Daniela Cocchi, Enrico Fabrizi*,
Carlo Trivisano

# A stratified model for the analysis of ozone trends in an urban area

Gruppo di Ricerca per le Applicazioni
della Statistica ai Problemi Ambientali

Working paper GRASPA n. 13

# Daniela Cocchi, Enrico Fabrizi, Carlo Trivisano

# A stratified model for the analysis of ozone trends in an urban area

**1**

## 1. Introduction

In this paper we face the problem of determining the trend of daily maxima of ozone concentrations measured in an urban area. Because of the close dependence between ozone and weather conditions, the major issue is to assess whether any observed trend in the data is due to variations in the emissions or to particular climatic conditions.

This subject has been widely discussed in literature and excellent reviews already exist (Thompson *et al.*, 2000).

We propose a new method for trend inspection that refines the tree based approach considered by Huang and Smith (1999). The analysis proposed by these authors can be described in two steps: first, a regression tree is introduced to partition the daily maxima into groups with an homogeneous ozone level, subsequently the trend is evaluated by means of a set of alternative linear homoschedastic random effects models, some of which allow for a different trend within each group. This approach, referred to as stratified (Thompson *et al.*, 2000), has the merit of allowing for the estimation of different trends at different levels of ozone concentration and weather conditions. The idea that trends may not be the same at each level of the process is not new. Smith (1989), modeling exceedances of daily maxima over a high threshold under an Extreme Value approach, considers a trend component for the location parameter of the Generalized Pareto distribution (Pickands, 1971), separately for groups of days approximately corresponding to months.

Similarly to Huang and Smith (1999), we propose a tree based partitioning of observations. We assume the daily maxima of ozone concentrations to be Weibull distributed and propose a random effects model for the natural logarithm of the quasi-scale parameter of this distribution, where the considered effects are represented by the year and the homogeneous ozone regimes resulting from tree partitioning. Our approach is free of any hypothesis about the shape of trend, allowing for trend inspection in short periods.

The Weibull distribution arises naturally in the context of maxima; moreover it represents a flexible assumption since, for different values of its parameters, it ranges from approximated normality to highly skewed forms. It may then be expected that our model would fit well on the tails

of distributions and provides a more precise estimation of high percentiles and exceedances.

Our modeling of the natural logarithm of the quasi-scale parameter is very close to the proposal of Cox and Chu (1993); nonetheless, we note that their approach is non stratified and their trend estimation is based on the assumption of a linear functional form for the trend component.

The models we propose are based on the assumption of conditional independence of observations given the quasi-scale parameter, that is within each group-by-year cluster. Ozone concentrations are known to be characterized, marginally, by strong serial autocorrelation. For this reason any independence assumption deserves careful investigation that is carried out by estimating a model including an autoregressive component for the Weibull quasi-scale parameters.

We apply this method of trend assessment to the series of daily maxima of hourly ozone concentrations measured from a single monitoring site located in the city of Bologna, Italy, in the the period 1994-1998.

We adopt a Bayesian approach for inference. Models are solved by means of Gibbs sampling routines, as they are implemented in the software BUGS (Spiegelhalter *et al.*, 1996).

The paper is organized as follows. Section 2 contains a description of the data, while in section 3 the proposed models are introduced. In section 4 results about trend estimation are discussed along with model checking in terms of percentiles and yearly number of exceedances prediction performances. The sensitivity analysis with respect to the assumption of conditional independence is treated in section 5.
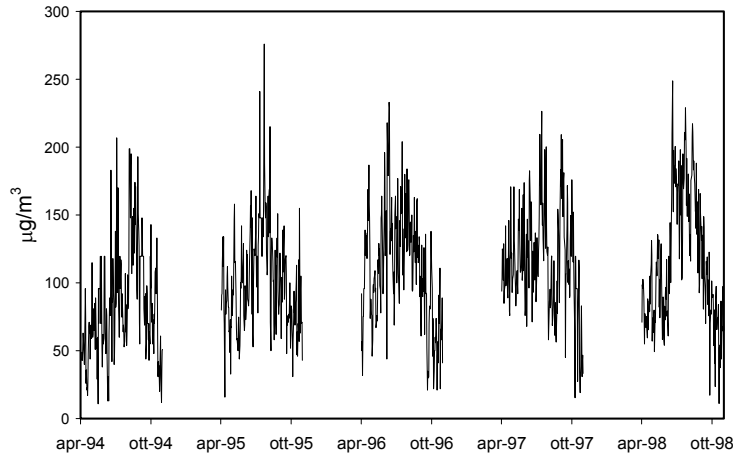
## 2. The data

In this paper we analyze the series of daily maxima of ozone concentrations over the metropolitan area of Bologna, in the North of Italy, in the period 1994-1998.

Data have been gathered, and their quality assessed, by the ARPA (*Agenzia Regionale per la Protezione Ambientale, i.e.* the Regional Environment Protection Agency) of the Emilia-Romagna region, that is charged by the Italian law with all monitoring and protection policies related to air pollution.

Maxima are calculated on the basis of hourly measurements from a single monitoring site situated in *Giardini Margherita*, a park close to the centre of the city. We consider exclusively data gathered during the ozone seasons (from April 1[st], to October 31[st]) yielding a gross total of 1067 observations. For 66 days the maximum is not computed because of measurement failures or other causes. We calculate regression trees by simply discarding the missing observations. Then, the 66 missing values of maxima are imputed using the selected regression tree model.

These imputations are needed in order to perform the comparison between models based on conditional independence and models including autoregressive components, as the one described in section 5.

The path of the series is shown in Figure 1.

*Figure 1*
*Daily maxima of hourly concentrations in Bologna, ozone seasons, from 1994 to 1998*



The meteorological variables are collected at the surface weather station of *Borgo Panigale*, a few kilometers away from the monitoring site of *Giardini Margherita*. Among the hourly available measurements, we consider the noon values of temperature (°C), dew point temperature(°C),

pressure (*mb*), relative humidity (%), visibility (*km*), wind speed (*m/s*) and wind direction (° from N).

Following Bloomfield *et al.* (1996) and Huang and Smith (1999) we calculate some additional meteorological variables:

Wind.U = *(-Wind.speed \*sin(2π\*Wind.dir/360)*,

Wind.V = *(-Wind.speed \*sin(2π\*Wind.dir/360)*,

Temperature2 = (*Temperature – 60)$^2$/10,*

Temperature3 = *Temperature – 60)$^3$/1000.*

Moreover, lagged values (1 and 2 days) for all the meteorological variables are considered.

Missing data are present in the series of meteorological variables as well. In this case we do not propose any imputation, since they are not needed in building the regression tree.

## 3. The proposed models

### 3.1 Identification of homogeneous regimes

Partitioning methods based on classification and regression trees (Breiman *et al.* 1984) gained a wide popularity in environmental statistics since they are particularly suitable to explore the non linear relationship between surface ozone and weather (Burrows *et al.*, 1995); moreover they perform rather well in forecasting occurrences of high concentrations. However, trees cannot be directly employed in trend assessment, since a trend component cannot be included into this class of models. A remarkable attempt to exploit the power of regression trees in meteorological adjustment for the purpose of trend investigation is carried out by Huang and Smith (1999). Working on the well-known Chicago ozone data, they first partition the observations into homogeneous clusters using a regression tree method (Clark and Pregibon, 1991). Once the tree is assessed, they apply ANOVA models based on the normality assumption to each cluster identified by the tree in which a year effect is included.

Following Huang and Smith, we build a regression tree for ozone daily maxima using the set of meteorological covariates described in section 2. We follow the CART approach (Breiman *et al.*, 1984), that does not

make use of any distributional assumption. At each node (parent) of the tree, data are partitioned into two homogeneous subsets (left child and right child) on the basis of the squared loss ($L^2$), that is by maximising

$$\Delta L^2 = L^2_{parent} - (L^2_{left\ child} + L^2_{right\ child})$$

The usual cost-complexity criterion described in Breiman *et al.* (1984) is considered for pruning the tree; cross-validation is employed to select the degree of pruning.

*3.2 Models for the trend*

Let's denote the vector of daily maxima with $y_i$. For each vector component we assume that

$$y_i \mid \lambda_i, \zeta \sim Weibull(\lambda_i, \zeta), \quad i = 1, ..., N = 1067, \tag{3.1}$$

*i. e.:*

$$f(y_i \mid \lambda_i, \zeta) = \zeta \lambda_i y_i^{\zeta-1} e^{-\lambda_i y_i^{\zeta}} . \tag{3.2}$$

Under the parameterization (3.2), the Weibull distribution has moments:

$$E(y_i \mid \lambda_i, \zeta) = \frac{\Gamma\left(1 + \dfrac{1}{\zeta}\right)}{\lambda_i^{1/\zeta}}$$

$$\tag{3.3}$$

$$V(y_i \mid \lambda_i, \zeta) = \frac{\Gamma\left(1 + \dfrac{2}{\zeta}\right) - \Gamma^2\left(1 + \dfrac{1}{\zeta}\right)}{\lambda_i^{2/\zeta}}$$

In (3.1) a different quasi-scale parameter $\lambda_i$ is considered in each "group by year" cluster, but a common shape parameter $\zeta$ across clusters is assumed. This is equivalent to the hypothesis of a common coefficient of

variation for all the conditional distributions, since it depends on the shape parameter only. In fact we have that:

$$CV\left(y_i \mid \lambda_i, \zeta\right) = \frac{\sqrt{\Gamma\left(1 + \frac{2}{\zeta}\right) - \Gamma^2\left(1 + \frac{1}{\zeta}\right)}}{\Gamma\left(1 + \frac{1}{\zeta}\right)} \tag{3.4}$$

We note that the assumption of a common shape parameter is made also by Cox and Chu (1993) and is consistent with Huang and Smith (1999), who suppose homoschedastic residuals for their group-specific random effects models. Moreover, since mean, median, and other basic summaries of the Weibull distribution depend on both parameters, the assumption of a common $\zeta$ allows for the influence of the group and year effects to be read directly in terms of different quasi-scale parameters. Otherwise, year- and group- specific shape parameters could incorporate part of the covariates effect in a way that makes comparisons of distributions across groups more difficult.

We consider two different models for the logarithm of the quasi-scale parameters: a model with interaction characterized by independent group by year random effects (referred to as model 1) and an additive model assuming independent year and group effects (model 2).

More formally, model 1 is characterized by

$$\ln\left(\lambda\right) = W\delta \tag{3.5}$$

where $\lambda = \left(\lambda_1, \lambda_2, ..., \lambda_N\right)'$ and $\delta$ is a $\left(G \times P\right)$ dimensional random vector ($P$ is the number of clusters identified by the tree, while $P$ the number of years in the period considered by this study). $W$ is defined, using a block notation, as $W = \left(W_1 \mid ... \mid W_P\right)$; the generic element of $W_p$ ($p = 1, ..., P$) can then be described as

$$\left(w_{ij(p)}\right) = \begin{cases} 1 & \text{if day } i \text{ falls in year } p \text{ and cluster } j, \quad j = 1, ..., G \\ 0 & \text{otherwise} \end{cases}$$

Model 2 is characterized by the equation:

$$\ln(\lambda) = \mathbf{1}\alpha + X\beta + Z\gamma \tag{3.6}$$

where $X$ is a $N \times (G-1)$ design matrix whose generic element is defined as

$$(x_{ij}) = \begin{cases} 1 & \text{if day } i \text{ falls in cluster } j+1 \\ 0 & \text{otherwise} \end{cases} \tag{3.7}$$

and Z is a $N \times (P-1)$ matrix for which we have

$$(z_{ij}) = \begin{cases} 1 & \text{if day } i \text{ falls in year } j+1 \\ 0 & \text{otherwise} \end{cases} \tag{3.8}$$

The general intercept $\alpha$ is introduced to avoid multicollinearity of year and group effects. As regards the problem of trend determination, we define the trend as the sequence of parameters associated to year effects, thus avoiding the assumption on any functional form. We note that model 1 corresponds to the hypothesis of a group specific trend, while model 2 assumes a common trend for all the groups identified by the tree.
We also consider a third benchmark model (model 3), formalizing the hypothesis of absence of trend; it is characterized by

$$\ln(\lambda) = \mathbf{1}\alpha + X\beta \tag{3.9}$$

For the specification of prior distributions, we assume, across all models, the following prior for the shape parameter:

$$\zeta \sim Gamma(3,1) \tag{3.10}$$

implying $E(\zeta) = 3$, $V(\zeta) = 3$. This distribution is centered on the value of $\zeta$ for which the Weibull has a shape similar to that of the

Normal, but is sufficiently diffuse to give support also to more skewed distributions.

Parameters associated to random effects are given in all cases a Normal prior distribution centered on 0, while assumptions on their precisions change according to the model chosen. In particular in model 1 we specify that

$$\delta \mid \tau \overset{ind}{\sim} N\left(\underline{0}, \Gamma\right) \tag{3.11}$$

where the precision matrix $\Gamma$ is block diagonal and can be written as

$$\Gamma = diagblock\left(\tau_1 I, ..., \tau_G I\right).$$

For the vector of hyperparameters $\tau = \left(\tau_1, \tau_2, ..., \tau_G\right)$ we assume that they are a priori independent and are distributed as

$$\tau_j \sim Gamma\left(0.01, 0.01\right), \quad j = 1, ..., G,$$

that is, we assume a common prior distribution for the year intercept related to the same group, while a diffuse Gamma is chosen for the parameters $\tau_j$; we remark that this non-informative "reference" solution is designed mostly for computational convenience.

In model 2 we introduce the following priors:

$$\alpha \sim N\left(0, 0.001\right)$$
$$\beta \mid \nu_1 \sim N\left(\underline{0}, \nu_1 I\right)$$
$$\nu_1 \sim Gamma(0.1, 0.1) \tag{3.12}$$
$$\gamma \mid \nu_2 \sim N\left(\underline{0}, \nu_2 I\right)$$
$$\nu_2 \sim Gamma(0.1, 0.1)$$

thus assuming common distributions for each year and group effect. Consistently, the same prior assumptions on $\alpha$ and $\beta$ are introduced in model 3.

We calculate all posterior distributions by using Markov chain Monte Carlo (McMC) sampling algorithms. In particular we use the software BUGS (Spiegelhalter, Thomas, Best and Gilks 1996) that is based on Gibbs sampling. The selected prior distributions have standard functional forms that lead to log-concave full conditionals in all cases. As regards the assessment of convergence we consider the multiple chain approach suggested by Gelman and Rubin (1992), running three different chains with well separated starting points for each model. The visual inspection of chains path and the modified Gelman and Rubin statistic (Brooks and Gelman, 1998) are our basic convergence assessement tools. We run 10,000 iterations for each chain, discarding on average a conservative "burn in" of 3,000, thus yielding an approximate 20,000 draws from the posterior of each model.

## 4. Model comparison and discussion of empirical results

### 4.1 The regression tree

On the basis of the CART method described in section 3.1, the tree shown in Figure 2 (see Appendix) is built.

As it should be expected from the photochemical nature of reactions leading to concentration of surface ozone, maximum daily temperature and other related variables are responsible for the most important splits in the tree, but we note also that relative humidity and visibility play a relevant role in determining homogeneous ozone regimes.

The optimal tree identifies eight groups with very different sizes, some of them being very small. The classification of daily maxima by group and year is reported in Table 1.

We note that group 8, characterized by extreme concentrations, includes only 13 maxima throughout the five years period, and group 5, even though characterized by an intermediate level of the process, includes only 23 cases. Once data are cross-classified by year, the resulting "group by year" clusters are constituted by very few observations.

| | | Year | | | | | |
|---|---|---|---|---|---|---|---|
| | | **1994** | **1995** | **1996** | **1997** | **1998** | |
| Group | **1** | 17 | 20 | 18 | 25 | 23 | 103 |
| | **2** | 40 | 66 | 60 | 45 | 63 | 274 |
| | **3** | 28 | 17 | 29 | 10 | 12 | 96 |
| | **4** | 14 | 16 | 10 | 3 | 6 | 49 |
| | **5** | 3 | 3 | 3 | 8 | 6 | 23 |
| | **6** | 61 | 61 | 81 | 90 | 62 | 355 |
| | **7** | 45 | 29 | 12 | 33 | 35 | 154 |
| | **8** | 4 | 2 | 0 | 0 | 7 | 13 |
| | | 212 | 214 | 213 | 214 | 214 | 1067 |

The hierarchical structure of priors for the models described in section 3.2 (see (3.9), (3.10) and (3.11)) is very helpful for the estimation of parameters as it allows, at the price of some shrinkage, for the borrowing of information across groups.

### 4.2 Model comparisons

Within the Bayesian framework, models are compared, in principle, by means of the Bayes Factors (BF). Since they are rather difficult to compute, a large sample approximation of *-2ln(BF),* given by

$$\Delta BIC = -2\ln\left[\frac{\sup_{M_0} f(y\,|\,\theta_0)}{\sup_{M_k} f(y\,|\,\theta_k)}\right] - \left(p_k - p_0 \ln n\right) \tag{4.1}$$

(see Schwartz, 1978), is commonly used. We compare the three proposed models by means of (4.1) that, moreover, has the merit of not referring to prior assumptions. We note that in (4.1) the subscripts $M_k$ *(k=1,…,K)* index the set of competing models and $\theta_k$ is the $p_k$ dimensional parameter indexing the likelihood associated to each model.

The null model $M_0$ against which the others are compared is model 3. Model comparisons are summarized in Table 2.

*Table 2*
*Model comparisons by means of $\Delta$ BIC*

|  | $\Delta$ BIC |
|---|---|
| Model 1 *vs* Model 3 | -21.014 |
| Model 2 *vs* Model 3 | 100.274 |
| Model 1 *vs* Model 2 | -121.288 |

The Bayes Information Criterion (4.1) indicates that both model 1 and 2 perform better than model 3, supporting the hypothesis that there is a significant variation of ozone levels  over the years, even within the homogeneous groups identified by the tree. It is also apparent that model 2 is to be preferred to model 3. This evidence suggests that the trends in the groups are not "so different" to justify the introduction of group-specific trends.

This can be clearly understood by looking at Figure 3 (in the Appendix) where posterior means of the Weibull distributions computed under model 1 and model 2 are compared by group and year.

From Figure 3 (in the Appendix) we can see that the only evident discrepancies between the common and the group-specific trends arise in group 5 and 8, that are particularly small. On the contrary, this discrepancy is almost negligible in larger groups.

The adequacy of the selected model is checked following a posterior predictive approach, that can be easily implemented on the basis of the McMC output. We check the consistency of data generated from the posterior predictive distribution with observed maxima in terms of the mean and the 95[th] percentile of each "group by year" cluster. Results about means are summarized in Figure 4 (in the Appendix) where predicted vs observed means are plotted by group and year (straight lines representing 0.95 probability intervals calculated on the basis of the posterior predictive distribution). The fit is good in almost all cases, with the only exception of groups 5 and 8, which however are included in the probability intervals calculated for the posterior predictive distributions.

Similarly, in Figure 5 (in the Appendix) predicted vs observed 95[th] percentiles are plotted, (straight lines represent 0.95 probability intervals for 95[th] percentiles).

Globally considered, the fit is good in this case as well, with exceptions being represented mostly by the smallest "group by year" clusters. The plot regarding group 8 is omitted because of the too small size of this group.

We check also whether the model is adequate in predicting the 95[th] percentiles of annual distributions. To this aim, we combine the eight annual groups to calculate the posterior predictive of the annual 95[th] percentiles. This is rather easy to do because of the nice form of the Weibull distribution survival function. In fact:

$$\Pr(y_{ijh} > z) = \exp(-\lambda_{ij} z^{\delta})$$ (4.2)

For each iteration of the Markov chain the following equation in $z_i$ is solved:

$$\sum_{i=1}^{8} \Pr(y_{ijh} > z_j) n_{ij} = 0.05 \sum_{i=1}^{8} n_{ij}, \quad \text{where } i = 1,...,8; j = 1,...,5.$$ (4.3)

Predicted and observed annual 95[th] percentiles (along with 0.95 posterior predictive intervals) are plotted in Figure 6a (see the Appendix), that shows a satisfying fit for all the years in the considered period.

On the basis of the survival function (4.2) we also compute the annual number of exceedances over the threshold of 180 μg/m$^3$ predicted by the model. This threshold is explicitly considered by the Italian law as a "warning level" for the ozone concentrations. The number of predicted versus actual exceedances are plotted in figure 6b (see the Appendix).

## 5. Model sensitivity

The models of section 3.2 assume that, conditionally on each "group by year" specific quasi-scale parameters, daily maxima are independently distributed. This assumption is consistent with the opinion of Piegorsch *et*

*al.* (1998) who notice that "in the case of ozone it is widely assumed that day to day values are conditionally independent given the meteorology". Nonetheless, since the series of ozone maxima exhibits strong auto-correlation, the assumption of conditional independence deserves further analysis.

To assess the sensitivity of inferences on long term trends with respect to the assumption of conditional independence, we generalize model 2 by introducing an AR(1) component in equation (3.6). That is, we replace (3.6) by

$$\ln(\lambda) = \mathbf{1}\alpha + X\beta + Z\gamma + U\varepsilon \qquad (5.1)$$

where $X$ and $Z$ are defined according to (3.7) and (3.8); the design matrix $U$ is a diagonal block one defined as

$$U = diagblock\left(U_1, ..., U_P\right)$$

where each $U_p$ is a $N_p \times N_p$ ( $N_p$ being the number of maxima in year $p$) whose generic element can be described as

$$\left(u_{ij(p)}\right) = \begin{cases} 0 & \text{if } j > i \\ \rho^{|i-j|} & \text{otherwise} \end{cases}$$

We leave the prior assumptions of (3.12) unchanged and specify the following distribution for ε:

$$\varepsilon \mid \omega \sim N\left(0, \omega I\right) \qquad (5.2)$$

As regards hyper-parameters we suppose that

$$\omega \sim Gamma(0.1, 0.1)$$
$$\rho \sim Uniform\left(-1, 1\right) \qquad (5.3)$$

We refer to this model as model 4. The introduction of the autoregressive component influences both the quasi-scale and the shape parameters.

15

To simplify the comparison between model 2 and 4, in the latter we set the shape parameter $\zeta$ equal to the value of its posterior mean in model 2. The posterior means and standard deviations of parameters of model 2 and 4 are listed in Table 3, where it is evident that all parameters, and in particular those associated to year effects, are very close in the two cases.

*Table 3*
*Posterior means and standard deviations of parameters from model 2 and model 4*

|  | Model 2 | | Model 4 | |
|---|---|---|---|---|
|  | *posterior mean* | *posterior st.dev.* | *posterior mean* | *posterior st.dev.* |
| $\alpha$ | -15,820 | 0,398 | -15,510 | 0,280 |
| $\gamma_2$ | -0,753 | 0,119 | -0,798 | 0,161 |
| $\gamma_3$ | 1,097 | 0,145 | 1,107 | 0,197 |
| $\gamma_4$ | 0,058 | 0,176 | -0,015 | 0,242 |
| $\gamma_5$ | -1,418 | 0,235 | -1,325 | 0,324 |
| $\gamma_6$ | -1,549 | 0,120 | -1,658 | 0,160 |
| $\gamma_7$ | -2,517 | 0,143 | -2,523 | 0,193 |
| $\gamma_8$ | -3,522 | 0,303 | -3,444 | 0,434 |
| $\delta_2$ | **-0,640** | **0,100** | **-0,711** | **0,174** |
| $\delta_3$ | **-1,238** | **0,105** | **-1,301** | **0,169** |
| $\delta_4$ | **-1,099** | **0,104** | **-1,128** | **0,164** |
| $\delta_5$ | **-0,963** | **0,102** | **-1,093** | **0,172** |
| $\zeta$ | 3,748 | 0,087 | 3,748 | |
| $\rho$ | | | 0,410 | 0,045 |

We note that the assumption of (3.10) as a prior distribution for $\zeta$ leads to similar conclusions, with the exception that parameters associated to random effects are expressed according to different scales and therefore are more difficult to compare.

## 6. Conclusions

In this work we propose a method for detecting ozone long term trend that can be seen as a development of the stratified approach of Huang and Smith (1999) and a special case of the regression model introduced by Cox and Chu (1993) that is still used as reference method for trends assessment of ozone concentrations by the US EPA (see Thompson *et al.*, 2000). We adopt a Bayesian viewpoint and base posterior summaries on McMC samples. We give special emphasis to testing for evidence in favour of separate trends at different process levels.

We apply our method to real data from a single monitoring station in Bologna, Italy, over the period 1994-1998. The analysis of these data does not provide sound evidence in favour of specific trends after cancelling out the effect of weather. As this may depend on the fact that we dispose of few observations for the study of trends at the highest level of the process, we can conclude that, at present, a common trend assumption provides the better description of long term time pattern of ozone daily maxima.

In particular our analysis highlight that there is a strong growth in standardized ozone concentrations from 1994 to 1996, while in the following two years they seems not to show relevant year to year variations.
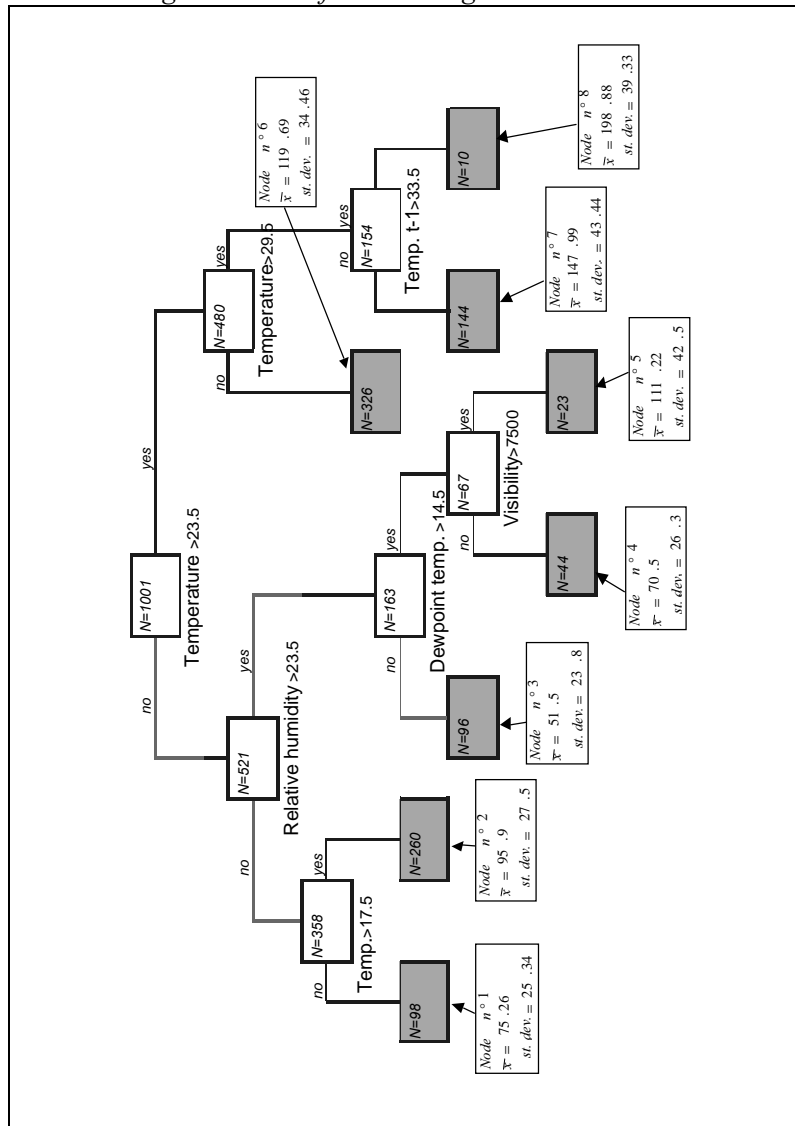
## References

P. Bloomfield, J.A. Royle L.J. Steinberg, Q. Yang, (1996) "Accounting for meteorological effects in measuring urban ozone levels and trends", *Atmospheric Environment*, n. 30, 3067-3077.

L. Breiman, J.H. Friedman, R.A. Olshen, C.J. Stone (1984) *Classification and regression trees*, Wadsworth, California.

W.R. Burrows, M. Benjamin, S. Beauchamp, E.R. Lord, D. Mccollor, B. Thompson (1995) "CART decision tree statistical analysis and prediction of summer season maximum surface ozone for the Vancouver, Montreal and Atlantic Regions of Canada", *Journal of Applied Meteorology*, n. 34, 1848-1862.

S.P. Brooks, A. Gelman (1998) "Alternative methods for monitoring convergence of iterative simulation", *Journal of Computational and Graphical Statistics*, n. 7, 434-455.

L.A. Clark, D. Pregibon (1991) "Tree-based Models", in *Statistical Models in S*, J.M. Chambers and T. J. Hastie (editors), Wadsworth, 377-420.

W.M. Cox, S.H. Chu (1993) "Meteorologically adjusted trends in urban areas: a probabilistic approach", *Atmospheric Environment,* n. 27B, 425-434.

A. Gelman, D.B. Rubin (1992) "Inference from iterative simulation using multiple sequence", *Statistical Science*, n. 7, 457-72.

L.S. Huang, R.L. Smith (1999) "Meteorologically-dependent trends in urban ozone", *Environmetrics*, n. 10, 103-118.

J. Picklands (1971) "The two dimensional Poisson process and extremal processes", *Journal of Applied Probability*, n. 8, 745-756.

W.W. Piegorsch, E.P. Smith., D. Edwards, R.L. Smith (1998) "Statistical advances in environmental science", *Statistical Science*, n. 13, 186-208.

G. Schwartz (1978), "Estimating the dimension of a model", *Annals of Statistics*, 6, 461-64.

R.L .Smith (1989) "Extreme value analysis of environmental time series: an application to trend detection in ground-level ozone", *Statistical Science*, 4, 367-393.

D.J. Spiegelhalter, A. Thomas, N. Best, W.R. Gilks (1996) *BUGS: Bayesian Inference Using Gibbs Sampling, version 0.50,* Technical

Report, Medical Research Council Biostatistic Unit, Institute of Public Health, Cambridge University.
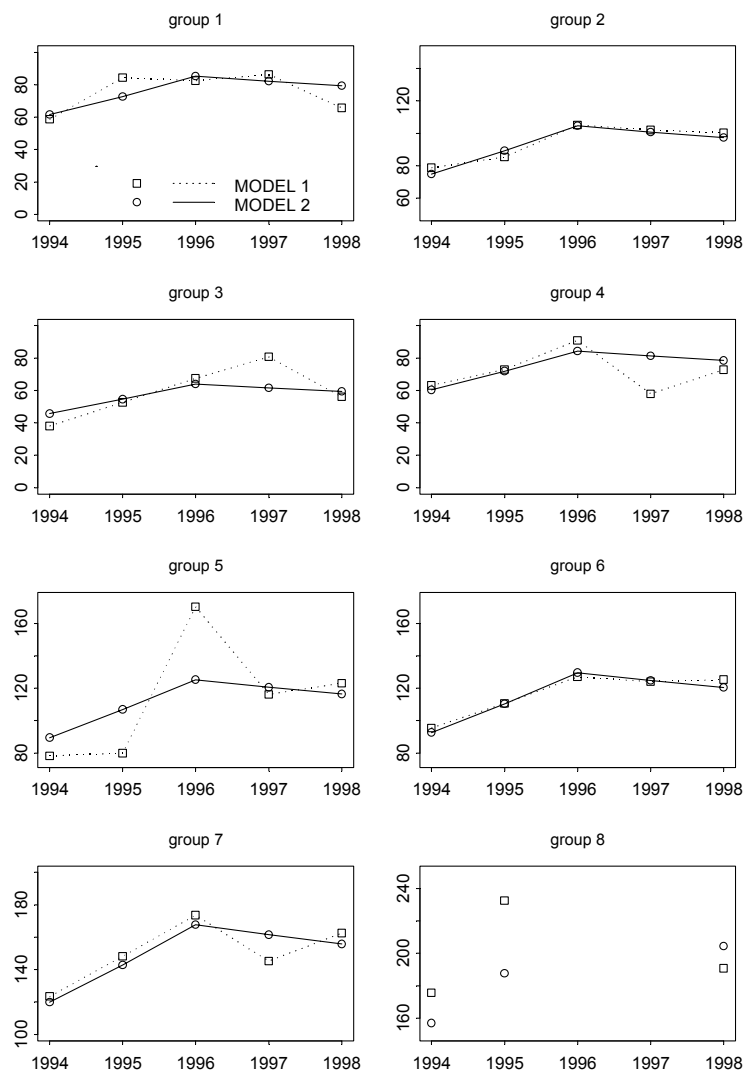
M.L. Thompson, J.Reynolds, L.H. Cox, P. Guttorp, P.D. Sampson (2000) *A review of statistical methods for the meteorological adjustment of tropospheric ozone (revised),* Technical Report Series n. 26, NRCSE-TRS.
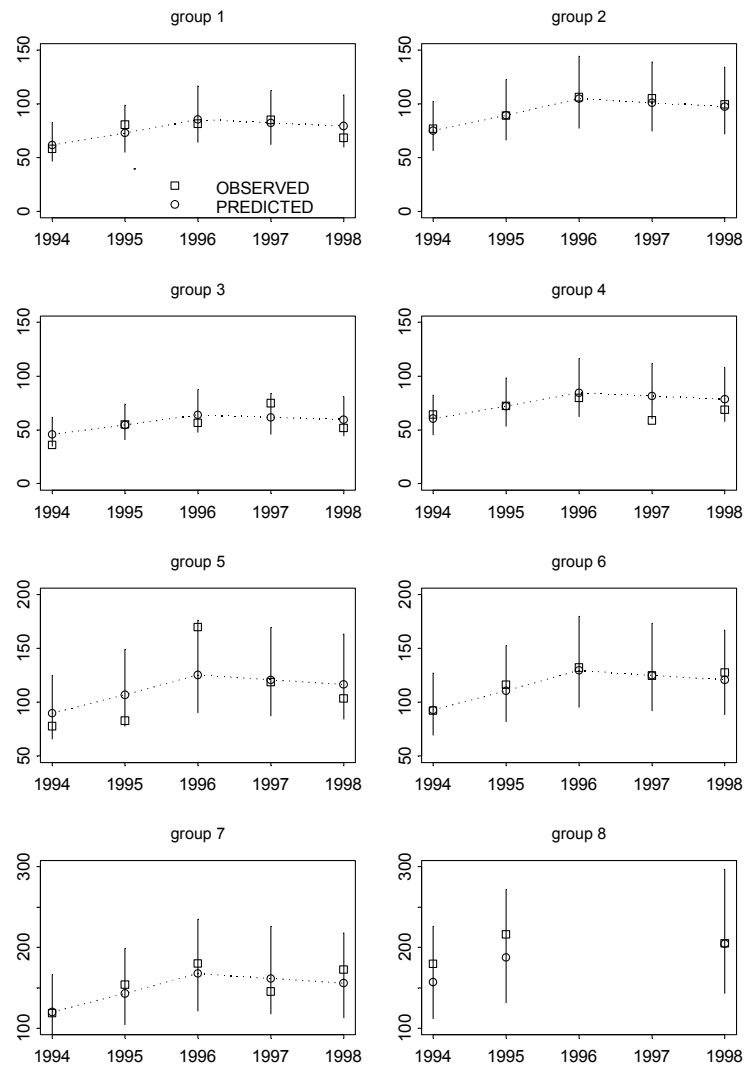
**Appendix**

*Figure 2*
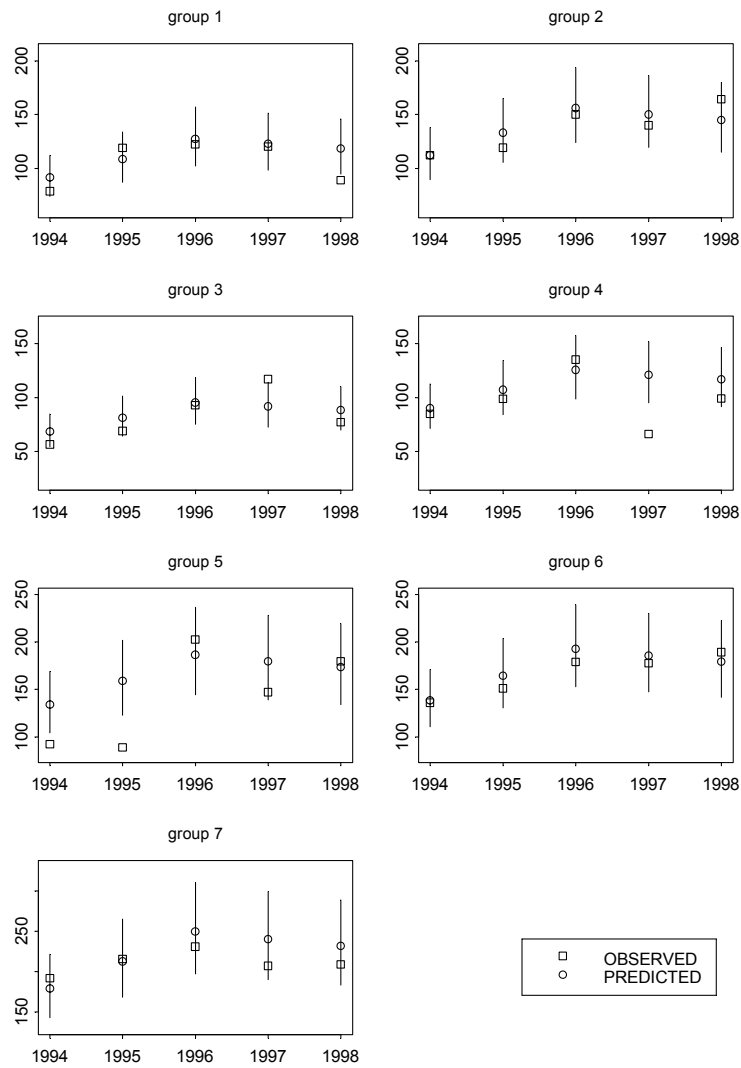*Regression tree for the Bologna ozone maxima*

*Figure 3*
*Posterior means under model 1 and model 2, by group and year*

*Figure 4*
*Predicted vs observed means by group and year and their 0.95*
*probability intervals*

*Figure 5*
*Predicted vs observed 95$^{th}$ percentiles by group and year and their 0.95
probability intervals*

*Figure 6*
*Predicted vs observed 95th percentiles and exceedances over the*
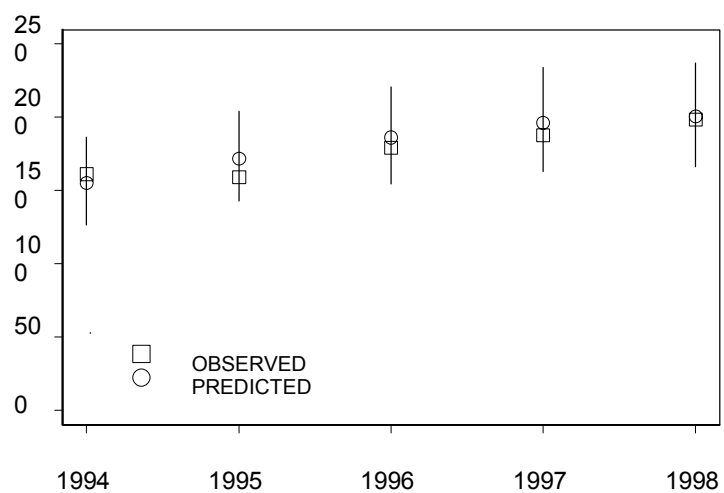*threshold of 180 μg/m3 by year*

Figure 6a



Figure 6b