

Parameter Estimation of Switching Piecewise Linear System

José Ragot, Gilles Mourot and Didier Maquin

Centre de Recherche en Automatique de Nancy, UMR CNRS 7039

Institut National Polytechnique de Nancy

2, avenue de la forêt de Haye, F 54 516 Vandoeuvre les Nancy

{Jose.Ragot,Gilles.Mourot,Didier.Maquin}@ensem.inpl-nancy.fr

Abstract—During the last years, a number of methodological papers on models with discrete parameter shifts have revived interest in the so-called regime switching models. Piecewise linear models are attractive when modelling a wide range of nonlinear system and determining simultaneously i) the data partition ii) the time instant of change iii) the parameter values of the different local models. This is a difficult problem for which no solution exists in the general case and we show here some aspects and particular results concerning the problem of off line learning of switching time series. We propose a method for identifying the parameters of the local models when choosing an adapted weighting function, this function allowing to select the data for which each local model is active. Indeed the proposed method is able to solve simultaneously the data allocation and the parameter estimation. The feasibility and the performance of the procedure is demonstrated using several academic examples.

I. INTRODUCTION

For the identification of nonlinear systems, there has been a large activity during the past years. In particular many interesting results have been reported in connection with multi-model [1] and/or multiple models [2], hinging hyperplanes [3], [4], hidden Markov models [5], mixture of regressions [6], segmented curves [7]. Most of these works refer to quasistationary or locally stationary systems characterised by abrupt changes between stationary segments with different statistical properties. Many formulations of this problem also appear in the field of fuzzy systems [8], [9].

In the following, we focus the attention on PieceWise Auto Regressive eXogeneous models (PWARX). As it will be pointed out latter, if the partition of piecewise mapping is known, the problem of identification can easily be solved by using standard techniques of estimation. However, when the partition is unknown the problem becomes much more difficult. Thus, there are two possibilities. Either a partitioning, defining the local domains in which the system is constant, is a priori defined or the partitioning has to be estimated along with the local models.

Our contribution is to illustrate this problem in the case where the structure and the number of the local models are known. Thus, we restrict the estimation problem to i) the estimation of switching between the local models, ii) the estimation of the parameters of the local models. Summarising, the main ideas of our contribution deal with

the use of adapted weights allowing a powerful classification of the data and a sequential estimation of the different local model parameters.

This paper is organised as follows : section 2 explains, through a simple example, what is the problem to solve and the foregoing difficulties. Section 3 constitutes the contribution of the paper and is followed by a conclusion. Some simulation examples provide an illustration of the proposed algorithms in section 4.

II. MODEL DESCRIPTION

To begin with, let us consider systems in regression form

$$y_k = \varphi_k^T \theta_j, \text{ if } H_j^T \varphi_k \leq 0, j = 1..s \quad (1)$$

where $\varphi_k \in \mathbb{R}^p$ is a regression vector depending on the input u and the output y of the system, $\theta_j \in \mathbb{R}^p$ the parameter vector associated with the j th local model and $H_j \in \mathbb{R}^p$ are unknown parameters. We then consider that the observations are generated by switching among s different AR or ARX models of orders p and parameters θ_j . Further, we will use also the notation:

$$y_k = \varphi_k^T \theta(\nu) \quad (2)$$

where ν is a key vector describing in what mode the system is for the time being ; ν can be a function of (k, u, y) or some external input and takes its values in a finite set $I_s = \{1, \dots, s\}$. Thus the time series is generated by the combination of s functions $\varphi_k^T \theta(\nu)$.

This is not the only way to describe switching system and the reader should refer to [4], [10], [11], [12], [13], [14] for other formulations using mixture of models, endogenous switching, structural break models, self exciting threshold autoregressions (SETAR), model smooth transition autoregressive model (STAR), neural network [15] and, at last, hybrid systems [16].

The regression vector φ could consist of old inputs and outputs. The sets $Z_j = \{H_j^T \varphi_k \leq 0\}, j = 1..s$ are polyhedral partitions of the φ -space. Our problem, when we are given y_k and $\varphi_k, k = 1..N$, consists in finding the PWARX model that best matches the given data, the number s being generally unknown. The model (1) can be identified by minimising the

optimisation criterion :

$$\Phi = \sum_{j=1}^s \sum_{k=1}^N (y_k - \varphi_k^T \theta_j)^2 \rho_j(\varphi_k) \quad (3)$$

subject to :

$$\rho_j(\varphi_k) = \begin{cases} 1 & \text{if } H_j^T \varphi_k \leq 0 \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

where θ_j and H_j , $j = 1..s$, are unknown.

An important problem is that of change-point detection, namely, detecting when the time-series has switched in some manner (1). At time k , for a given regression vector φ_k , only one condition is satisfied and exactly one of the functions $\rho_j(\varphi_k)$ equals 1. This means that all the terms in (3) are zero excepted one, those corresponding to the active local model. Consequently, if we known that a set of data is belonging to a particular local model, minimising (3) is straightforward. However, in general, the data partitioning is a priori unknown and the parameter estimation problem becomes difficult.

In our case, we limit the estimation problem to the one of θ_j ; however, we need to simultaneously estimate the values of the function φ_k in order to know the time switching i.e. the data useful to estimate the parameters of the j th local model. In fact we are not involved with the explanation of the switching, i.e. the estimation of the H_j parameters.

III. THE MAIN ALGORITHM

Here we present our main contribution. Let y_k represents the output measurements of the underlying system and $y_{k,j}$ the output of the j th local model. To fit the local model to the data, we attempt to minimise the error function:

$$\Phi = \sum_{j=1}^s \sum_{k=1}^N (y_k - y_{k,j})^2 p_{k,j} \quad (5)$$

$$y_{k,j} = \varphi_k^T \theta_j$$

where the weights $p_{k,j}$ have to be designed such that the local model j is adapted only with the input-output data for which it is concerned. It can be seen that the cost function (5) represents a trade-off between local and global learning. Indeed, when the model output $y_{k,j}$ is closed to the measurement y_k then model j matches the measurements and $p_{k,j}$ must be smaller than $p_{k,l}$, $\forall l \neq j$. In general, this performance index has to be minimised with respect to the parameter vectors θ_j for all possible disjoint partitions of the measurement set and for all possible numbers of submodels. Here we restrict the identification problem, the number of submodels being a priori chosen. Obviously, the key point is the design of these weights. In the following a non parametric estimation is used because there is no need to parameterize the weighting functions, only their values being useful to separate the data according to the s local models. The ideal

situation deals with the knowledge of the partition of the data into s groups, the first one gathering the data in accordance with the first model, and similarly for the other groups. These s sets are noted S_j , $j = 1..s$:

$$S_j = \{(x_k, y_k), k = 1..N / (x_k, y_k) \text{ satisfy model } j\}$$

Thus, the optimal weights are defined by:

$$p_{k,j} = \begin{cases} 1 & \text{if } (x_k, y_k) \in S_j \\ 0 & \text{if } (x_k, y_k) \notin S_j \end{cases} \quad k = 1..N, j = 1..s \quad (6)$$

In fact our algorithm tries to adapt the weights as closed as possible to the optimal ones.

The complete iterative algorithm is now described. Each iteration consists of two steps. The first one (step 1) is to determine an estimation of the weighting functions $p_{k,j}$ given the local models. The second step (step 2) is to identify the local models given the weights. Note that in [17] a similar algorithm is used in the context of weighted combination of local linear state-space systems but using a different approach based on an extended Kalman smoother allowing to estimate changes in the weights. It should be noted that the proposed algorithm estimates sequentially these local models (and use a serial data allocation). More precisely, the algorithm estimates the first local model, the second local model explains the residuals of the first local model and so on. The algorithm uses an adapted weighting function allowing the clustering of the data automatically.

Algorithm : sequential estimation

Step 0. Initialisation

Select s the number of local models

Select a set of weighting matrices W_j for the s local models.

Select a threshold ϵ for the convergence test.

Define the matrices:

$$X = (\varphi_1, \dots, \varphi_N)^T, y = (y_1, \dots, y_N), \Delta y_0 = y, \hat{\theta}_0 = 0$$

Set $r = 0$

Step 1. Parameter computation

For $j = 1..s$

$$\text{residual regression: } \Delta \hat{\theta}_j = (X^T W_j X)^{-1} X^T W_j \Delta y_{j-1}$$

$$\text{residual estimation: } \Delta \hat{y}_{j-1} = X \Delta \hat{\theta}_j$$

$$\text{residual: } \Delta y_j = \Delta y_{j-1} - \Delta \hat{y}_{j-1}$$

$$\text{local model parameters: } \hat{\theta}_j = \hat{\theta}_{j-1} + \Delta \hat{\theta}_j$$

$$\text{residual criterion: } \Phi_j = \|\Delta y_j\|_{W_j}^2$$

Step 2. Weight computation

For $j = 1..s$

$$\text{weights: } p_j = \prod_{q=1, q \neq j}^s (\Delta y_q)^{2r}$$

$$\text{normalised weights: } p_j = p_j / \prod_{j=1}^s p_j$$

$$W_j^{(r)} = \text{diag}(p_j)$$

The operator \prod is used for evaluating the Hadamard product of vectors. The $/$ operator allows to divide two

vectors component by component. The "diag" operator allows to construct a diagonal matrix from a vector.

Step 3. Convergence test

Check for termination in some convenient matrix form. If $\|W^{(r)} - W^{(r+1)}\| \leq \epsilon$ go to step 4, otherwise set $r = r + 1$ and return to step 1.

Step 4. Classification

The fuzzy data allocation is naturally given by the values of the weights. It is also possible to transform these weights into a binary representation involving only the values 0 and 1.

Remarks

- For the implementation issues, in step 2, the coefficient r enforced the weight and in our experience, it must be chosen between 2 and 4.
- The preceding algorithm supposes that matrices $X^T W_j X$ are regular. Indeed, it rarely occurs in practice that particular data and weights will cause singularity of $X^T W_j X$.
- The convergence of the algorithm is not discussed here and the reader may refer to [18] in which the use of EM and FCRM algorithms encounter the same convergence problem. As an evidence, the initialisation is a key point. It is important to note that, generally speaking, classification algorithms may terminate at extrema different from the true value. In our case, the proposed algorithm is quite enough insensitive to the initialisation used. However it is always possible to generate particular data for which the algorithm will trap at a local solution.

IV. EXAMPLES

To verify the validity of the proposed algorithm and test its performance, we conducted several Monte-Carlo experiments with simulated data most of them being collected from systems used as benchmark in the literature. Here some results are presented.

A. Example 1

The data have been generated by a PWARX model (this structure has been proposed and analysed in [19]):

$$y_{k+1} = au_k + b + e_k \quad (7)$$

$$\begin{cases} a = -1 & b = 0 & \text{if } u_k \in [-4, 0] \\ a = 1 & b = 0 & \text{if } u_k \in [0, 2] \\ a = 3 & b = -2 & \text{if } u_k \in [2, 4] \end{cases}$$

The input $u_k \in \mathbb{R}$ is a random sequence with uniform distribution on $[-4, 4]$ and the noise e_k is a random sequence with standard deviation $\sigma = 0.1$. In that example, the three clusters have respectively, 30, 15 and 15 data. Ideally, clusters 1, 2 and 3 (corresponding to models 1, 2 and 3) would respectively contain indices 1 to 30, indices 31 to 45 and indices 46 to 50. We have applied the proposed algorithm

to these data when the number of local models is fixed to 3 (which corresponds to the exact number of clusters in the data). The identified parameters are:

$$\text{Model} \begin{cases} 1 : & a = -0.999 & b = 0.011 \\ 2 : & a = 0.905 & b = 0.133 \\ 3 : & a = 2.917 & b = -1.754 \end{cases}$$

Figure 1 (left) presents the data of the simulated system in the plane $\{y_k, u_k\}$. Figure 1 (right) simultaneously displays the data and the estimated model for each class. The vicinity of the data in regard to the local models provides a good approximation of the PWARX system (7).

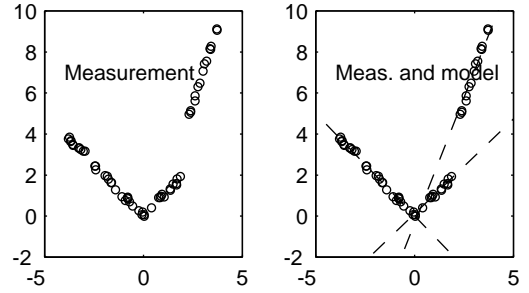


Fig. 1. The PWARX data sets in the plane y_{k+1}, u_k

As an example of classification performance, the left part of figure 2 shows the estimation errors while the right part presents the values of the weighting functions for local models 1, 2 and 3 at iteration 7 (after convergence of the algorithm). These weights allow to perform the classification of the data which is given here as membership coefficient normalised between 0 and 1 (it would be also possible to define a boolean classification).

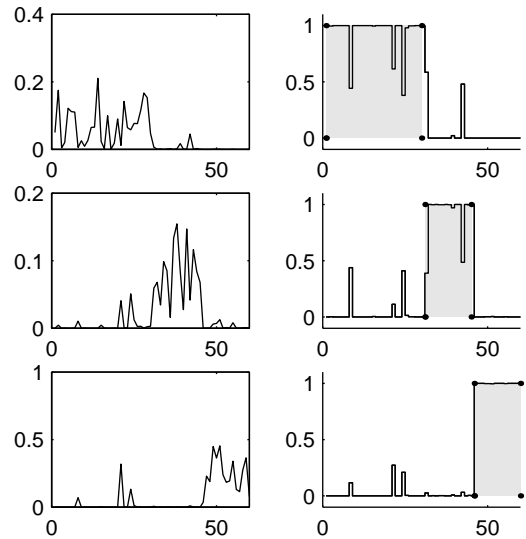


Fig. 2. Estimation errors and weights versus time

On the same figure, the true data allocation is indicated by shading areas ; in the computation, of course, the class labels of each data point are not known by the algorithm which "see" all the data simply as points $\{y_{k+1}, u_k\} \in \mathbb{R}^2$.

B. Example 2

The time series y_k is defined by [11]:

$$\begin{aligned} y_{k+1} &= au_k + be_k \\ y_k &= x_k + \epsilon_k \end{aligned} \quad (8)$$

$$\begin{cases} \text{if } |y_k| \leq 0.7 & a = 0.9 & b = 0.2 \\ \text{otherwise} & a = -0.9 & b = 0.3 \end{cases}$$

The random term e_k is a normally distributed white noise process with zero mean and variance 0.0625 ; the noise ϵ_k is also normally distributed with variance 0.002. Figure 3 shows the data and the result. Rows 1 to 3 present the output evolution and those of the parameter a and b only taking one of the two values -0.9 or 0.9 and 0.2 or 0.3 . Row 4 at left shows the data in the plane $\{y_k, y_{k-1}\}$ while the right part compare the data with the obtained model.

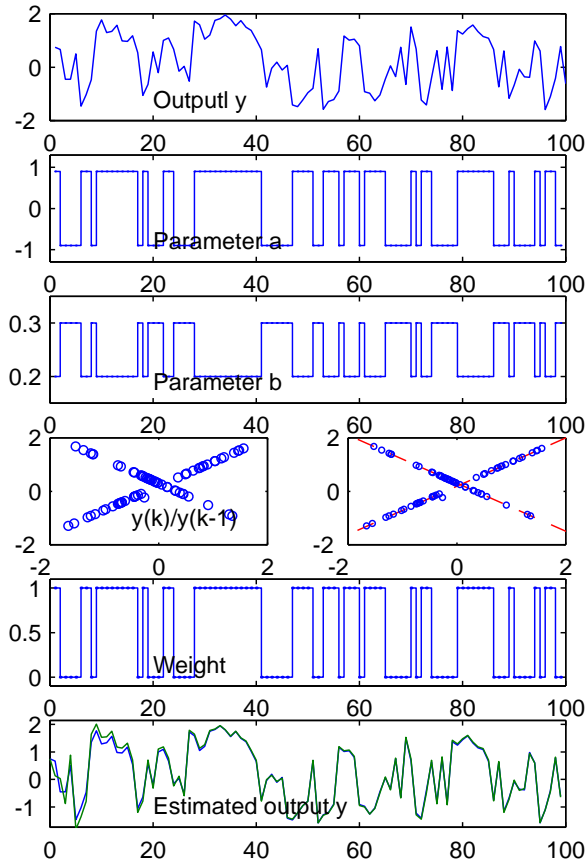


Fig. 3. Data and estimations

The normalised and rounded weights (taking only the values 0 and 1) given by the algorithm are drawn on row 5 and perfectly agree with the evolution of the a and b parameters ; these weights may be used to allocate the data to the two local models. The final clusters represented are correctly defined. Note that there is no hyperplan that separate the data sets (in the coordinate space $\{y_k, y_{k-1}\}$) because the clusters are not convex (due to the presence of an absolute value in the switching condition in (8)). In our approach, the clusters are defined through the weights that have been estimated simultaneously with the model parameters. Row 6 of figure 3 allows to compare the measured and the reconstructed output using the estimated parameters and time switching. Excepted in the vicinity of time origin for which bad initial conditions justify discrepancy, reconstructed state agrees with the true one.

C. Example 3: An hybrid tank system

The identification procedure is now applied to the simulation data of a tank system shown in figure 4. The input flow q_1 is time varying whereas the output flow q_2 linearly depends on the level h in the tank. The system's hybrid nature results from the interaction of the continuous dynamics and the discrete event dynamics and vice versa. The continuous dynamic depends on the liquid level in the tank while the dynamics switches if the levels rise above or fall beneath the height h_s for which the section of the tank is changing.

The model of the system is piecewise linear:

$$\begin{aligned} h_{k+1} &= h_k + S(\nu_k) (q_{1,k} - q_{2,k}) \\ q_{2,k} &= K h_k \end{aligned} \quad (9)$$

$$\nu(k) = \begin{cases} 0 & \text{if } h_k > h_s \\ 1 & \text{if } h_k \leq h_s \end{cases}, S(\nu_k) = S_1 + (S_2 - S_1)\nu_k$$

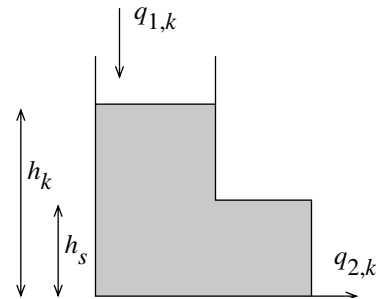


Fig. 4. An "hybrid" tank

Figure 5 gathers the data and the results. Rows 1 to 3 respectively indicate the input, the level and the output of the process while rows 4 and 5 respectively present the switching according to the section modification and the estimated switching; the estimated switching perfectly agrees with the true ones.

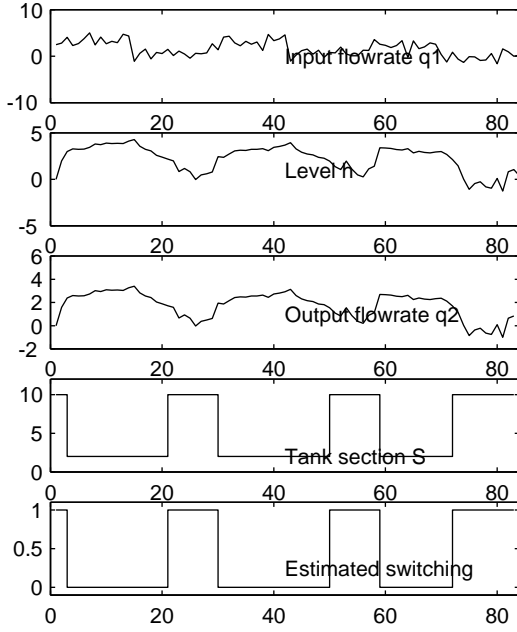


Fig. 5. Data and estimation versus time

D. Example 4: System with unknown a priori number of models

All the previous examples use an a priori knowledge about the number of local models. Here, the given example shows how a false hypothesis influences the identification results. The simulated system is described by a first order model taking three different sets of parameters. The measurements y_k are corrupted by a white noise process ϵ_k with zero mean and standard deviation 0.1:

$$\begin{aligned} x_{k+1} &= ax_k + bu_k \\ y_k &= x_k + \epsilon_k \end{aligned} \quad (10)$$

$$\begin{aligned} a &= 1 & b &= 2 \\ a &= 2 & b &= 8 \\ a &= -1 & b &= 3 \end{aligned}$$

The procedure is performed with 4 local models. Figure 6 shows the measurements (left) and the obtained models superposed to the measurements (right). The data classification is given in figure 7 which also shows the true allocation (grey boxes). The estimated parameters are collected in table I.

TABLE I
ESTIMATED PARAMETERS

| | model 1 | model 2 | model 3 | model 4 |
|---|---------|---------|---------|---------|
| a | 1.056 | 0.983 | 1.975 | -0.999 |
| b | 2.135 | 1.973 | 8.002 | 3.125 |

Analysing the results clearly points out that the structure of the model is not optimal. There are several ways to analyse

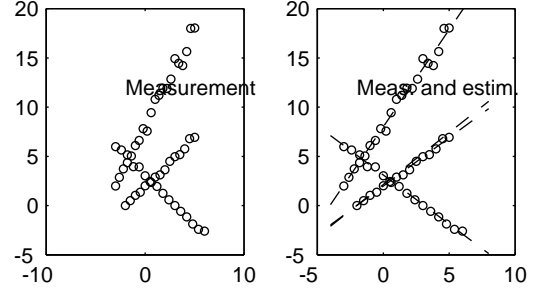


Fig. 6. Data and estimated model in plane $\{y_{k+1}/u_k, y_k/u_k\}$

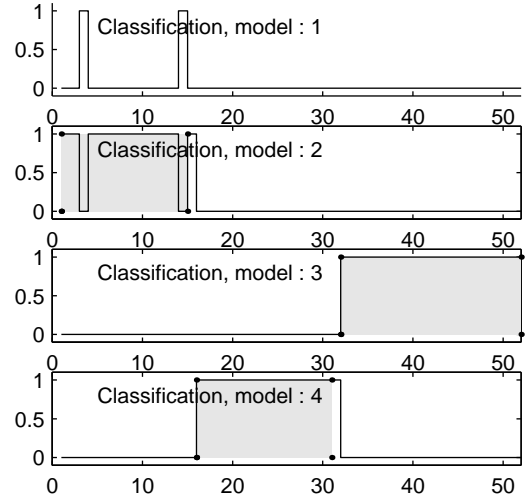


Fig. 7. Classification versus time

the problem of structure optimisation of the complete model [1], [20], [21]. The general idea is to detect the presence of neighbouring local models or the presence of local model with a very few number of associated data. In the present situation, according to the model coefficients in table I, two local models 1 and 2 have to be merged. For that purpose the merging may be performed by constraining models 1 and 2 to have the same behaviour, i.e. the same parameter vector. The new results are not presented, but they confirm the optimal structure with only three local models.

V. CONCLUSION

The proposed approach combines the identification of the parameters of a piecewise linear or affine form and the clustering of the data. This allows to identify both the affine local models and the partition of the domain in which each local model is valid ; in other words we have solved the data allocation task which consists in discovering that several local models exist and separated the data into groups corresponding to each model. We have successfully applied

the proposed approach to static and dynamic switching regressions.

Further investigations will firstly focus on the performances of the method namely in the case of noisy measurement. Secondly, the order selection of the local models together with their number needs to be further analysed and optimised. Thirdly, so far all the data sets we have deal with have a relatively low dimension, a large scale simulation will provide more insight into the robustness of this approach.

VI. REFERENCES

- [1] Gasso G., Mourot G. Ragot J. "Structure identification in multiple model representation : elimination and merging of local models", *40th IEEE Conference on Decision and Control*, Orlando, FL, USA, December 4-7, 2001.
- [2] Mihaylova L., Lampaert V., Bruyninckx H., Sweters J. "Hysteresis functions identification by multiple model approach", *IEEE Conference on Multi Sensor Fusion and Integration for Intelligent Systems, MFI'2001*, Baden Baden, Germany, August 20-21, 2001.
- [3] Pucar P., Sjoberg J. "On the hinge-finding algorithm for hinging hyperplanes, *IEEE Transactions on Information Theory*, 4 (3), p. 1310-1319, 1998.
- [4] Breiman L. Hinging hyperplans for regression, classification and function approximation, *IEEE Transactions on Information Theory*, 39 (3), p. 999-1013, 1993.
- [5] Ding Z., Hong L. An interactive multiple model algorithm with a switching markov chain, *Math. Comput. Modelling*, 25 (1), p. 1-9, 1997.
- [6] Quandt R.E. The estimation of the parameters of a linear regression system obeying two separate regimes, *Journal of the American Statistical Association*, p. 873-880, 1958.
- [7] Hudson D.J. Finding segmented curves whose join points have to be estimated, *Journal of the American Statistical Association*, 61, p. 1097-1129, 1966.
- [8] Kim E., Park M., Seunghwan L., Park M. A new approach to fuzzy modeling, *IEEE Transactions on Fuzzy Systems*, 5 (3), p. 328-337, 1997.
- [9] Yu J.R., Tzeng G.H., Li H.L. General fuzzy piecewise regression analysis with automatic change point detection, *Fuzzy Sets and Systems*, 119, p. 247-257, 2001.
- [10] Krolzig H.M. *Markov-switching vector autoregressive modellin, statistical inference and application to business and analysis*. Lecture Notes in Economics and Mathematical Systems, 454, Springer.
- [11] Medeiros M., Veiga A., Resenda M.G.C. A combinatorial approach to piecewise linear time series analysis, *Journal of Computational and Graphical Statistics*, 11 (1), p. 236-258, 2000.
- [12] Roll J. Robust verification and identification of piecewise affine systems. Thesis 899, Linköping, 2001.
- [13] Chua L.O., Kang S.M. "Section-wise piecewise linear functions : canonical representation, properties and applications", *Proceedings of IEEE*, 65, p. 915-929, 1977.
- [14] Fontaine L., Mourot G., Ragot J. "Segmentation d'électrocardiogrammes par réseau de modèles locaux", *Troisième Conférence Internationale sur l'Automatisation Industrielle*, Montréal, Canada, 7-9 juin 1999 (in French).
- [15] Kehagias A., Petridis V. Predictive modular neural networks for unsupervised segmentation of switching time series : the data allocation problem, *IEEE Transactions on Neural Networks*, 13 (6), p. 1432-1449, 2002.
- [16] Münz E., Krebs V. "Identification of hybrid systems using a priori knowledge", *15th IFAC World Congress*, Barcelona, Spain, July 21-26, 2002.
- [17] Verdult V., Verhaegen M. "Identification of a weighted combination of multivariable local linear state-space systems from input and output data", *40th IEEE Conference on Decision and Control*, Orlando, FL, USA, December 4-7, 2001.
- [18] Hathaway R.J., Bezdek C. Switching regression models and fuzzy clustering, *IEEE Transactions on Fuzzy Systems*, 1 (3), p. 195-204, 1993.
- [19] Ferrari-Trecate G., Muselli M., Liberati D., Morari M. "Identification of piecewise affine and hybrid systems", *American Control Conference*, Arlington, VA, USA, June 25-27, 2001.
- [20] Rao A.V., Miller J., Rose K., Gersho A. A deterministic annealing approach for parcimonious design of piecewise regression models, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21 (2), p. 159-173, 1999.
- [21] Strikholm B., Teräsvirta T. "Determining the number of regimes in a threshold autoregressive model using smooth transition autoregression", *13th EC2 Conference, Model Selection and Evaluation*, Bologna, Italy, December 13-14, 2002.