

OZONE CONCENTRATION MODELLING USING A MULTIPLE MODEL APPROACH

Komi Gasso, Gilles Mourot and José Ragot

Centre de Recherche en Automatique de Nancy - CNRS UPRESA 7039

Institut National Polytechnique de Lorraine - Ecole Nationale Supérieure de Géologie

2, Avenue de la Forêt de Haye - 54516 Vandoeuvre-lès-Nancy

Phone: 03 83 59 56 84 - Fax: 03 83 59 56 44 - Email: {ggasso, gmourot, jragot}@ensem.u-nancy.fr

Abstract: In this paper, the modelling of ozone concentration in an urban area is considered. The mechanisms of the phenomenon are non-linear, multivariable and time-varying. To solve the problem, we adopt a multiple model approach which is based on a decomposition principle. The main difficulty related to the identification of this structure concerns the partition of the feature space. The decomposition of the feature space along a decision tree is applied. The application to ozone concentration modelling is reported as well as the expected outcome.

Keywords: non-linear dynamic system; system identification; multiple model; k-d partition; ozone modelling.

1. INTRODUCTION

The problem under investigation in this paper is the modelling of ozone concentration using a multiple model approach. Ozone is a pollutant in the lower troposphere. It has detrimental effects on human health and on environment when its concentration reaches excessive values. In situations of high ozone level, appropriate decisions must be taken by authorities to inform the public and possibly to control the phenomenon. Therefore, the development of models is necessary to estimate ozone level in order to anticipate the decisions or to diagnose the abnormal operation of the sensors in air quality monitoring network.

Ozone is a secondary pollutant produced by complex photochemical reactions between nitrogen oxides (mainly NO and NO₂) and Volatile Organic Compounds emitted into the atmosphere. These reactions depend highly on the precursors emissions level and on the vertical and horizontal movements of the atmosphere that are linked to the meteorological conditions. By way of their interaction, these physical and chemical elements constitute a dynamic non-linear, multivariable and time-varying process. The theoretical models of ozone comprise the description of the physico-chemical mechanisms of ozone production and destruction and combine a large number of equations. Unfortunately, there is incomplete knowledge of the overall mechanisms. Furthermore, these models are computationally costly and they need measurements which are rarely available in air quality monitoring network. Therefore, they are hardly use in practice and black-box modelling of the ozone concentration has to be performed. To solve the problem, a multiple model approach is chosen because this representation can provide an accurate approximation of a non-linear system using a relative simple structure.

The basic principle of the multiple model approach to system modelling is the decomposition of the input space of the system into a set of small operating areas. Simple local models (usually linear models) are used to describe the system in each operating area. The global model is an interpolation of the local models using validity functions. These functions assess the relative validity of the local models and provide a smooth transition between the different operating ranges. The approximation abilities of this representation are well established (Johansen and Foss, 1993). Identify such a structure involves several tasks among which one of the most important is certainly the partition of the input space. Indeed, if little knowledge is available on the process, a grid partition is employed. But this leads to the curse of dimensionality, i.e. the number of partitions becomes overwhelmingly large if the number of inputs and/or the number of division per input increases. The structure thus obtained is overparameterised and this may result in ill-conditioned matrices that may affect the model performances. Therefore, a non-lattice partition is desirable.

In this paper, we describe a constructive technique for structure optimisation. The idea is to realise the input space decomposition using a decision tree. The parameters defining the validity functions are directly deduced and those of the local models are optimised in least squares sense. The paper is organised as follows: in section 2, the multiple modelling framework is described, as well as the inherent structure identification problem. Section 3 presents the input space partition algorithm and section 4 illustrates the application of the technique to model ozone concentration. Conclusions and perspectives of the present work are drawn in the last section.

2. MULTIPLE MODEL REPRESENTATION

2.1. MATHEMATICAL FORMULATION

Identification of a non-linear MISO dynamic system, on the basis of experimental data, is equivalent to the following function approximation problem:

$$\hat{y}(t) = F[\varphi(t), \theta] \quad (1)$$

where the unknown non-linear function $F[.]$ parameterised by the vector θ defines the model of the system. This function can be structured in many ways. The multiple model approach introduced by Johansen and Foss (1992) represents the function $F[.]$ as a local model network described by the equation:

$$\hat{y}(t) = \sum_{i=1}^M \omega_i(Z(t), \beta) f_i(\varphi(t), \theta_i) \quad (2)$$

Each local model or submodel f_i depends on a regression vector $\varphi(t)$ and on a vector of local parameters θ_i . The weighting function ω_i associated with each submodel acts as a validity measure of the corresponding submodel according to the regime under which the system is currently operating. These validity functions are defined over the feature space \mathbf{Z} (spanned by the vector $\mathbf{Z}(t) \in \mathfrak{R}^{n_z}$) and their number and position determine the partition of the domain \mathbf{Z} . The feature variables can include lagged inputs and output of the system or any auxiliary variable allowing the non-linearities of the process to be taken into account. Commonly, the validity functions are chosen as normalised gaussians

$$\omega_i(Z(t), \beta) = \frac{\rho_i(Z(t), \beta_i)}{\sum_{j=1}^M \rho_j(Z(t), \beta_j)} \quad (3)$$

with:

$$\rho_i(Z(t), \beta_i) = \prod_{h=1}^{n_z} \exp\left(-\frac{(Z_h(t) - c_{i,h})^2}{2\sigma_{i,h}^2}\right) = \exp\left(-\sum_{h=1}^{n_z} \frac{(Z_h(t) - c_{i,h})^2}{2\sigma_{i,h}^2}\right) \quad (4)$$

$$\beta_i = [c_{i,1}, \dots, c_{i,n_z}, \sigma_{i,1}, \dots, \sigma_{i,n_z}]^T \quad i = 1, \dots, M \quad (5)$$

$c_{i,h}$ et $\sigma_{i,h}$ are respectively the centre and the dispersion of the gaussian ρ_i along the axis h .

The local models are approximation of the system in their corresponding operating ranges. They could be of any form but for sake of simplicity, they are all chosen in this paper as dynamic linear model with orders p and q :

$$f_i(t) = \varphi^T(t) \theta_i \quad (6)$$

$$\varphi(t) = [y(t-1), \dots, y(t-p), u_1(t-1), \dots, u_1(t-q), \dots, u_m(t-1), \dots, u_m(t-q)]^T \quad (7)$$

$$\theta_i = \left[a_1^{(i)}, \dots, a_p^{(i)}, b_{1,1}^{(i)}, \dots, b_{1,q}^{(i)}, \dots, b_{m,1}^{(i)}, \dots, b_{m,q}^{(i)} \right]^T \quad i = 1, \dots, M \quad (8)$$

It can be noticed that the equations (2),(3),(4) and (6) form also the well-known Takagi-Sugeno model (Takagi and Sugeno, 1985) in which the membership functions are gaussians and the conjunction operator is product.

2.2. MULTIPLE MODEL IDENTIFICATION ISSUES

Building a multiple model involves the tasks of structure identification and parameter estimation. The former determines the local model orders and the validity function number and position. The later is related to the local model and validity function parameters tuning, given the structure.

2.2.1. Parameter estimation

Considering the vectors $\beta = [\beta_1^T, \beta_2^T, \dots, \beta_M^T]^T$ and $\theta = [\theta_1^T, \theta_2^T, \dots, \theta_M^T]$, parameter estimation consists to optimise β and θ in order to minimise the error between the target output $y(t)$ and the model $\hat{y}(t)$. To perform the optimisation, the following prediction error criterion is considered:

$$J = \sum_{t=1}^{N_A} \left(y(t) - \sum_{i=1}^M \omega_i(Z(t), \beta) \varphi^T(t) \theta_i \right)^2 \quad (9)$$

with N_A the number of training data.

Simultaneous optimisation of β and θ requires non-linear programming procedures like gradient descent method. The weak points of these techniques are the initialisation and the convergence (to local rather than global minimum) problems which are encountered, mainly in high-dimensional parametric space. As the criterion J is non-linear in β and quadratic with respect to θ , a two stages algorithm is used to alleviate these drawbacks (Tan et al., 1994). The vector θ is estimated by least squares for fixed value of β ; θ being known, the vector β is optimised using a non-linear optimisation method.

2.2.2. Structure optimisation

Besides the local model orders, the problem is to choose the feature variables $Z(t)$ and then to decompose the feature space Z into a number of small areas. Grid partition of the input space is unsuitable for high-dimensional problem due to the curse of dimensionality previously mentioned. To overcome the problem, many constraints must be imposed in order to limit the flexibility of the structure. This was performed by Mourot and Ragot (1997) who use a Takagi-Sugeno model with a grid partition to solve a similar ozone modelling problem. For instance, Mourot and Ragot (1997) consider no more than two modalities per premise variable and the widths of the membership function are equal whatever the modality of the considered variable. Despite these constraints, the resultant structure was refined by discarding the irrelevant rules.

Clustering techniques which seek to cluster the data in a number of groups can deal with high-dimensional problem. But these techniques have their own difficulties. First, the number of clusters is predefined by the user. Cluster validity measures and/or compatible cluster merging techniques are used further to refine the structure. Second, to avoid that the clusters are solely based on the spatial distribution of the samples, the clustering is performed in the product space of Z and the output space. Therefore, all the regressors are automatically the features variables. It results the presence of irrelevant variables in the both parts of the structure.

For these reasons, we adopt a k-d partition that permits to fulfil more easily the parsimony principle. Indeed, it has this appealing feature: the model structure is gradually extended by allocating new models in areas where most needed according to the model accuracy. The partition technique is presented in the next section. We do not discuss the local model orders estimation task. We consider each local model to be a first order model because with the multiple model approach, a trade-off must be achieved between the number and the complexity of the local models. So, they are assumed simple and the structure optimisation algorithm is let to find the appropriate size of the local model network.

3. PARTITIONING THE FEATURE SPACE

3.1. CRITERION FOR STRUCTURE OPTIMISATION

Structure identification is comparable to the determination of an appropriate model which strikes a balance between model accuracy and complexity, i.e. a model that has good generalisation capacities. The straightforward manner to evaluate model generalisation ability is the cross-validation. So, the following criterion is used as structure performance index:

$$J_{\text{STRUC}} = \sum_{t=1}^{N_B} (y(t) - \hat{y}(t))^2 \quad (10)$$

where N_B is the number of testing data and $\hat{y}(t)$ is the estimation obtained by running the model in parallel. Notice that, unlike the criterion J (eq.9), the above objective function is an output error criterion.

3.2. DECISION TREE DECOMPOSITION

The framework is similar to those of local search algorithm (LSA) of Johansen and Foss (1995) and local linear model tree (LOLIMOT) of Nelles (1997). The k-d partition results from an iterative procedure of splitting the feature space in hyper-rectangles. Given the feature variables $Z(t)$, the feature space Z is reduced to a hyper-box since the measurements are bounded:

$$Z = [Z_{1,\min} \quad Z_{1,\max}] \times [Z_{2,\min} \quad Z_{2,\max}] \times \dots \times [Z_{n_Z,\min} \quad Z_{n_Z,\max}] \quad (11)$$

where $Z_{i,\min}$ and $Z_{i,\max}$ are the lower and the upper bounds of the range of the variable Z_i .

At the n^{th} step of the procedure, the feature space is already divided in n areas $Z^{(1)}$, $Z^{(2)}$, ..., $Z^{(n)}$. The accuracy of the structure is tested in each subspace using a local index performance expressed as the weighting sum of the quadratic output error, i.e.:

$$\begin{aligned} \varepsilon_i &= \frac{1}{N_i} \sum_{t=1}^{N_B} \omega_i(Z(t), \beta) (y(t) - \hat{y}(t))^2 \\ N_i &= \sum_{t=1}^{N_B} \omega_i(Z(t), \beta) \end{aligned} \quad i = 1, \dots, n \quad (12)$$

where N_i is the pseudo-measure of the number of observations belonging to the subspace $Z^{(i)}$. Let $Z^{(k)}$ the subspace where the model has the most lack of accuracy. It is divided into two subspaces $Z^{(k_1)}$ and $Z^{(k_2)}$ by a cut along a hyperplane orthogonal to one of the axis of the feature space. The two new hyper-rectangles are defined as follows:

$$Z^{(k_1)} = \{Z(t) \in Z^{(k)} / Z_j(t) < \xi_j\} \quad \text{and} \quad Z^{(k_2)} = \{Z(t) \in Z^{(k)} / Z_j(t) \geq \xi_j\} \quad j \in \{1, 2, \dots, n_Z\} \quad (13)$$

The splitting point ξ_j must belong to the range of the feature variable Z_j in the subspace $Z^{(k)}$. The validity functions of the new hyper-rectangles $Z^{(k_1)}$ and $Z^{(k_2)}$ are directly deduced from that of $Z^{(k)}$. Indeed, as the cut is made only in the dimension Z_j , the gaussians defining the two child subspaces have the same parameters as the validity function of the parent area excepted the centre and the dispersion along the axe Z_j which are given by the equations below

$$\begin{aligned}
c_{k_1,j} &= \frac{\xi_j + Z_{j,\min}^{(k)}}{2} & \sigma_{k_1,j} &= \frac{\xi_j - Z_{j,\min}^{(k)}}{2} \gamma_j \\
c_{k_2,j} &= \frac{Z_{j,\max}^{(k)} + \xi_j}{2} & \sigma_{k_2,j} &= \frac{Z_{j,\max}^{(k)} - \xi_j}{2} \gamma_j
\end{aligned} \tag{14}$$

where γ_j determines the overlap between the two new submodels.

The validity function parameters of the other hyper-rectangles remain unchanged. Thus, the global vector θ of local model parameters is optimised by minimising the criterion J . The split in all the dimensions of the feature space is tested. This leads to n_Z different structures \mathbf{M}_{n+1} which have $n+1$ local models. The selected structure is the one that yields the best performance index J_{STRUC} . An illustration of the procedure is shown in figure 1.

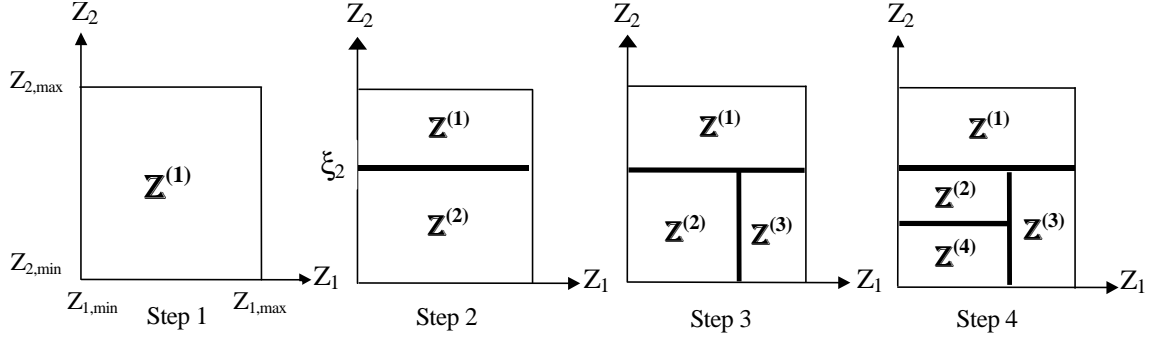


Figure 1: K-d partition of a two dimensional feature space

The algorithm is summarised by the following steps:

1. Choose the local model inputs. Identify the global (linear) model. Let $\mathbf{Z}^{(k)} = \mathbf{Z}$ and $n=1$.
2. For each feature space variable Z_j ($j=1, \dots, n_Z$)
 - find the splitting point ξ_j , the overlap parameter γ_j and split $\mathbf{Z}^{(k)}$ in the dimension Z_j ,
 - calculate the parameters of the validity function of the new submodels using equation (14) and optimise θ by least squares,
 - evaluate the corresponding local model network by computing J_{STRUC} (eq. 10).
3. Select the structure \mathbf{M}_{n+1} with the lowest criterion J_{STRUC} .
If $J_{\text{STRUC}}(\mathbf{M}_{n+1}) > J_{\text{STRUC}}(\mathbf{M}_n)$, go to step 5.
4. Find the subspace $\mathbf{Z}^{(k)}$ where the structure has the worst approximation performance with $k = \arg \min\{\epsilon_1, \epsilon_2, \dots, \epsilon_n\}$. Increase n , i.e. let $n=n+1$ and go to step 2.
5. End of the procedure

It must be pointed out that the great difficulty of the technique is the determination at each step of the splitting point and the overlap parameter. These parameters can be optimised by minimising the criterion J , subject to the constraints that the centres of the resultant gaussians belong to the new hyper-rectangles and the dispersion are proportional to the size of the new areas. This will ensure the local effect of the validity function and will avoid reactivation (Shorten and Murray-Smith, 1997) that is linked to the normalisation of the validity functions (eq. 3). But the computational cost becomes high. For the application and to simplify the problem, the overlap parameter is set to a fixed value (according to Johansen and Foss (1995), γ_j varies between 0.25 and 2 but the typical values are around 1). The splitting point is estimated by scanning the range of the considered feature variable in $\mathbf{Z}^{(k)}$. It is given by the formula:

$$\xi_j = Z_{j,\min}^{(k)} + \alpha(Z_{j,\max}^{(k)} - Z_{j,\min}^{(k)}) \quad \alpha \in [0.1, 0.9] \tag{15}$$

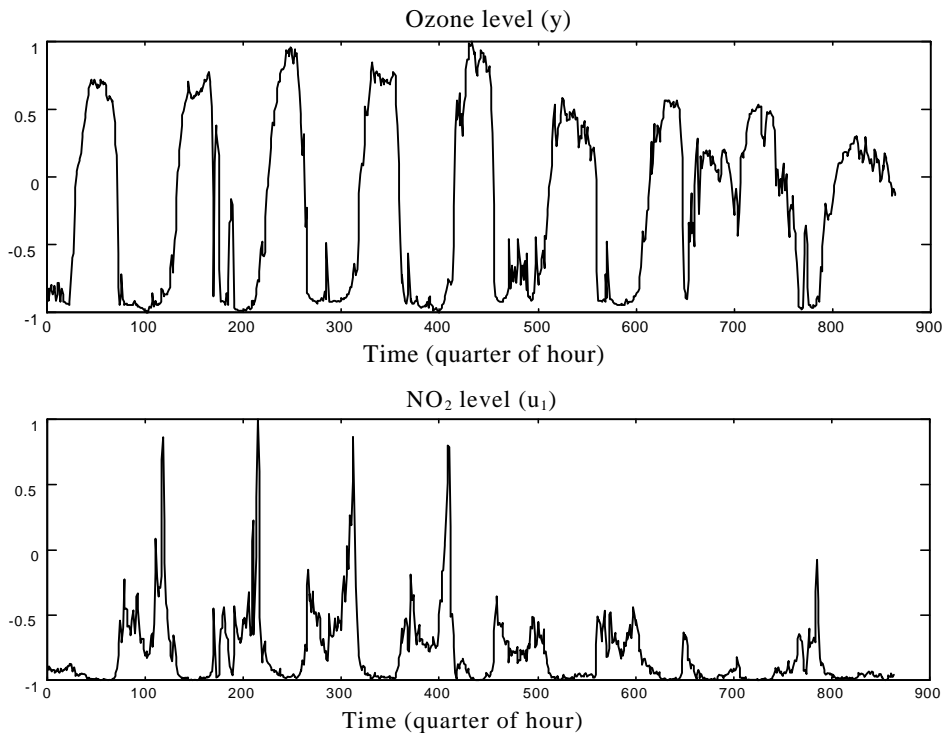
Notice that Nelles (1997) in his algorithm LOLIMOT, divides the hyper-rectangle into two halves. This provides a more simpler solution but in the ozone modelling problem, the best results have been obtained using the above procedure.

4. APPLICATION: OZONE CONCENTRATION MODELLING

The procedure is illustrated under the problem of ozone concentration modelling in the city of Nancy (France). Measurements of ozone and eight other variables are available. These variables are: nitrogen oxide and dioxide concentration (NO and NO₂), temperature, relative humidity, pressure, solar radiation, wind speed and wind direction. The variables are sampled over a period of quarter of hour that is 96 samples per day. We must select among the eight variables, those that are relevant to explain the variations of ozone concentration. Preliminary analysis of the data reveals that the relevant variables to describe the phenomenon are NO₂ (u₁), temperature (u₂), solar radiation (u₃) and wind speed (u₄) and relative humidity (u₅). Figure 2 shows the temporal evolution of these variables and ozone concentration during a few days. We can remark that the ozone level which is low during the nocturnal period rises progressively at the beginning of the day and reaches maximum values in the afternoon. Then, it decreases and attains again low value during night and the process continues. So, we can distinguish a nocturnal and a diurnal behaviour of the phenomenon. Low values of solar radiation correspond to nocturnal period.

The training data are recorded over a period of nine days. To test the generalisation capacity of the model, two validation data sets are used: the first concerns six consecutive days and the second fourteen other consecutive days. As the variables have different ranges and physical units, they are normalised so that their domain becomes [-1, 1]:

$$\tilde{x} = \frac{2x - (\max(x) + \min(x))}{\max(x) - \min(x)} \quad x = \{z, y, u\} \quad (16)$$



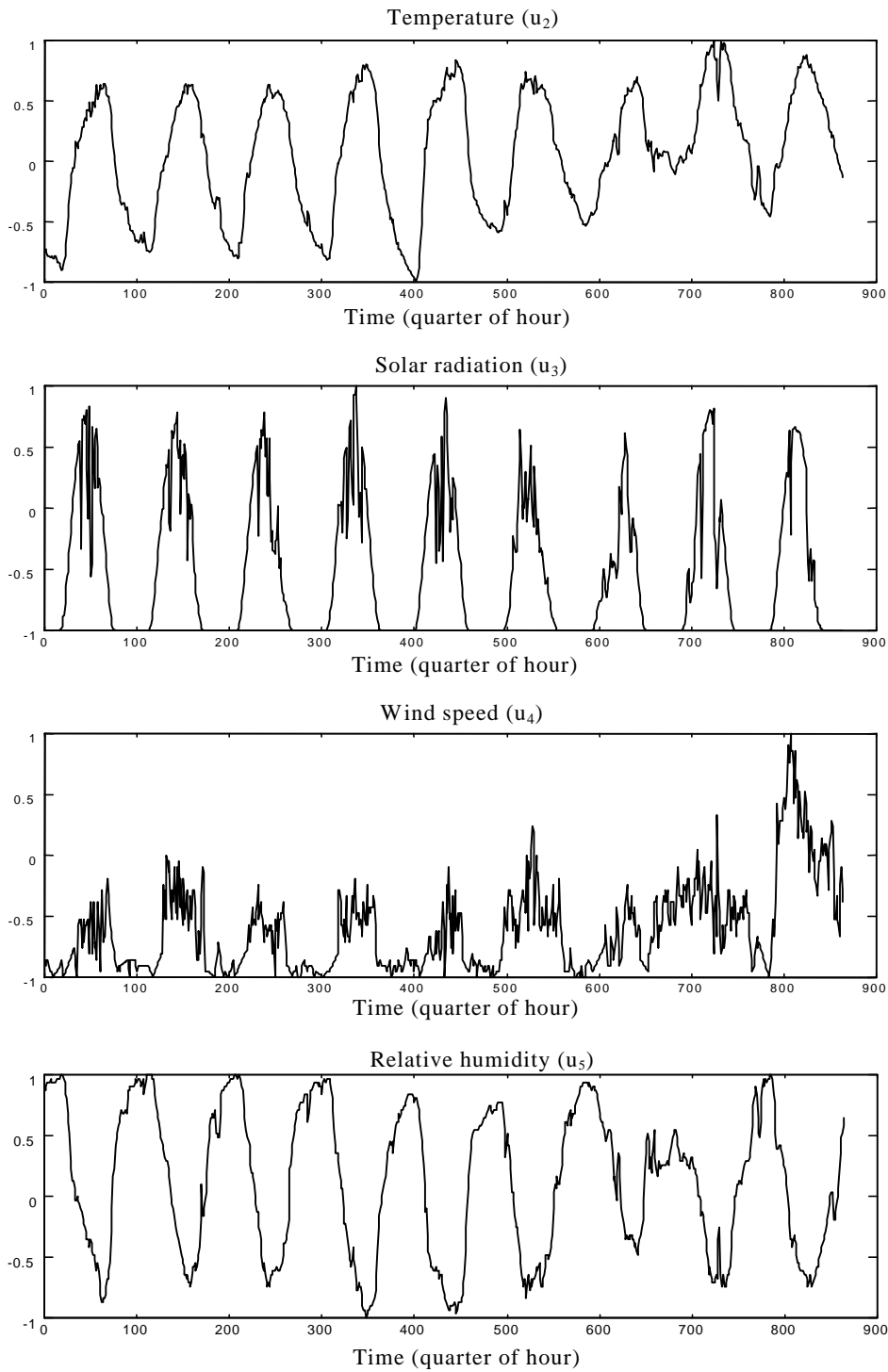


Figure 2: Input-Output signals of the identification data set

To test the decomposition procedure, we have to choose the input variables u and the feature variables Z . Using an idea of Tanaka et al. (1995) the selected inputs of the local models are the variables intervening in the optimal linear structure. Thus, the four following inputs are retained: NO_2 , temperature, solar radiation and wind speed. Without other information, the feature variables are chosen as $Z(t) = [u_1(t-1) \ u_2(t-1) \ u_3(t-1) \ u_4(t-1) \ u_5(t-1)]$. The partition algorithm is run for $\gamma=0.9$. Figure 3 shows the evolution of the output error criterion versus the number of local models for the three data sets.

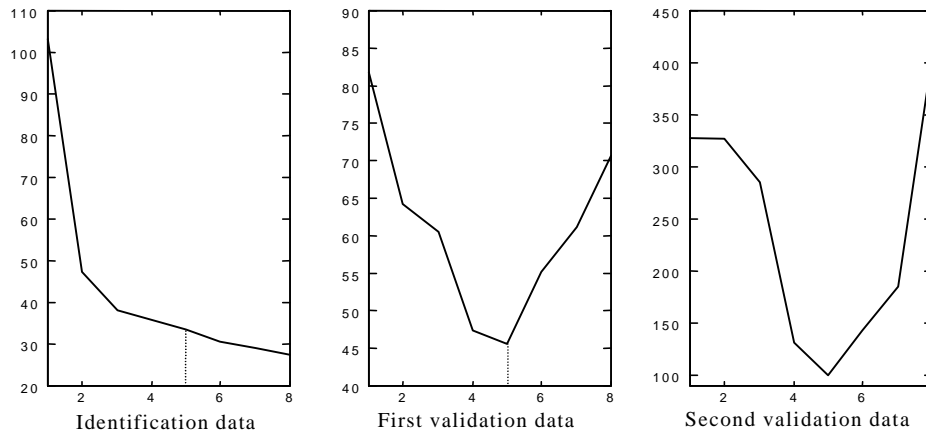


Figure 3 : Performance index decreasing as a function of the number of local models

The selected structure has five local models. Indeed, there is slightly improvement of the performance index on the identification data set when adding new submodels in the structure. Moreover the criteria on the two validation data stop to decrease and growth for a number of local models superior to five. The decision tree corresponding to the retained structure is shown in figure3. We have indicated on this figure the dimension split at each step of the decomposition procedure.

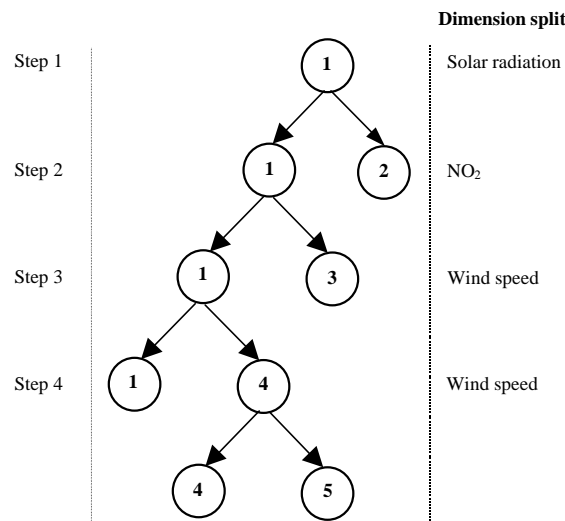


Figure 4: Decision tree of the structure selected

It can be noticed that no cut has been made in the dimension defined by the temperature and the relative humidity. The conclusion to be drawn is that for the data used, these two variables are not characteristic of the non-linear behaviour of the phenomenon. Hence, the multiple model structure is simplified since the five-dimensional feature space is reduced to a three-dimensional domain. It is also interesting to notice that only one cut is made in the dimension defined by the solar radiation. Examining the obtained ranges, it is found that the two ranges correspond approximately to the nocturnal and diurnal periods. This clearly indicates the partition algorithm ability to capture this typical behaviour of the ozone phenomenon. The two remaining variables defining the feature space, say wind speed and NO₂ are representative of the chemical and physical aspects of the phenomenon. The wind speed has been divided in three ranges while NO₂ level is split into two intervals. However, it is difficult to give a more detailed physical meaning to the obtained model. Notice that if a grid partition is applied with the above number of ranges for each feature variable, we will get a structure with twelve local models instead of actually five local models. It is the reason that the lattice partition is unsuitable for a relative high number of inputs mainly if there is no further refinements of the obtained structure.

The model performances on the training data are plotted in figure 5 while figures 6 and 7 depict the performances on the validation data sets. We precise that all the plots show simulations and not one-step-ahead predictions. The measurements are plotted in solid lines while the dashed lines indicate the estimations. We notice that the multiple model provides a good approximation of high ozone level – that is important for the application – as well as the nocturnal low concentration on all data set. Although their respective structures are different, the results of the multiple model are compared to those obtained with the best linear ARX model which has orders $p=2$ $q=1$ and contains the four aforementioned inputs variables. For illustration purpose, the outputs provided by the linear model on the training data and on the second validation data set are shown in figure 8. We remark that the linear model gives fairly good approximation on the training data but unfortunately, it predicts poorly when applied to the validation data as we can see in figure 8-b. We deduce that the linear model is unsuitable to deal with this ozone concentration modelling problem.

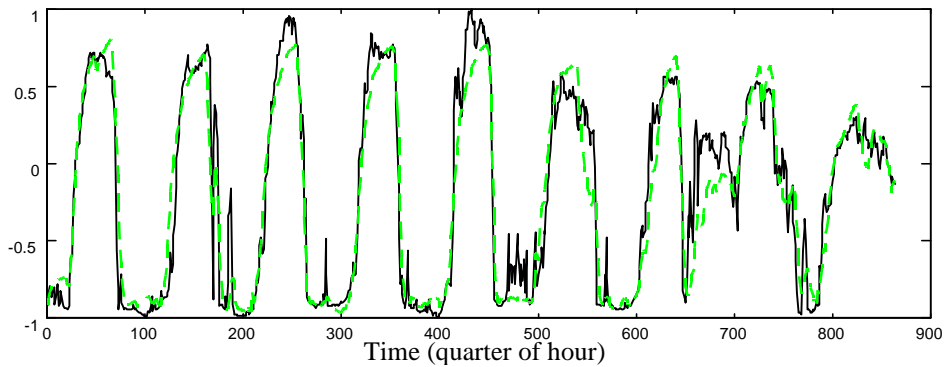


Figure 5: Comparison measurements - multiple model estimation on the training data

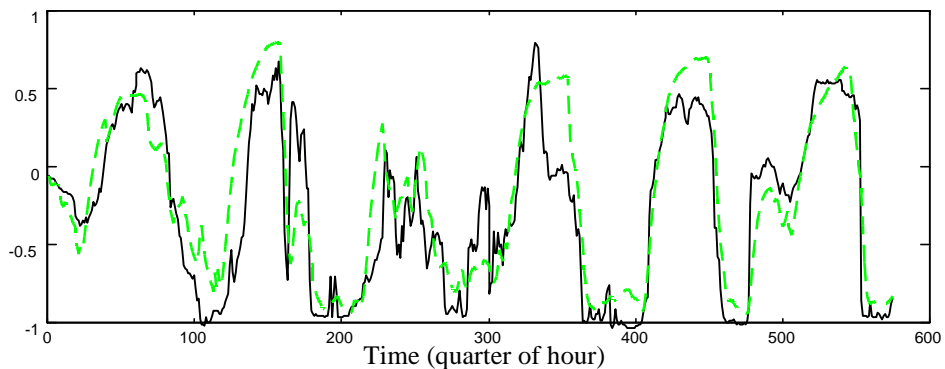


Figure 6: Multiple model performances on the first validation data

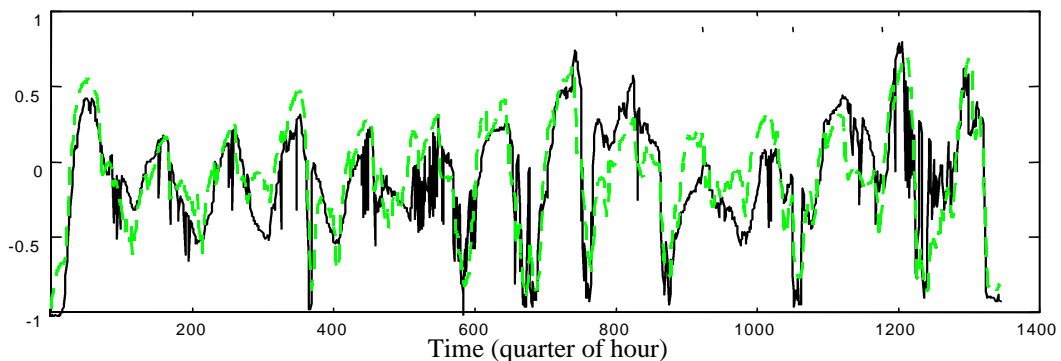
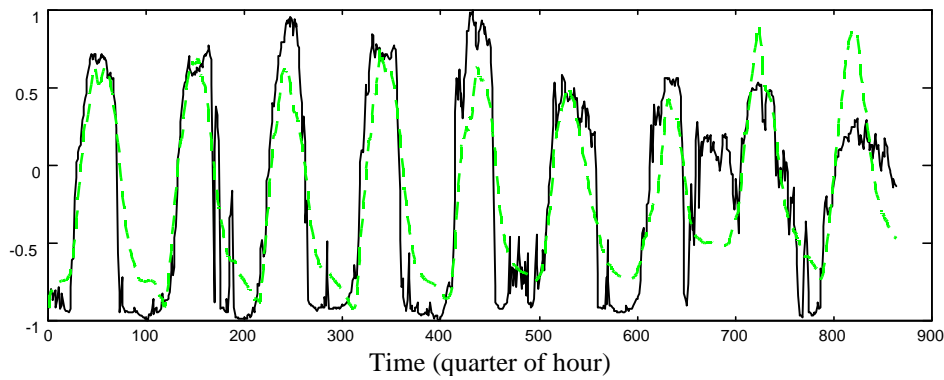
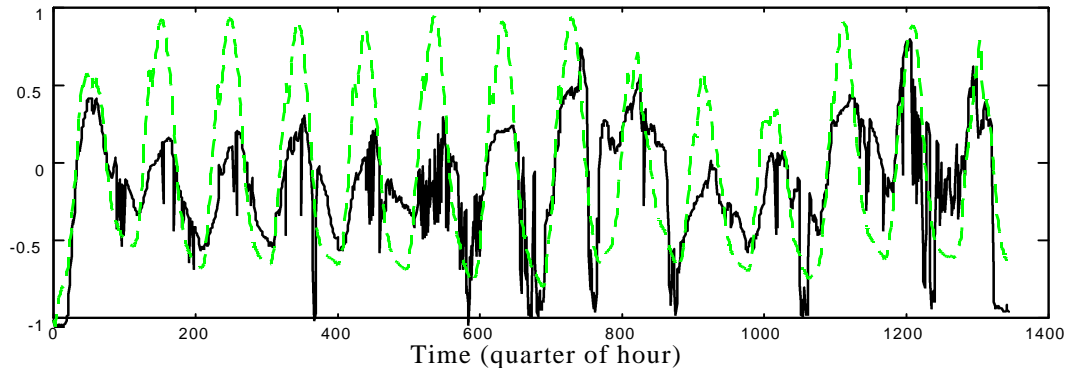


Figure 7: Comparison measurements - multiple model estimation on the second validation data



(a)- Training data



(b)- Second validation data

Figure 8: Performances of the best linear ARX model

We would like to emphasize that the obtained results are better than those of (Mourot and Ragot, 1997). Indeed the model presented in this reference predicts well the maximum values of ozone concentration but shows important differences or lags between estimation and measurements during nocturnal periods. It is not the case of our model and this must be underlined even if less accurate estimations of low ozone level are not a real drawback. Moreover, to deal with this multivariable modelling problem, we neither impose a lot of constraints nor discard irrelevant local models because a constructive structure optimisation technique is employed.

5. CONCLUSION

This paper has investigated the modelling of ozone concentration. The phenomenon is known non-linear and multi-dimensional. The multiple model approach used to solve the problem has been presented and we have discussed the related difficulties. In order to reduce the complexity of the problem, the structure identification is performed by a hierarchical technique. The results concerning the ozone modelling are satisfactory because the obtained structure is simple and takes into account the main peculiarities of the phenomenon. The multiple model structure has been validated with success on two relative short periods (one to two weeks). The validation on much longer periods (say one to three months) remains to be carried out to test the efficiency of the model. This will certainly underscore the inadequacies of the actual model; the future work will thus focus on the extension of the structure to take account of them.

Acknowledgements: We would like to thank the staff of AIRLOR who provide us the data used in this application.

REFERENCES

- Johansen T.A. and Foss B. (1992) "Nonlinear local model representation for adaptive systems". In *Proc. IEEE Int. Conf. On Intelligent Control and Instrumentation*, Singapore, vol 2, pp 667-682.
- Johansen T.A. and Foss B. (1993) "Constructing NARMAX using ARMAX". *Int. Journal of Control*, vol 58, N°5, pp 1125-1153.
- Johansen T.A. and Foss B. (1995) "Identification of non-linear systems structure and parameters using regime decomposition". *Automatica*, vol 31, N° 2, pp 321-326.
- Mourot G. and Ragot J. (1997) "Identification de modèle de Takagi-Sugeno- Application à la modélisation de la concentration d'ozone". *European Journal of Automation. RAIRO-API-JESA* vol 31, N°9 - 10/1997, pp 1587-1608
- Nelles O. (1997) "Nonlinear system identification with neuro-fuzzy methods". Chapter in : da Ruan(ed.) : *Intelligent Hybrid Systems*. Kluwer Academic Publishers, Dordrecht.
- Shorten R. and Murray-Smith R. (1997) "Side-effects of normalising basis functions in local model networks". Chapter Eight in: *Multiple Model Approaches to Modelling and Control*, Murray-Smith R. and Johansen T.A, Ed. Taylor and Francis, London.
- Takagi and Sugeno (1985) "Fuzzy identification of systems and its application to modelling and control". *IEEE Trans. On Systems Man and Cybernetics*, 15, pp 116-132.
- Tanaka K., Sano M. and Watanabe H. (1995) "Modelling and control of carbon monoxide concentration using a neuro-fuzzy technique". *IEEE Trans on Fuzzy Systems*, Vol 3. N°3, pp 271-279.
- Tan E., Mourot G. Maquin D. and Ragot J (1994) "Identification of fuzzy models". *International workshop on Fuzzy Technologies and in Automation and Intelligent System. Fuzzy Duisburg'94*. Duisburg, Germany, pp 190-199.