the data. Information criteria can be easily calculated for single regimes, but then, these values have to be combined into one index in order to obtain one value for the whole model. Secondly, there exists no general rule if and how additional parameters of the nonlinear model have to be incorporated into the penalty terms of information criteria. Hence, the application of information criteria to nonlinear models may not result in the selection of optimal models (cf. Clements and Krolzig 1998).

Nonetheless, in the literature there exist some examples of the application of information criteria to nonlinear time series models. Emiliano, Vivanco, and De Menezes (2014) assess the performance of information criteria for lag order selection in linear time series models and model selection in nonlinear growth models. Lag order selection is treated in Kapetanios (2001) for SETAR and MSAR models, in Smith, Naik, and Tsai (2006) for MSAR models and in Tong (1983), Wong and Li (1998) and Li (1988) for SETAR models. Hamaker (2009) and Gonzalo and Pitarakis (2002) use information criteria to distinguish between linear AR models and single- and multiple-regime SETAR models. Further selection of the model class is considered among others in Psaradakis et al. (2009) and in Kapetanios (2001). Though, in contrast to Psaradakis et al. (2009) and Gonzalo and Pitarakis (2002), Kapetanios (2001) incorporates the threshold parameters of the SETAR models into the penalty terms of information criteria. Additional to the usual information criteria Hamaker (2009) also considers a version for change-point estimation, which penalizes thresholds as additional parameters.

In most works the lag orders and other relevant parameters, like the delay parameter of SETAR and STAR models, are treated as given. Thus, the results are obtained under ideal conditions. In fact, if information criteria are applied to empirical data, the parameter values are unknown in advance and have to be estimated first. The resulting estimation errors can influence further calculations and hence, deteriorate the performance of the information criteria so that former results are not valid anymore. Therefore, in this work we apply different information criteria to select the optimal model class, the corresponding lag order and additional parameters simultaneously. The performance of the information criteria in different scenarios is assessed in several simulation studies. There, we will take three different versions of the respective criteria into account. Special focus will be on the fact, whether the criteria are able to successfully distinguish between linear and nonlinear time series models.

The rest of the paper is organized as follows. In Section 2 the nonlinear models which are considered in the simulation studies are explained. In Section 3 we shortly repeat the intuition of information criteria and introduce the four criteria we use. In Section 4 the simulation set-up and the simulation results are presented. Finally, Section 5 concludes.

## 2 Nonlinear time series models

There exists a variety of different nonlinear time series models in the literature. In our simulation studies we focus on regime-switching models with switches in the mean equation. Hence, the class of ARCH (cf. Engle 1982) and GARCH (cf. Bollerslev 1986) models is not considered. The selection of ARCH/GARCH orders is treated e.g. in Hughes and King (2003) and Hughes, King, and Kwek (2004).

Instead, we focus on SETAR and STAR models since they are frequently applied and have a similar switching mechanism. Actually, the LSTAR model nests the SETAR model for  $\gamma \rightarrow \infty$ . In addition to SETAR and STAR models there also exists the class of Markov-switching autoregressive (MSAR) models (cf. Hamilton 1989, 1994). However, in contrast to SETAR and STAR models, where the change of regime is governed by an endogenous variable, the regime shift in MSAR models is controlled by an exogenous, unobservable state variable. Due to this difference we do not consider MSAR models in our simulation studies but focus on the ability of the information criteria to discriminate between linear and nonlinear models and to detect the correct form of the transition function.

## 2.1 The SETAR model

Self-exciting threshold autoregressive models were introduced in Tong and Lim (1980) and Tong (1983) (cf. also Tong 1990). Since linear AR models are not able to capture certain nonlinear features of the data, but are easy to specify, the SETAR model is a natural extension of the linear model to the nonlinear case. SETAR models combine multiple piecewise linear regimes, which are separated by threshold parameters, into one model. Hence, SETAR models may be locally linear, but due to the regime changes, which can be interpreted as breaks in the series, the model is globally nonlinear. A two regime SETAR model with  $p_1$  and  $p_2$  lags respectively can be written as

$$y_t = \left( \phi_{0_1} + \sum_{i=1}^{p_1} \phi_{1_i} y_{t-i} \right) \mathbb{1}_{\{y_{t-d} > c\}} + \left( \phi_{0_2} + \sum_{i=1}^{p_2} \phi_{2_i} y_{t-i} \right) \mathbb{1}_{\{y_{t-d} \leq c\}} + \varepsilon_t, \quad (1)$$

where  $\varepsilon_t \sim iid(0, \sigma^2)$  and  $\mathbb{1}_{\{ \cdot \}}$  denotes the indicator function. The dependent variable  $y_t$  falls into the first regime, if the threshold variable  $y_{t-d}$  exceeds the threshold  $c$ . Otherwise  $y_t$  falls into the second regime. In SETAR models the threshold variable is a lagged value of the dependent variable. The lag  $d$  is called the delay parameter and does not exceed the largest lag length (cf. Pitarakis 2006).

## 2.2 The STAR model

In contrast to the SETAR models, in STAR models the regime switches are not discrete jumps but smooth transitions. Therefore, the indicator function in Equation (1) is replaced by the transition function  $F_t(\cdot)$ , which takes values in the unit interval. This implies that each observation does not lie in one single regime but is a weighted mixture of both regimes, where the transition function attaches the weights  $1-F_t(\cdot)$  and  $F_t(\cdot)$  to the first and second regime respectively. Depending on the form of the transition function, the STAR model is a logistic STAR (LSTAR) or an exponential STAR (ESTAR) model,

$$F_t(y_{t-d}, \gamma, c) = \begin{cases} 1 / (1 + \exp(-\gamma(y_{t-d} - c))), & \text{for LSTAR;} \\ 1 - \exp(-\gamma(y_{t-d} - c)^2), & \text{for ESTAR.} \end{cases} \quad (2)$$

Like in SETAR models  $c$  denotes the threshold and  $d$  the delay parameter. The parameter  $\gamma$  regulates the speed of the transition between the regimes. For a value of  $\gamma=0$  there exists no regime shift, instead the model is linear, for  $\gamma \rightarrow \infty$  the LSTAR model becomes a two-regime SETAR model, whereas the ESTAR model reduces to a linear model (cf. Luukkonen, Saikkonen, and Teräsvirta 1988a; Teräsvirta 1994; van Dijk, Teräsvirta, and Franses 2002).

## 3 Information criteria

The idea of information criteria is to balance the goodness of fit and the complexity of a model using a loss function  $L = G(\hat{\sigma}^2) + P(n, p)$  (cf. Wu and Sepulveda 1998). The first term of the loss function accounts for the goodness of fit and depends on an estimate of the unknown error term variance. The smaller the estimated variance of the error terms, the better is the model fit. The second term is the penalty term which depends on the sample size  $n$  and on the number of parameters  $p$ . Minimizing the loss function guarantees that if two models yield the same model fit, the model which contains fewer parameters is preferred. This is also known as the principle of parsimony (cf. Akaike 1974; Schwarz 1978).

### 3.1 The traditional information criteria

Different choices of the penalty term yield different information criteria. The Akaike Information Criterion (AIC) (cf. Akaike 1974) and the Schwarz Information Criterion (SIC) (cf. Schwarz 1978) are the most commonly used information criteria,

$$\text{AIC} = n(\log(\hat{\sigma}^2) + 1) + 2(p+1), \quad (3)$$

$$\text{SIC} = n\log(\hat{\sigma}^2) + p\log(n). \quad (4)$$

The AIC has a rather weak penalty term which can result in overfitting (selecting too many parameters) in finite samples (cf. Hurvich and Tsai 1989), whereas the SIC has a stronger penalty term in order to prevent overfitting which can lead to underfitting (selection of too few parameters) in small samples. Since the AIC is biased and therefore tends to overfit in finite samples, Hurvich and Tsai (1989) introduced a bias-corrected version of the AIC, the corrected Akaike Information Criterion (AICc),

$$\text{AICc} = n\log(\hat{\sigma}^2) + \frac{n(n+p)}{n-p-2}. \quad (5)$$

For small samples the AICc has a stronger penalty term than the AIC to solve the problem of overfitting. Asymptotically both versions are equivalent. Wu and Sepulveda (1998) introduced the Weighted Average Information Criterion (WIC) in order to combine the strengths of different criteria to obtain a criterion which performs well not depending on the sample size,

$$\text{WIC} = n\log(\hat{\sigma}^2) + \frac{(2n(p+1)/(n-p-2))^2 + (p\log(n))^2}{2n(p+1)/(n-p-2) + p\log(n)}. \quad (6)$$

The WIC is a weighted version of the AICc and the SIC. Setting the weights equal to the penalty terms of the respective criteria guarantees that the WIC behaves like AICc in small samples and like SIC in large samples. Since both criteria perform well in the respective sample size, these properties of WIC may be very valuable if we apply WIC separately to the regimes e.g. of a SETAR model with a dominant regime. In this case, one regime contains significantly more observations than others. But due to the independence of the sample size, WIC should perform well in all regimes. We will further discuss and evaluate this point in Section 4.

### 3.2 The versions of information criteria

As already mentioned in Section 1, it is not clear how to calculate the value of an information criterion for multiple-regime models. Since each information criterion depends on an estimated error term variance and the number of parameters of the model under consideration, it would be straightforward to estimate the error term variance using the residual sum of squares of the whole model and add the number of parameters of all regimes in order to obtain the number of parameters of the whole model (cf. Pitarakis 2006). We call the resulting information criteria **Overall Model Criteria**. However, it is not clear whether information criteria maintain their optimality properties of linear specification, when they are applied to nonlinear models (cf. Clements and Krolzig 1998). Thus, we also follow another approach. We consider two additional versions of information criteria, where we separate the models into their regimes. Due to the fact that the single regimes are linear, information criteria are supposed to select optimal lag orders of the regimes. Following the approach of Tong (1983), for the **Equally Weighted Criteria** we calculate the information criteria separately for each regime and then combine these values into one model information criterion, where each regime gets the equal weight

$$IC_{\text{model}} = \frac{1}{m} \sum_{i=1}^m IC_{\text{reg},i}. \quad (7)$$

For the **Regime Weighted Criteria** the weighting of the regimes is proportional to the dominance of the regime. The more dominant the respective regime is, the higher is the attached weight  $w_i$

$$IC_{model} = \sum_{i=1}^m w_i IC_{reg.i} \quad \text{with } w_i \in [0, 1]. \quad (8)$$

For SETAR models the weights can simply be determined by the the number of observations that fall into the respective regime divided by the total number of observations. In STAR models, the value of the transition function can be used as a weight, since in STAR models, the observations do not fall in one regime only but are a weighted sum of both regimes.

The differentiation between these three versions is only meaningful for regime-switching models. Hence, for the linear AR model, which only consists of one single regime, all three versions are equivalent.

### 3.3 The role of the error term variance estimator

All previously introduced information criteria depend on an estimate of the error term variance. Generally, the Maximum Likelihood estimator

$$\hat{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n}, \quad (9)$$

where  $\hat{\varepsilon}$  is the vector of residuals and  $n$  denotes the sample size, is used to calculate the values of the information criteria. However,  $\hat{\sigma}^2$  is a biased estimator of the true error term variance. Therefore, McQuarrie, Shumway, and Tsai (1997) suggest to use

$$\tilde{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n-p-1} \quad (10)$$

or the sample variance  $s^2$  to evaluate different models using information criteria (cf. also Wen and Tu 2001). According to McQuarrie, Shumway, and Tsai (1997) the information criteria involving  $\tilde{\sigma}^2$  have a stronger penalty term and avoid overfitting.

## 4 Simulation study

In the following simulation studies we evaluate the information criteria and their different versions presented in Section 3.

### 4.1 Data generation

In order to identify factors that influence the performance of the information criteria we generate data from different linear and nonlinear models and vary the sample size, the persistence parameters, the degree of dominance and the number of regimes. The models of the basic set-up are tabulated in the Supplementary Material, where the  $\varepsilon_t$  form a Gaussian white noise process. Following Kapetanios (2001) for every sample size, we simulate 200 additional observations for each DGP, which are discarded afterwards to avoid a starting value bias. All initial values are set to zero. The simulation results are based on 1000 replications.

### 4.2 Model estimation

After generating the data, we fit different linear and nonlinear models to the data, calculate the different versions of the information criteria and choose the model, which minimizes the respective information criterion.

Throughout, we assume that the error terms  $\varepsilon_t$  are Gaussian. In order to minimize the computational effort, we set a maximum lag length of  $p_{\max}=4$  for the models fitted to the data. Thus, the largest AR model is an AR(4) and the largest regime switching models consist of an AR(4) specification in each regime. According to Luukkonen, Saikkonen, and Teräsvirta (1988a) a lag order exceeding  $p=3$  is rather unlikely for a small sample size of  $n=50$ , but probable for larger samples. In order to make the results among different sample sizes comparable, we choose  $p_{\max}=4$  for all sample sizes (cf. also Tong and Lim 1980; Pitarakis 2006). The effective sample size used to fit the models to the data and in all further calculations is thus  $n-p_{\max}$  (cf. Tong and Lim 1980; Wong and Li 1998).

Parameter estimation is done by (conditional) least squares. For the threshold and delay parameter grids are constructed and the remaining parameters are estimated for each grid point. Following Hansen (1997) the grid of the threshold consists of the interval from the 15% to the 85% quantile of  $y_t$ . Disregarding the lower and upper 15% quantiles should guarantee that at least 15% of the data lie in each regime and hence, the number of observations in each regime is sufficient for persistence parameter estimation. The grid of the delay simply consists of integer values from 1 to  $p_{\max}$ . For each model the parameter combination which minimizes the residual sum of squares is selected and the corresponding values of the information criteria are calculated. In STAR models the grid search is done conditional on  $\gamma$ . According to Teräsvirta (1994) it is possible to standardize the exponent of the transition function and choose  $\gamma=1$  as a starting value. After determining all parameters the value of  $\gamma$  is adjusted by minimizing the residual sum of squares with respect to  $\gamma$ . For ESTAR models the exponent of the transition function is divided by the variance of  $y_t$ , whereas for LSTAR models the standard deviation of  $y_t$  is appropriate.

As already mentioned there exists no general rule if and how additional parameters of nonlinear models like the threshold and the delay should be incorporated into the penalty terms of information criteria. We decide to follow the approach of Kapetanios (2001) and add all additional parameters of the nonlinear models to the number of parameters. The intuition is that if the true DGP is nonlinear, then a nonlinear model will provide a better fit to the data. Hence, the values of the information criteria will decrease. However, the computational effort will increase. Since the information criteria are supposed to balance model fit and complexity, additional parameters are incorporated. As a result, the nonlinear model will only be selected if it provides a substantially better fit than a simpler model. For all three nonlinear models the additional parameters are the threshold parameter  $c$  as well as the delay parameter  $d$ . In STAR models we also consider the transition parameter  $\gamma$  as an additional parameter in the penalty term.

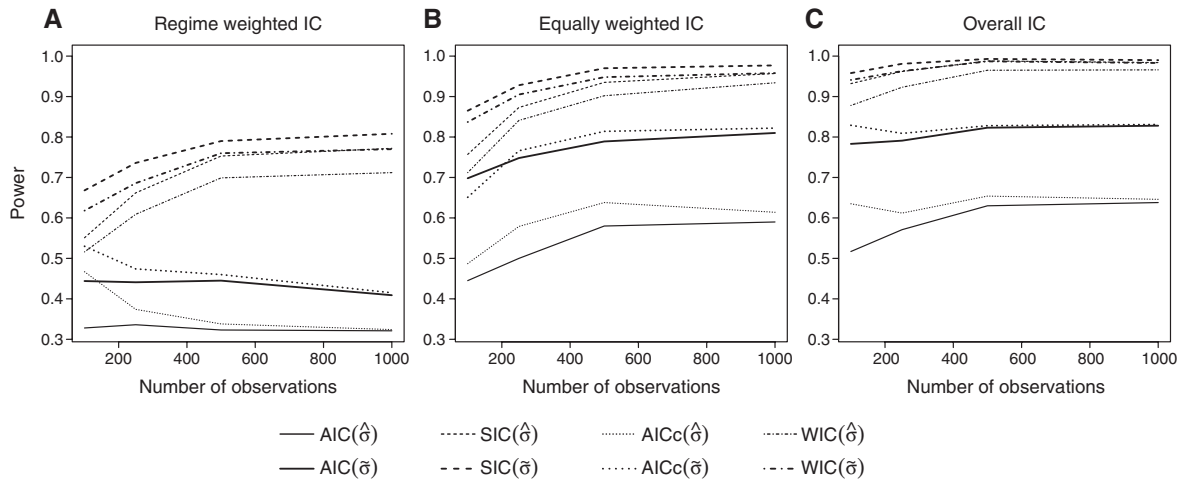
### 4.3 Simulation results

In the following scenarios we assess the ability of the different information criteria to select the correct model. The figures and tables display the respective selection frequencies. In order to evaluate the performance let the power of an information criterion be defined as the relative frequency of selecting the correct model.

#### 4.3.1 Lag order selection

This paragraph focuses on lag order selection within a certain model class. Hence, the following figures display the power of the information criteria when only the lag order (combination) within the true model class has to be determined. Although the assumption of knowing the correct model class will not be met if the information criteria are applied to real data, the power results will be helpful to correctly interpret the performance of the information criteria in further simulation studies. This is due to the fact, that if the criteria are not able to determine the correct lag order (combination) within the true model class, we cannot expect them to point to the correct model when also the model class has to be selected. In Figure 1 the power of the three variants of information criteria are depicted when the true DGP is the LSTAR(1,1) model. It shows first characteristics of the information criteria. The regime weighted criteria perform worse than their equally weighted





**Figure 1:** Power of the information criteria (IC) for LSTAR models (LSTAR(1,1)).

and overall counterparts. So, the selection frequencies vary between 30% and 80% in large samples. In contrast, the equally weighted and overall versions have correct selection frequencies of about 60% to 100% in large samples. Thus, the regime weighted information criteria are clearly outperformed which is also true for small samples.

Furthermore, all versions of AIC and AICc cannot select the correct model with a probability approaching 1 when the sample size increases. This is due to the fact, that AIC and AICc are not consistent information criteria (cf. Shibata 1986). Instead of the correct lag order combination AIC and AICc tend to overfit and choose larger lag order combinations. However, AICc performs better than AIC in small and moderate samples (cf. also Wong and Li 1998). Using the unbiased error term variance estimator  $\tilde{\sigma}^2$  reduces the probability of overfitting. So, the versions of AIC and AICc involving  $\tilde{\sigma}^2$  improve up to 20 percentage points (cf. Figure 1B and C).

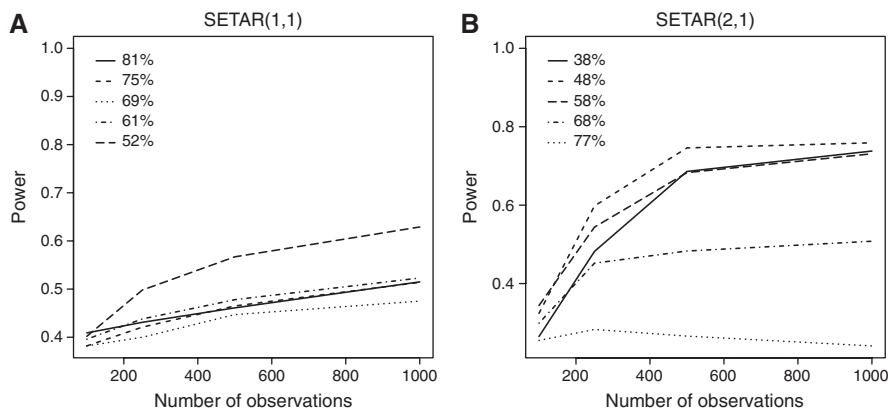
In general, all information criteria perform better in small models. With increasing lag order combinations, there is a tendency to underestimate the lag order of one regime independent of the error term variance estimator. Naturally, AIC and AICc outperform SIC and WIC in these cases due to their tendency to overfit. However, underfitting can occur due to weakly identifiable models (cf. McQuarrie, Shumway, and Tsai 1997). If the largest lag order only has a minor influence, it may be neglected and a smaller model is preferred.

For lag order selection we recommend to apply equally weighted and overall SIC and WIC. In small samples AICc( $\tilde{\sigma}^2$ ) also performs well.

#### 4.3.2 Dominant regimes

The idea to include the regime weighted information criteria and also the WIC into the simulation study is that they may lead to better results if one regime is dominant in the model. Therefore, we consider some modified DGPs. From the SETAR(1,1) and the SETAR(2,1) model we generate data with five different thresholds, resulting in a first regime with a share of observations varying from 38% to 81%. Afterwards, we fit SETAR models to the data and assess the effect of dominant regimes.

In Figure 2 we see that the power of the regime weighted AIC with the unbiased error term variance estimator  $\tilde{\sigma}^2$  is highest if the first regime is not more dominant than the second. So, the best performance is achieved if 52% of the observations fall into the first regime for the SETAR(1,1) and 48% for the SETAR(2,1), respectively. But in fact, the size of the effect is quite different. Although for the SETAR(1,1) the fair separation of regimes is definitely the best, the difference between the best and the worst separation (with 69% of the observations in the first regime) only amounts to approximately 15 percentage points in large samples. In contrast, for the SETAR(2,1) model the fair separation yields similar results as if 58% and asymptotically similar results as if 38% of the observations lie in the first regime. Nevertheless, the difference between the



**Figure 2:** Power of the information criteria in dominant regimes: regime weighted AIC with unbiased error term variance estimator.

best and the worst separation (with 77% of the observations in the first regime) amounts to 50 percentage points in large samples. In small samples, the difference between dominant and fair regimes is rather small. The comparison points out another fact: The worst separation is not always the one with one very dominant regime (cf. Figure 2A). These considerations are valid for all regime weighted information criteria. Though, the effects are more pronounced for AIC and AICc. The equally weighted and overall information criteria are more independent from the dominance of the regimes. So asymptotically all separations yield the same power results. For moderate samples the effect is mostly pronounced.

Hence, the performance of the regime weighted information criteria is not proportional to the degree of dominance of one regime. Instead, the equally weighted and overall criteria outperform their regime weighted counterparts.

#### 4.3.3 Size of discriminating linear and nonlinear model

The main part of this paper focuses on power properties of the information criteria. In fact, good power results are only meaningful if the size of the information criteria is under control. Therefore, we briefly illustrate the size properties of the information criteria under consideration. Figure 3 displays the selection frequencies of a nonlinear model if the true DGP is linear. It points out that weighted versions of AIC and AICc never select the linear model. The overall AIC and AICc are still oversized. All regime weighted information criteria are heavily oversized (size  $\geq 0.5$ ) and thus not very reliable. The size even increases with an increasing number of observations. The equally weighted SIC and WIC are oversized but perform better than their regime weighted counterparts. The size difference amounts to 15 percentage points. The overall SIC and WIC have a size of below 0.2 in small samples which further decreases with an increasing number of observations and finally approaches zero.

Hence, the equally weighted information criteria may have a tendency to favor nonlinear models whereas the overall versions may select the linear models more often. The following simulation results will further assess this point.

According to the size properties, the overall SIC and WIC should be preferred for model class selection. Equally weighted SIC and WIC can be applied, whereas AIC, AICc and regime weighted criteria are clearly outperformed.

#### 4.3.4 Power of discriminating linear and nonlinear models

In this paragraph we assess the performance of the information criteria for selecting between the linear and the correct nonlinear model. Due to the fact that the regime weighted information criteria are inferior in terms of size and power and cannot outperform the other versions in dominant regimes, we will focus on the per-



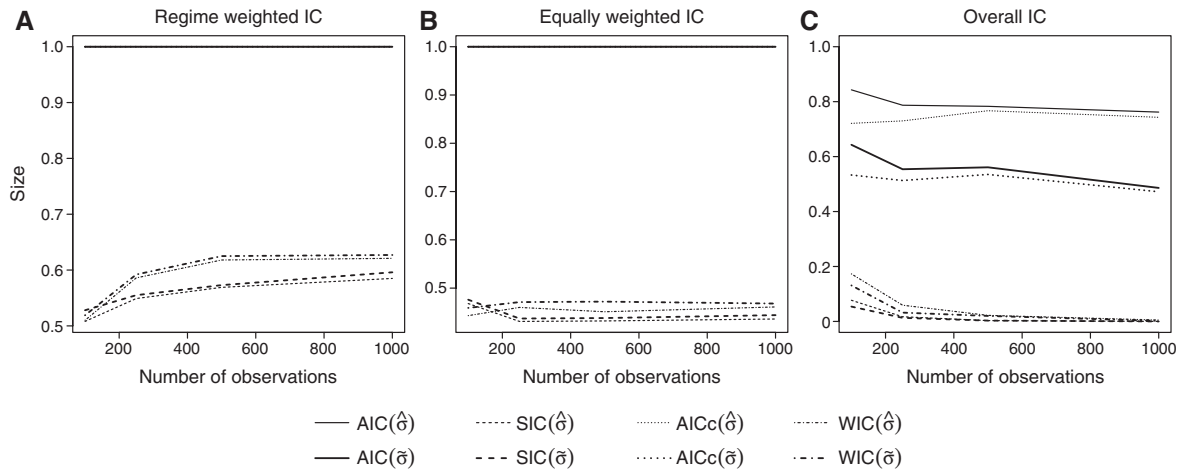


Figure 3: Size of the information criteria (IC) for AR models (AR(3)).

formance of the equally weighted and the overall information criteria. We only give an extract of the whole tables for  $n=100$  and  $n=1000$ . Further results for  $n=250$  and  $n=500$  and for the regime weighted versions can be found in the Supplementary Material. In the following tables the blue row indicates the correct model, whereas the bold numbers mark the models with the highest selection frequency for the respective criterion. In Table 1 the selection frequencies are given when the data is generated by a SETAR(1,1) process. It strikes again that weighted versions of AIC and AICc never select a linear model. Although it is the correct decision in this setting, this behavior is spurious. This fact is also pointed out by Gonzalo and Pitarakis (2002) and Pitarakis (2006). In contrast, SIC and WIC have the tendency to select the linear model for a small sample of  $n=100$  (cf. Psaradakis et al. 2009). However, taking into account the results for  $n=1000$ , it becomes obvious that this is a problem in small samples due to the fact that for the increased sample size all criteria prefer the correct model.

Comparing these results with Table 2, it becomes evident that the tendency to select linear models also depends on the true model structure. So, in symmetric models, where the lag order is equal in both regimes, SIC and WIC tend to prefer the linear over the nonlinear model. If the lag orders differ among regimes, it

Table 1: Selection frequencies of the information criteria: AR vs. SETAR models for SETAR(1,1) DGP.

$n$	AIC( $\hat{\sigma}$ )		AIC( $\bar{\sigma}$ )		SIC( $\hat{\sigma}$ )		SIC( $\bar{\sigma}$ )		AICc( $\hat{\sigma}$ )		AICc( $\bar{\sigma}$ )		WIC( $\hat{\sigma}$ )		WIC( $\bar{\sigma}$ )	
	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000
(a) Equally weighted information criteria																
AR(1)	0.000	0.000	0.000	0.000	<b>0.492</b>	0.106	<b>0.516</b>	0.106	0.000	0.000	0.000	0.000	<b>0.461</b>	0.090	0.488	0.091
SETAR(1,1)	<b>0.252</b>	<b>0.465</b>	<b>0.525</b>	<b>0.703</b>	0.299	<b>0.838</b>	0.380	<b>0.868</b>	<b>0.473</b>	<b>0.483</b>	<b>0.693</b>	<b>0.710</b>	0.319	<b>0.813</b>	<b>0.405</b>	<b>0.849</b>
SETAR(1,2)	0.115	0.130	0.114	0.095	0.049	0.030	0.031	0.016	0.124	0.123	0.107	0.093	0.052	0.047	0.039	0.030
SETAR(1,3)	0.093	0.066	0.054	0.040	0.026	0.002	0.013	0.002	0.059	0.067	0.020	0.037	0.022	0.009	0.007	0.002
SETAR(1,4)	0.105	0.053	0.045	0.020	0.021	0.002	0.007	0.000	0.045	0.049	0.018	0.020	0.018	0.003	0.005	0.002
SETAR(2,1)	0.111	0.105	0.105	0.076	0.043	0.015	0.027	0.003	0.124	0.106	0.078	0.077	0.040	0.026	0.029	0.018
(b) Overall information criteria																
AR(1)	0.050	0.000	0.178	0.000	<b>0.806</b>	0.012	<b>0.863</b>	0.019	0.100	0.000	0.238	0.000	<b>0.591</b>	0.002	<b>0.710</b>	0.003
SETAR(1,1)	<b>0.272</b>	<b>0.501</b>	<b>0.479</b>	<b>0.728</b>	0.143	<b>0.952</b>	0.108	<b>0.970</b>	<b>0.342</b>	<b>0.506</b>	<b>0.491</b>	<b>0.734</b>	0.272	<b>0.935</b>	0.232	<b>0.958</b>
SETAR(1,2)	0.105	0.110	0.086	0.081	0.012	0.020	0.008	0.008	0.103	0.111	0.079	0.080	0.030	0.029	0.014	0.019
SETAR(1,3)	0.088	0.064	0.037	0.039	0.003	0.001	0.001	0.000	0.081	0.064	0.028	0.037	0.012	0.003	0.004	0.002
SETAR(1,4)	0.071	0.047	0.028	0.023	0.000	0.000	0.000	0.000	0.058	0.047	0.016	0.022	0.006	0.001	0.000	0.000
SETAR(2,1)	0.106	0.101	0.072	0.069	0.008	0.011	0.006	0.002	0.099	0.100	0.063	0.070	0.032	0.024	0.014	0.014

Bold numbers mark the models with the highest selection frequency for the respective criteria.

**Table 2:** Selection frequencies of the equally weighted information criteria: AR vs. SETAR models for SETAR(2,3) DGP.

$n$	AIC( $\tilde{\sigma}$ )		AIC( $\hat{\sigma}$ )		SIC( $\hat{\sigma}$ )		SIC( $\tilde{\sigma}$ )		AICc( $\hat{\sigma}$ )		AICc( $\tilde{\sigma}$ )		WIC( $\hat{\sigma}$ )		WIC( $\tilde{\sigma}$ )	
	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000
AR(1)	0.000	0.000	0.000	0.000	0.166	0.000	0.199	0.000	0.000	0.000	0.000	0.000	0.136	0.000	0.172	0.000
AR(2)	0.000	0.000	0.000	0.000	0.109	0.001	0.091	0.001	0.000	0.000	0.000	0.000	0.134	0.001	0.109	0.001
SETAR(1,1)	0.077	0.000	0.210	0.000	0.176	0.001	<b>0.258</b>	0.001	0.182	0.000	<b>0.322</b>	0.000	0.169	0.001	<b>0.264</b>	0.001
SETAR(1,2)	0.142	0.000	<b>0.238</b>	0.000	<b>0.182</b>	0.006	0.195	0.008	<b>0.214</b>	0.000	0.262	0.000	<b>0.195</b>	0.003	0.205	0.006
SETAR(2,2)	<b>0.148</b>	<b>0.481</b>	0.149	<b>0.686</b>	0.108	<b>0.846</b>	0.090	<b>0.879</b>	0.149	<b>0.509</b>	0.115	<b>0.695</b>	0.112	<b>0.812</b>	0.082	<b>0.856</b>
SETAR(2,3)	0.088	0.215	0.040	0.173	0.027	<b>0.103</b>	0.014	0.067	0.044	0.205	0.020	0.169	0.020	0.133	0.009	0.094

Bold numbers mark the models with the highest selection frequency for the respective criteria.

becomes easier to detect the nonlinearity in small samples. However, the overall SIC and WIC still tend to select the linear model. Table 2 also picks up the problem of underfitting already mentioned in the previous paragraph. The larger lag order combination is not estimated correctly. However, this might be due to identification problems.

The previous remarks on SETAR models are also valid for the selection between AR and STAR models. However, the differentiation between linear and nonlinear models is slightly better for LSTAR than for ESTAR models. This is due to the asymptotic behavior concerning the transition parameter  $\gamma$ . ESTAR models converge to linear models for both extremes  $\gamma \rightarrow 0$  and  $\gamma \rightarrow \infty$ , whereas LSTAR models only become linear for  $\gamma \rightarrow 0$ . Thus, the parameter estimate of the transition variable plays an important role. This will be further discussed in the context of discrimination between nonlinear models.

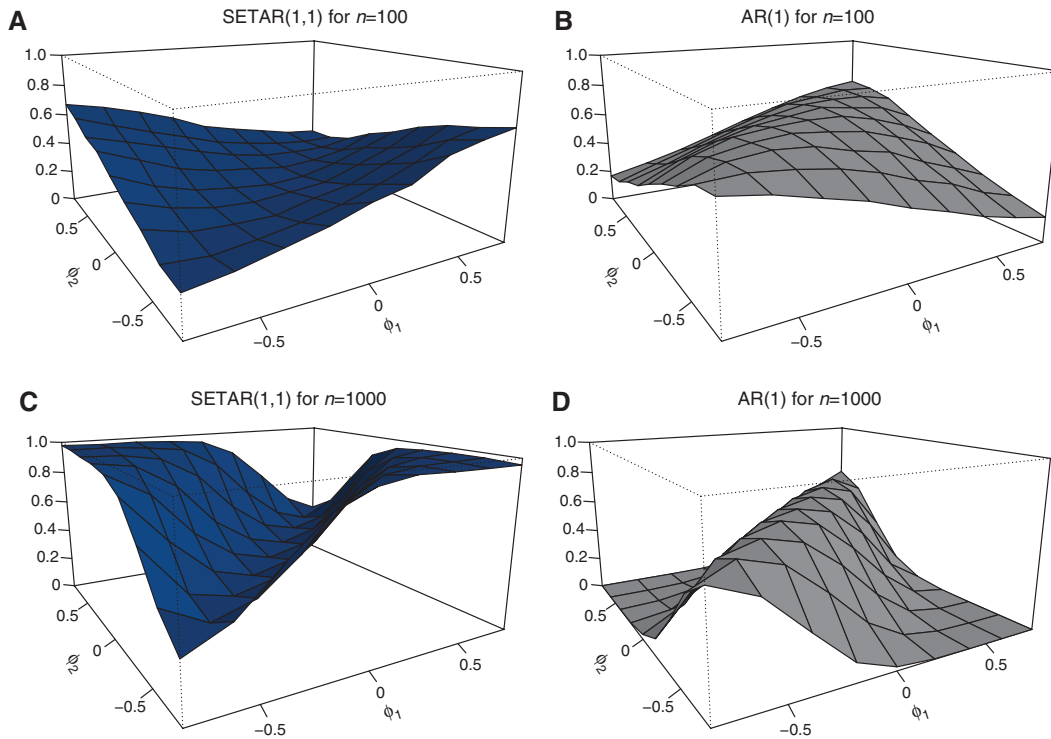
According to the power properties, equally weighted SIC and WIC perform the best. In small samples the overall AIC( $\tilde{\sigma}^2$ ) and AICc( $\tilde{\sigma}^2$ ) can be applied as well. Overall SIC and WIC also perform well in large samples, but tend to select linear models in small samples. This is in line with Emiliano, Vivanco, and De Menezes (2014).

#### 4.3.5 Effects of the persistence parameters

According to Psaradakis et al. (2009) the performance of the information criteria to select between the linear and the nonlinear model is better when the persistence parameters among the regimes differ. The following figures display this fact and the dependence on the sample size. The true DGP is a SETAR(1,1) model with persistence parameters varying from  $-0.8$  to  $0.8$  by  $0.2$ . Since the error terms of both regimes are iid normally distributed, the SETAR(1,1) reduces to an AR(1) model if the persistence parameters are equal. For a small sample size of  $n=100$  the selection frequency of a SETAR(1,1) model increases with the distance between the persistence parameters (cf. Figure 4A). This confirms the findings of Psaradakis et al. (2009) that the differentiation between linear and nonlinear models becomes easier the more the regimes differ. On the diagonal the persistence parameters are equal and the model reduces to an AR(1) process. On this diagonal the selection frequency for the SETAR(1,1) is lowest. However, the selection frequency of the SETAR(1,1) model is already relatively low when the regimes are not equal but quite similar.

This is confirmed by Figure 4B which depicts the respective selection frequency of an AR(1) model. On the diagonal the true model is actually the AR(1) model and there the selection frequency is the highest. Close to the diagonal, the selection frequency is still rather high. This implies that the linear model is preferred to the nonlinear model. For distinct regimes the AR(1) model is clearly inferior to the SETAR(1,1) model.

With an increasing sample size the differentiation between linear and nonlinear models is more reliable. In Figure 4C the selection frequency of the SETAR(1,1) model approaches 1 if the regimes are distinct. The more similar the regimes become, the lower is the selection frequency. Again, it is minimal on the diagonal where the model reduces to the linear case. Additionally, in Figure 4D the respective selection frequencies of the AR(1) model are presented. For distinct persistence parameters the linear model is never selected. The more similar the regimes become, the higher is the selection frequency. Nevertheless, the regimes have to



**Figure 4:** Selection frequency for equally weighted SIC with  $\bar{\sigma}^2$ .

be very similar, otherwise the nonlinear model is superior to the linear model. Hence, there are more correct selections if the number of observations increases.

All versions of SIC and WIC have these properties. The overall versions of AIC and AICc are able to detect linearity but especially in small samples have a tendency to prefer the nonlinear over the linear model. The weighted versions of AIC and AICc cannot detect linearity. They spuriously select the nonlinear model and never the linear model independent of the sample size and the distance between regimes. Thus, these results emphasize the previous findings that weighted versions of AIC and AICc should not be applied if the collection of models includes linear models.

#### 4.3.6 Discriminating nonlinear models

So far we have only considered the selection between linear and nonlinear models, focussing on one nonlinear model class. Now, the collection of models includes also other nonlinear model classes in order to assess the ability of the information criteria to determine the form of the transition function. In a first step we only consider the LSTAR and the ESTAR models (cf. Tables 3 and 4). Our results point out that there are two important factors that influence the information criteria: The sample size and the transition parameter  $\gamma$ . As already pointed out by Psaradakis et al. (2009), it is difficult to determine the switching mechanism when only a small number of observations is available. Moreover, the selection results are better for a small value of  $\gamma=1$  than for  $\gamma=20$ . In fact, in small samples the equally weighted and the overall criteria tend to a spurious selection of the ESTAR model class for larger  $\gamma$  (cf. Table 4), whereas the regime weighted criteria favor LSTAR models. However, the selected lag orders are always similar to the lag order combination preferred if only the correct model class is considered. With an increasing sample size both, lag order and model class selection, improve for all versions of information criteria. Nevertheless, the selection frequencies of the correct model are higher if  $\gamma=1$ .

In the next step we also allow for SETAR models in the collection of models. As a result, we find another spurious behavior of the information criteria. In small samples all equally weighted and overall information

**Table 3:** Selection frequencies of the information criteria: ESTAR vs. LSTAR models for LSTAR(2,2) DGP with  $\gamma=1$ .

$n$	AIC( $\hat{\sigma}$ )		AIC( $\tilde{\sigma}$ )		SIC( $\hat{\sigma}$ )		SIC( $\tilde{\sigma}$ )		AICc( $\hat{\sigma}$ )		AICc( $\tilde{\sigma}$ )		WIC( $\hat{\sigma}$ )		WIC( $\tilde{\sigma}$ )	
	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000
(a) Equally weighted information criteria																
ESTAR(1,2)	0.064	0.000	0.108	0.000	0.101	0.000	0.146	0.000	0.094	0.000	0.125	0.000	0.124	0.000	0.158	0.000
ESTAR(2,1)	0.021	0.000	0.041	0.000	0.047	0.000	0.066	0.000	0.091	0.000	0.125	0.000	0.078	0.000	0.106	0.000
ESTAR(2,2)	0.225	0.061	0.274	0.083	0.295	0.132	0.277	0.136	0.172	0.063	0.171	0.081	0.232	0.120	0.206	0.126
LSTAR(1,1)	0.004	0.000	0.022	0.000	0.031	0.000	0.048	0.000	0.008	0.000	0.019	0.000	0.025	0.000	0.042	0.000
LSTAR(2,1)	0.017	0.000	0.029	0.000	0.046	0.000	0.071	0.000	0.037	0.000	0.058	0.000	0.038	0.000	0.058	0.000
LSTAR(2,2)	<b>0.237</b>	<b>0.601</b>	<b>0.314</b>	<b>0.740</b>	<b>0.271</b>	<b>0.836</b>	<b>0.285</b>	<b>0.850</b>	<b>0.317</b>	<b>0.611</b>	<b>0.347</b>	<b>0.746</b>	<b>0.345</b>	<b>0.820</b>	<b>0.346</b>	<b>0.845</b>
(b) Overall information criteria																
ESTAR(1,2)	0.059	0.000	0.114	0.000	0.176	0.000	0.212	0.000	0.082	0.000	0.126	0.000	0.143	0.000	0.187	0.000
ESTAR(2,1)	0.017	0.000	0.043	0.000	0.065	0.000	0.076	0.000	0.031	0.000	0.049	0.000	0.054	0.000	0.068	0.000
ESTAR(2,2)	0.247	0.059	0.285	0.074	0.266	0.089	0.230	0.092	0.279	0.060	0.292	0.075	0.284	0.087	0.261	0.089
LSTAR(1,1)	0.004	0.000	0.014	0.000	0.052	0.000	0.086	0.000	0.009	0.000	0.022	0.000	0.028	0.000	0.060	0.000
LSTAR(2,1)	0.014	0.000	0.031	0.000	0.056	0.000	0.070	0.000	0.020	0.000	0.039	0.000	0.043	0.000	0.062	0.000
LSTAR(2,2)	<b>0.263</b>	<b>0.605</b>	<b>0.323</b>	<b>0.755</b>	<b>0.309</b>	<b>0.897</b>	<b>0.281</b>	<b>0.902</b>	<b>0.292</b>	<b>0.611</b>	<b>0.329</b>	<b>0.758</b>	<b>0.321</b>	<b>0.882</b>	<b>0.304</b>	<b>0.895</b>

Bold numbers mark the models with the highest selection frequency for the respective criteria.

**Table 4:** Selection frequencies of the overall information criteria: ESTAR vs. LSTAR models for LSTAR(2,2) DGP with  $\gamma=20$ .

$n$	AIC( $\hat{\sigma}$ )		AIC( $\tilde{\sigma}$ )		SIC( $\hat{\sigma}$ )		SIC( $\tilde{\sigma}$ )		AICc( $\hat{\sigma}$ )		AICc( $\tilde{\sigma}$ )		WIC( $\hat{\sigma}$ )		WIC( $\tilde{\sigma}$ )	
	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000
ESTAR(1,2)	0.157	0.000	0.257	0.001	<b>0.374</b>	0.003	<b>0.439</b>	0.004	0.190	0.000	0.292	0.001	0.327	0.003	<b>0.385</b>	0.003
ESTAR(2,1)	0.005	0.000	0.008	0.000	0.011	0.000	0.013	0.000	0.007	0.000	0.008	0.000	0.008	0.000	0.012	0.000
ESTAR(2,2)	<b>0.304</b>	0.228	<b>0.349</b>	0.287	0.318	0.359	0.275	0.366	<b>0.355</b>	0.233	<b>0.351</b>	0.290	<b>0.339</b>	0.354	0.311	0.359
LSTAR(1,1)	0.000	0.000	0.004	0.000	0.017	0.000	0.030	0.000	0.000	0.000	0.004	0.000	0.004	0.000	0.019	0.000
LSTAR(2,1)	0.002	0.000	0.004	0.000	0.014	0.000	0.017	0.000	0.002	0.000	0.009	0.000	0.012	0.000	0.017	0.000
LSTAR(2,2)	0.176	<b>0.425</b>	0.205	<b>0.518</b>	0.201	<b>0.614</b>	0.193	<b>0.618</b>	0.192	<b>0.429</b>	0.208	<b>0.519</b>	0.207	<b>0.607</b>	0.199	<b>0.614</b>

Bold numbers mark the models with the highest selection frequency for the respective criteria.

criteria as well as the regime weighted SIC and WIC select the SETAR model. This might be due to the fact that SETAR models have a smaller penalty term since the transition parameter  $\gamma$  does not have to be estimated. Hence, this spurious behavior can be interpreted as some kind of underfitting. The regime weighted versions of AIC and AICc prefer LSTAR models, even if the true DGP is a SETAR model. In larger samples the information criteria select the true model class. Again, the transition parameter  $\gamma$  is a key factor of the performance. For large values of  $\gamma$  the LSTAR model becomes a SETAR model. Then, it is impossible to distinguish these two model classes. The information criteria tend to favor the SETAR model in these cases, since the penalty term is smaller due to the missing transition parameter. For smaller values of  $\gamma$  the performance of the information criteria improves and especially for larger samples the differentiation between the different types of transition functions is reliable. Although in some cases a wrong model class is selected, the lag order coincides with the one preferred if only the correct model class is considered.

Finally, linear models are included into the collection of models. Then, the problem of distinguishing between linear and nonlinear models becomes relevant again. As already shown, weighted versions of AIC and AICc cannot detect linearity. SIC and WIC tend to favor linear models although the model is nonlinear. This is especially a problem in small samples. The equally weighted and the overall versions of the information criteria often select the SETAR models. However, as already mentioned earlier, the performance improves with an increasing sample size. Regime weighted AIC and AICc spuriously select LSTAR models. Lag order selection works well. Even if the correct model class is not selected, the lag order combination corresponds to the artificial case with the correct model class.

In order to correctly determine the form of the transition function, equally weighted and overall SIC and WIC should be applied. In small samples also  $AIC(\hat{\sigma}^2)$  and  $AICc(\hat{\sigma}^2)$  can be used. However, it should be kept in mind that the results are not very reliable in small samples.

### 4.3.7 Discriminating regimes

Finally, we evaluate whether the information criteria can also be used to select the number of regimes. In this paragraph we only consider SETAR models with two or three regimes (cf. Clements and Krolzig 1998; Gonzalo and Pitarakis 2002; Hamaker 2009) and reduce the maximum lag order to two, which yields eight lag order combinations for the three regime model. Furthermore, we also change the grid for thresholds. In this simulation study the grid consists of all quantiles from the 10% to the 90% quantile of  $y_t$ . Although this grid might be rough, we choose it because it guarantees that even if two consecutive grid points are chosen as thresholds, there will lie 10% of the observations in the middle regime. The estimation procedure of the thresholds is the sequential approach from Gonzalo and Pitarakis (2002). We consider the 1-step and the 2-step approach, i.e. we estimate the first threshold and given this threshold the second threshold. In the 1-step approach we keep these two estimates. In contrast, in the 2-step approach the first threshold is reestimated given the second threshold and finally the second threshold is reestimated given the refined first threshold (cf. Gonzalo and Pitarakis 2002). For the three regime SETAR models the number of additional parameters increases to three (two thresholds, one delay parameter). Gonzalo and Pitarakis (2002) do not account for the number of thresholds in the penalty term, whereas Liu, Wu, and Zidek (1997) incorporate thresholds as additional parameters into the penalty term. Hamaker (2009) uses one criterion which accounts for additional parameters.

Table 5 presents the selection frequencies of the information criteria if the true model is an AR(2) process and the 1-step approach is applied. The corresponding results for the 2-step approach can be found in the Supplementary Material. The weighted versions of the 1-step approach lead to more correct selections of the AR(2) model than the weighted criteria of the 2-step approach. The weighted versions of AIC and AICc never select the true model, but versions of SIC and WIC detect the linearity. The corresponding information criteria of the 2-step algorithm prefer a three regime SETAR model. Both approaches work well if the overall information criteria are applied. Then, the 2-step approach is even slightly superior. For the 1-step approach the selection frequencies of the true model vary between 0.492 for the overall  $AIC(\hat{\sigma})$  and 0.826 for the overall  $WIC(\hat{\sigma})$  in the small sample. In the large sample the selection frequencies further increase to 0.695 for the overall  $AIC(\hat{\sigma})$  and 1 for overall SIC and WIC.

In Table 6 the results for the SETAR(1,1) process are tabulated. In small samples weighted versions of SIC and WIC and the overall information criteria of the 1-step approach favor the linear models, whereas the respective versions of AIC and AICc favor the SETAR(1,1,1). For the 2-step approach equally weighted and overall SIC and WIC select the linear model. The other criteria prefer the SETAR(1,1,1). With an increasing sample size the weighted versions of SIC and WIC and all overall information criteria select the correct model when the 1-step approach is applied. For the 2-step approach only the overall criteria favor the correct model. The other versions prefer the SETAR(1,1,1).

**Table 5:** Selection frequencies of the equally weighted information criteria: AR vs. SETAR(2; ·, ·) vs. SETAR(3; ·, ·, ·) models for AR(2) DGP (1-step Estimation).

<i>n</i>	AIC( $\hat{\sigma}$ )		AIC( $\hat{\sigma}$ )		SIC( $\hat{\sigma}$ )		SIC( $\hat{\sigma}$ )		AICc( $\hat{\sigma}$ )		AICc( $\hat{\sigma}$ )		WIC( $\hat{\sigma}$ )		WIC( $\hat{\sigma}$ )	
	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000
AR(2)	0.000	0.000	0.000	0.000	<b>0.470</b>	<b>0.637</b>	<b>0.410</b>	<b>0.637</b>	0.000	0.000	0.000	0.000	<b>0.547</b>	<b>0.609</b>	<b>0.500</b>	<b>0.607</b>
SETAR(2,2)	0.000	0.000	0.000	0.000	0.043	0.036	0.027	0.030	0.001	0.000	0.001	0.000	0.057	0.042	0.043	0.037
SETAR(1,2,1)	0.104	0.002	0.165	0.004	0.051	0.011	0.057	0.018	<b>0.191</b>	0.002	<b>0.208</b>	0.005	0.018	0.007	0.020	0.011
SETAR(1,2,2)	0.194	0.034	<b>0.170</b>	0.059	0.031	0.053	0.030	0.051	0.179	0.035	0.157	0.060	0.014	0.043	0.014	0.057
SETAR(2,2,2)	<b>0.236</b>	<b>0.904</b>	0.146	<b>0.841</b>	0.034	0.158	0.019	0.137	0.100	<b>0.900</b>	0.074	<b>0.844</b>	0.017	0.220	0.012	0.187

Bold numbers mark the models with the highest selection frequency for the respective criteria.

**Table 6:** Selection frequencies of the information criteria: AR vs. SETAR(2; ·, ·) vs. SETAR(3; ·, ·, ·) models for SETAR(1,1) DGP (1-step Estimation).

$n$	AIC( $\hat{\sigma}$ )		AIC( $\tilde{\sigma}$ )		SIC( $\hat{\sigma}$ )		SIC( $\tilde{\sigma}$ )		AICc( $\hat{\sigma}$ )		AICc( $\tilde{\sigma}$ )		WIC( $\hat{\sigma}$ )		WIC( $\tilde{\sigma}$ )	
	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000	100	1000
(a) Equally weighted information criteria																
AR(1)	0.000	0.000	0.000	0.000	<b>0.504</b>	0.110	<b>0.532</b>	0.110	0.000	0.000	0.000	0.000	<b>0.468</b>	0.090	<b>0.489</b>	0.091
SETAR(1,1)	0.000	0.000	0.000	0.000	0.117	<b>0.524</b>	0.145	<b>0.536</b>	0.003	0.000	0.004	0.000	0.245	<b>0.506</b>	0.272	<b>0.524</b>
SETAR(1,1,1)	<b>0.327</b>	<b>0.469</b>	<b>0.500</b>	<b>0.634</b>	0.119	0.261	0.159	0.282	<b>0.501</b>	<b>0.504</b>	<b>0.567</b>	<b>0.650</b>	0.062	0.256	0.064	0.290
SETAR(1,2,1)	0.179	0.151	0.157	0.124	0.064	0.035	0.041	0.020	0.168	0.145	0.150	0.114	0.053	0.042	0.050	0.027
SETAR(2,1,1)	0.160	0.137	0.142	0.114	0.040	0.019	0.033	0.018	0.129	0.138	0.121	0.115	0.016	0.024	0.018	0.019
(b) Overall information criteria																
AR(1)	0.147	0.000	<b>0.296</b>	0.000	<b>0.842</b>	0.016	<b>0.879</b>	0.021	<b>0.222</b>	0.000	<b>0.358</b>	0.000	<b>0.686</b>	0.003	<b>0.784</b>	0.005
SETAR(1,1)	0.167	<b>0.306</b>	0.288	<b>0.535</b>	0.096	<b>0.953</b>	0.090	<b>0.967</b>	0.217	<b>0.321</b>	0.315	<b>0.546</b>	0.171	<b>0.940</b>	0.147	<b>0.963</b>
SETAR(1,1,1)	<b>0.200</b>	0.254	0.158	0.206	0.012	0.003	0.005	0.001	0.174	0.252	0.117	0.200	0.030	0.008	0.018	0.004
SETAR(1,2,1)	0.094	0.094	0.046	0.044	0.001	0.000	0.000	0.000	0.070	0.090	0.029	0.043	0.007	0.000	0.002	0.000
SETAR(2,1,1)	0.088	0.079	0.036	0.041	0.002	0.000	0.001	0.000	0.057	0.077	0.024	0.040	0.007	0.000	0.004	0.000

Bold numbers mark the models with the highest selection frequency for the respective criteria.

In the Supplementary Material the results for the SETAR(2,2,2) model can be found. Even in a small sample all versions of information criteria identify the correct model when the 1-step algorithm is applied. In large samples the selection frequencies converge towards 1. In case of the 2-step approach the weighted versions select the correct model class though not the correct lag order combination when the sample size is small. For a larger sample the correct model is selected. The overall information criteria favor a two regime SETAR model in small samples and an underfitted three regime SETAR model in the larger sample.

Taking into account all these results, we cannot find one superior approach. In general, we would expect the 2-step algorithm to outperform the 1-step approach. However, this is not always case, which is probably due to the grid. Li (1988) points out that the choice of the grid has an effect on the threshold estimates. Thus, a finer grid might improve the results, but the computational effort would increase enormously. This problem may be solved by applying a 2-step grid search algorithm. In the first step a rough grid is used to find a first threshold estimate. In the second step, a finer grid is built around this point and the estimate is refined. As a result, it would be unnecessary to do a global fine grid search. Instead, the computational effort would only increase locally.

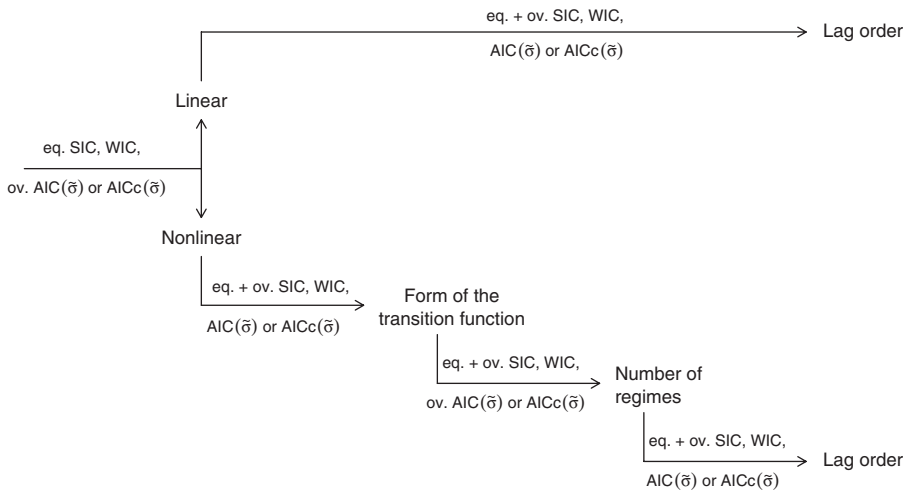
Generally, information criteria can be used to distinguish between SETAR models with a different number of regimes. In particular the overall criteria perform well but in small samples there is the possibility of underfitting, especially for SIC and WIC. Therefore, we recommend to apply overall AIC( $\tilde{\sigma}^2$ ) and AICc( $\tilde{\sigma}^2$ ) or equally weighted SIC and WIC in small samples. In large samples equally weighted and overall SIC and WIC should be used.

However, since the computational effort of estimating nonlinear models with more than two regimes is rather high, Gonzalo and Pitarakis (2002) recommend first to use a collection of models consisting only of linear and two regime models. If the information criteria select the nonlinear model, in a second step the two regime models and the three regime models are evaluated. If in contrast the linear model is preferred, it is not necessary to estimate the three regime models. Thus, model selection can be understood as an iterative process. Figure 5 illustrates the steps of the model selection approach and the respective information criteria which should be applied.

## 5 Conclusion

In this paper we evaluate the performance of different information criteria for simultaneous lag order and model class selection of nonlinear models. We focus on SETAR and STAR models due to the fact that they





**Figure 5:** The model selection approach.

have a similar switching mechanism and are frequently applied. Our set of information criteria consists of the commonly known criteria AIC, SIC and AICc. Furthermore, we also apply WIC which is supposed to perform well independent of the sample size. All in all, we consider 24 different information criteria with varying penalty terms, error term variance estimators and regime weightings. Our aim is to identify one or more criteria that can be used to select a best fitting model among different nonlinear model classes. Strictly speaking, information criteria cannot be employed in order to select between different model classes, because they are developed under the assumption that all models under consideration belong to the same parametric family (cf. Kapetanios 2001). Nevertheless, this approach can be a valuable alternative to linearity tests since it does not depend on a (possibly misspecified) lag order estimate under the assumption of linearity and will lead to a definite model choice.

Our results show that the information criteria perform well in general. However, there are some key factors that seriously influence the performance of the criteria: the sample size, the identifiability of the true model, the distance between regimes, and the shape of the transition functions. In small samples some criteria adopt a spurious behavior. Especially overall SIC and WIC tend to select a simple model. Depending on the collection of models, this results in underfitting, the selection of linear models, although the true model is nonlinear, or a model with fewer additional parameters. In large samples the information criteria perform well. The identifiability of the true model influences the lag order selection. If the largest lag has only a minor influence, the model is weakly identifiable (cf. McQuarrie, Shumway, and Tsai 1997) and underfitting occurs. Hence, underfitting is not necessarily a drawback of the information criteria but is due to the true DGP. The distance between regimes concerns the identifiability as well. If the regimes are very similar, a simpler model will be preferred. Especially in small samples the selection frequencies of the correct model are quite low. Asymptotically there are only few incorrect selections. The distance of regimes plays an important role for the discrimination between linear and nonlinear models. We have shown that weighted versions of AIC and AICc cannot detect linearity and therefore, should not be applied, if linear models are among the collection of models, since they spuriously point to the nonlinear model (cf. also Gonzalo and Pitarakis 2002; Pitarakis 2006). Finally, the shape of the transition functions of STAR models depends on the transition parameter  $\gamma$ . For small values both STAR models converge to AR models, whereas for large values the LSTAR model converges to a SETAR model and the ESTAR model reduces to a linear model. In these extreme cases it is difficult to distinguish the different model types. Due to the smaller penalties AR and SETAR models will be preferred to STAR models. Again a large sample size facilitates the selection process. For a larger number of observations, the differentiation of the type of transition function becomes more reliable (cf. also Psaradakis et al. 2009). The selection of the lag order combination does not suffer from misspecified model classes. Instead, it corresponds to the one selected if only the true model class is considered.

**Table 7:** Application of information criteria in small samples.

To determine	Information criteria
Lag orders	Equally weighted and overall SIC, WIC, AIC( $\hat{\sigma}^2$ ) and AICc( $\hat{\sigma}^2$ )
(Non-)Linearity	Equally weighted SIC and WIC, overall AIC( $\hat{\sigma}^2$ ) and AICc( $\hat{\sigma}^2$ )
Transition function	Equally weighted and overall SIC, WIC, AIC( $\hat{\sigma}^2$ ) and AICc( $\hat{\sigma}^2$ )
Number of regimes	Overall AIC( $\hat{\sigma}^2$ ), AICc( $\hat{\sigma}^2$ ) and equally weighted SIC, WIC

To conclude, we have shown that the regime weighted information criteria are outperformed by their equally weighted and overall counterparts independent of the dominance of the regimes and thus should not be applied. If the collection of models includes linear models, the weighted versions of AIC and AICc will spuriously point to the nonlinear model. So, other criteria should be preferred. Finally, the probability of overfitting of AIC and AICc can be reduced by employing the unbiased variance estimator.

Table 7 briefly summarizes which information criteria should be employed in small samples depending on the aim of application. In large samples the equally weighted and overall SIC and WIC perform well in all application areas.

It should be kept in mind that the discrimination of the transition function is not very reliable in small samples (cf. also Psaradakis et al. 2009).

Although there are some criteria performing better in certain situations, we cannot generally advise the use of one single information criterion. Instead, several criteria should be applied in order to balance the individual strengths and weaknesses, like the tendency to over- or underfit. However, our results show that the performance of the information criteria is not deteriorated if model specific parameters like threshold, delay or transition parameters are unknown and have to be estimated. Hence, the application of information criteria to empirical data in order to identify the best fitting model is an alternative approach to linearity tests.

**Acknowledgments:** The authors are grateful to two anonymous referees, the associate editor and the participants of the Statistische Woche in Hannover for helpful comments and suggestions which improved the paper. Financial support by the Deutsche Forschungsgemeinschaft (<http://www.dfg.de/index.jspDFG>) under grant SI 745/9-1 is gratefully acknowledged.

## References

- Akaike, H. 1974. "A New Look at The Statistical Model Identification." *IEEE Transactions on Automatic Control* 19: 716–723.
- Bollerslev, T. 1986. "Generalized Autoregressive Conditional Heteroskedasticity." *Journal of Econometrics* 31: 307–327.
- Clements, M. P., and H.-M. Krolzig. 1998. "A Comparison of The Forecast Performance of Markov-Switching and Threshold Autoregressive Models of US GNP." *The Econometrics Journal* 1: 47–75.
- Emiliano, P. C., M. J. Vivanco, and F. S. De Menezes. 2014. "Information Criteria: How Do They Behave in different Models?" *Computational Statistics & Data Analysis* 69: 141–153.
- Engle, R. F. 1982. "Autoregressive Conditional Heteroscedasticity with Estimates of the Variance of United Kingdom Inflation." *Econometrica* 50 (4): 987–1007.
- Gonzalo, J., and J.-Y. Pitarakis. 2002. "Estimation and Model Selection Based Inference in Single and Multiple Threshold Models." *Journal of Econometrics* 110: 319–352.
- Hamaker, E. 2009. Using Information Criteria to Determine The Number of Regimes in Threshold Autoregressive Models." *Journal of Mathematical Psychology* 53: 518–529.
- Hamilton, J. D. 1989. "A New Approach To The Economic Analysis of Nonstationary Time Series and The Business Cycle." *Econometrica* 57: 357–384.
- Hamilton, J. D. 1994. *Time Series Analysis*. Princeton University Press.
- Hansen, B. E. 1997. "Inference in TAR Models." *Studies in Nonlinear Dynamics & Econometrics* 2.
- Hughes, A. W., and M. L. King. 2003. "Model Selection Using AIC in the Presence of One-Sided Information." *Journal of Statistical Planning and Inference* 115: 397–411.
- Hughes, A. W., M. L. King, and K. T. Kwek. 2004. "Selecting the Order of An ARCH model." *Economics Letters* 83: 269–275.
- Hurvich, C. M., and C.-L. Tsai. 1989. "Regression and Time Series Model Selection in Small Samples." *Biometrika* 76: 297–307.

- Kapetanios, G. 2001. "Model Selection in Threshold Models." *Journal of Time Series Analysis* 22: 733–754.
- Li, W. 1988. "The Akaike Information Criterion in Threshold Modelling: Some Empirical Evidences." In *Nonlinear Time Series and Signal Processing, Lecture Notes in Control and Information Sciences*, vol. 106, Berlin Heidelberg: Springer, pp. 88–96, URL <http://dx.doi.org/10.1007/BFb0044277>.
- Liu, J., S. Wu, and J. V. Zidek. 1997. "On Segmented Multivariate Regression," *Statistica Sinica* 7: 497–525.
- Luukkonen, R., P. Saikkonen, and T. Teräsvirta. 1988a. "Testing Linearity Against Smooth Transition Autoregressive Models." *Biometrika* 75: 491–499.
- Luukkonen, R., P. Saikkonen, and T. Teräsvirta. 1988b. "Testing Linearity in Univariate Time Series Models." *Scandinavian Journal of Statistics* 15: 161–175.
- McQuarrie, A., R. Shumway, and C.-L. Tsai. 1997. "The Model Selection Criterion AICu." *Statistics & Probability Letters* 34: 285–292.
- Pitarakis, J.-Y. 2006. "Model Selection Uncertainty and Detection of Threshold Effects." *Studies in Nonlinear Dynamics & Econometrics* 10.
- Psaradakis, Z., M. Sola, F. Spagnolo, and N. Spagnolo. 2009. "Selecting Nonlinear Time Series Models Using Information Criteria." *Journal of Time Series Analysis* 30: 369–394.
- Schwarz, G. 1978. "Estimating the Dimension of A Model." *The Annals of Statistics* 6: 461–464.
- Shibata, R. 1986. "Consistency of Model Selection and Parameter Estimation." *Journal of Applied Probability* 23: 127–141.
- Smith, A., P. A. Naik, and C.-L. Tsai. 2006. "Markov-switching Model Selection using Kullback–Leibler Divergence." *Journal of Econometrics* 134: 553–577.
- Teräsvirta, T. 1994. "Specification, Estimation, and Evaluation of Smooth Transition Autoregressive Models." *Journal of the American Statistical Association* 89: 208–218.
- Tong, H. 1983. *Threshold Models in Non-Linear Time Series Analysis. Lecture Notes in Statistics, No. 21*. Springer-Verlag.
- Tong, H. 1990. *Non-linear Time Series: A Dynamical System Approach*. Oxford University Press.
- Tong, H., and K. S. Lim. 1980. "Threshold Autoregression, Limit Cycles and Cyclical Data." *Journal of the Royal Statistical Society. Series B (Methodological)* 245–292.
- van Dijk, D., T. Teräsvirta, and P. H. Franses. 2002. "Smooth Transition Autoregressive Models – A Survey of Recent Developments." *Econometric Reviews* 21: 1–47.
- Wen, M.-J., and Y.-H. Tu. 2001. "Modified WIC for Order Selection in Autoregressive Model." Technical Report No. 40, National Cheng-Kung University, Institute of Statistics.
- Wong, C. S., and W. K. Li. 1998. "A Note on the Corrected Akaike Information Criterion for Threshold Autoregressive Models." *Journal of Time Series Analysis* 19: 113–124.
- Wu, T.-J., and A. Sepulveda. 1998. "The Weighted Average Information Criterion for Order Selection in Time Series and Regression Models." *Statistics & Probability Letters* 39: 1–10.

---

**Supplemental Material:** The online version of this article (DOI: 10.1515/snde-2015-0026) offers supplementary material, available to authorized users.