



Deposited via The University of Leeds.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/43634/>

Version: Accepted Version

---

**Article:**

Gyarmati-Szabo, J, Bogachev, L and Chen, H (2011) Modelling threshold exceedances of air pollution concentrations via non-homogeneous Poisson process with multiple change-points. *Atmospheric Environment*, 45 (31). 5493 - 5503 . ISSN: 1352-2310

<https://doi.org/10.1016/j.atmosenv.2011.06.049>

---

© 2011, Elsevier. Licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International <http://creativecommons.org/licenses/by-nc-nd/4.0/>

**Reuse**

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

# Modelling threshold exceedances of air pollution concentrations via non-homogeneous Poisson process with multiple change-points

János Gyarmati-Szabó<sup>a,\*</sup>, Leonid V. Bogachev<sup>b</sup>, Haibo Chen<sup>a</sup>

<sup>a</sup> *Institute for Transport Studies, University of Leeds, Leeds LS2 9JT, UK*

<sup>b</sup> *Department of Statistics, University of Leeds, Leeds LS2 9JT, UK*

## ABSTRACT

A Bayesian multiple change-point model is proposed to analyse violations of air quality standards by pollutants such as nitrogen oxides (NO<sub>2</sub> and NO) and carbon monoxide (CO). Threshold exceedance occurrences are modelled by a step rate Poisson process fitted after short-range correlations in the exceedance data are removed via declusterisation. The change-points are identified, and the rate function is estimated, using an MCMC algorithm adapted from Green (1995). This technique is applied to the daily concentration data collected in Leeds, UK (1993–2009). Results are validated by running the MCMC estimator on the posterior-replicated data. Findings are discussed in the context of the past environmental actions and events. The proposed methodology may be useful for the air quality management by providing quantitative means to measure the efficacy of pollution control programmes.

*Keywords:* Air quality action plan; Urban air pollution; Threshold exceedances; Non-homogeneous Poisson process; Change-points; MCMC

## 1. Introduction

Due to harmful effects of pollution on human health and the environment, air quality in urban areas is a matter of worldwide concern amongst scientists, policy makers and public alike. Significant health problems (e.g., respiratory and cardiovascular) can be caused or worsened by exposure to the air pollution, from mild effects such as increased use of inhalers by asthmatics to more serious including advancement of death (Brimblecombe, 1998; Ayres, 2002).

The UK National Air Quality Strategy, adopted in 1997 (see AQS, 2007), determined statutory standards for eight main air pollutants, including nitrogen dioxide (NO<sub>2</sub>), nitrogen monoxide (NO) and carbon monoxide (CO). In particular, UK local authorities are required to identify Air Quality Management Areas (AQMA) where appropriate, along with Air Quality Action Plans in pursuit of the air quality objectives (AQMA, 1997).

It is well documented that the pollutant emission to the atmosphere is dominated by road transport, estimated to contribute 54% of CO and 32% of NO<sub>x</sub> emissions in

---

\* Corresponding author. Tel.: +44 (0)113 343 1788; fax: +44 (0)113 343 5334.  
*E-mail address:* tsjs@leeds.ac.uk (János Gyarmati-Szabó).

1 the UK (NAEI, 2010). In zones close to roadside, road traffic may be responsible for  
2 over 90% of the total emission of harmful pollutants such as NO<sub>2</sub> (AQAP, 2004). The  
3 UK Traffic Management Act 2004 (TMA, 2004) has put in place legislation aiming to  
4 promote a systematic integration of air quality management with transport planning.  
5 As a result of these measures, including tough regulations for industry and tightened  
6 emission standards for vehicles (EUR, 2003), an improved air quality has been  
7 achieved in the UK, even though elevated levels of pollution still occur (AQS, 2007;  
8 Faulkner and Russell, 2010).

9 Leeds City Council is taking a proactive stance in monitoring and managing air  
10 quality by deploying the state-of-the-art instrumentations in the city (with the data ac-  
11 cessible to transportation and environmental researchers and modellers) and by im-  
12 plementing traffic control schemes. Two AQMAs were declared in small areas of the  
13 inner city (AQAP, 2004), with the action plans including measures such as traffic de-  
14 mand management, improvements to the highways network, reducing emissions  
15 from industrial and domestic sources, and raising awareness.

16 In order to monitor and assess the efficacy of these and future policies, it is impor-  
17 tant to develop adequate statistical methods to measure the impact of the regulations  
18 on the dynamics of various pollutants, especially with regard to the set standards.  
19 Clearly, any such methods should be linked to exceedance episodes, when the con-  
20 centrations of one or more pollutants violate certain thresholds. This underpins the  
21 importance of developing reliable tools for statistical modelling of exceedances, in-  
22 cluding their prediction and control. Progress in this direction would provide the nec-  
23 essary feedback helping to assess and improve the current environmental policies,  
24 thus avoiding a grossly inefficient and costly “trial and error” approach. Furthermore,  
25 in addition to the environmental and health benefits, advancement of research meth-  
26 odologies in the statistical modelling of pollutant concentration extremes could have a  
27 significant economic value, potentially leading to massive savings in the public  
28 spending related to the government’s environmental policies.

29 To date, a number of statistical methods have been developed to model violations  
30 of the air quality standards. The common approach is based on extreme value theory  
31 (see review papers by Thompson et al., 2001; Smith, 2004, and further references  
32 therein). Exceedances of air pollution concentrations can also be modelled using  
33 non-homogeneous Poisson processes (see Raftery, 1989; Smith, 1989; Shively,  
34 1991; Smith and Shively, 1995; Achcar et al., 2009). This approach is motivated by  
35 the asymptotic theory of stationary point processes stating that the series of occur-  
36 rences of high threshold exceedances can be approximated by a Poisson process  
37 (Leadbetter et al., 1983). For nonstationary data, the approximating Poisson process  
38 has to be non-homogeneous, with the intensity rate depending on time. Estimation of  
39 the unknown rate function  $\lambda(t)$  is greatly facilitated by choosing certain parametric  
40 classes, such as  $\lambda(t) = \alpha t^{\beta-1}$ ,  $\lambda(t) = \beta/(t + \alpha)$ ,  $\lambda(t) = \alpha e^{-\beta t}$  ( $\alpha, \beta > 0$ ), etc. (Achcar  
41 et al., 2009). Note that these and other conventional formulas are restricted by a  
42 monotone behaviour of  $\lambda(t)$ . More complex models are obtained by incorporating the  
43 dependence of parameters on “external” factors (covariates) such as meteorological

1 conditions (Smith and Shively, 1995). However, such models are not suitable if the  
 2 rate may change abruptly (e.g., following an updated air quality action plan). To cap-  
 3 ture this type of variation in the exceedance dynamics, Achcar et al. (2009) proposed  
 4 a parametric non-homogeneous Poisson model allowing to handle up to three  
 5 change-points using complicated iterative updates of the prior information.

6 In the present paper, we use a non-homogeneous Poisson model with no prior re-  
 7 striction on the number of change-points. For simplicity, the rate  $\lambda(t)$  is assumed to  
 8 be a step function (i.e., constant between adjacent change-points). Note that statis-  
 9 tics of threshold exceedances is likely to involve short-range correlations due to  
 10 patches of elevated concentration levels. Therefore, to ensure a better fit of the Pois-  
 11 son model it is essential to apply a suitable thinning (declusterisation) of the thresh-  
 12 old exceedance data (Shively, 1991). An efficient MCMC estimator is then used to  
 13 automatically determine the number and locations of change-points, as well as the  
 14 respective values (heights) of the rate function. These techniques are applied to NO<sub>2</sub>,  
 15 NO and CO data collected in Leeds, UK, as daily concentration maxima over nearly  
 16 17 years (1993–2009).

17 The paper is organised as follows. In Section 2, we review basic properties of  
 18 non-homogeneous Poisson processes and set out the step rate model; formal statis-  
 19 tical tests to compare different Poisson models are also outlined there. Section 3  
 20 deals with the model specification and preprocessing, including thresholding, imput-  
 21 ing missing values, deseasonalisation and declusterisation. The results of fitting a  
 22 step rate Poisson model to the NO<sub>2</sub>, NO and CO threshold exceedance data, using a  
 23 suitable MCMC estimator, are reported in Section 4, followed by the model verifica-  
 24 tion and validation. Section 5 contains the discussion, including some preliminary in-  
 25 terpretation of the results in terms of past air quality related actions and events. Fi-  
 26 nally, a brief summary of the work is given in Section 6.

## 27 **2. Poisson model**

### 28 *2.1. Non-homogeneous Poisson process*

29 Poisson process (Cox and Lewis, 1966) is a statistical model describing a series  
 30 of uncorrelated random events (“arrivals”), whereby new arrivals are statistically in-  
 31 dependent of any past arrivals. Rate function  $\lambda(t)$  of the Poisson process determines  
 32 its local intensity. The number  $N(t)$  of arrivals on  $[0, t]$  is distributed according to a  
 33 Poisson law,

$$34 \quad P\{N(t) = n\} = \frac{\Lambda(t)^n}{n!} e^{-\Lambda(t)}, \quad n = 0, 1, 2, \dots, \quad (1)$$

35 where  $\Lambda(t) = \int_0^t \lambda(u) du$  is the expected number of arrivals. More generally, the in-  
 36 crements  $N(t) - N(s)$  ( $t > s$ ) have Poisson distribution with parameter  $\Lambda(t) - \Lambda(s)$ .  
 37 The waiting time  $\tau$  to first arrival after time  $s$  satisfies

1  $P\{\tau > t - s\} = P\{N(t) - N(s) = 0\} = e^{-\Lambda(t) + \Lambda(s)} \quad (t > s).$

2 Hence, the joint density of  $n$  consecutive arrivals restricted to  $[0, T]$  reads

3  $f(t_1, \dots, t_n) = e^{-\Lambda(T) + \Lambda(t_n)} \prod_{i=1}^n \lambda(t_i) e^{-\Lambda(t_i) + \Lambda(t_{i-1})} = e^{-\Lambda(T)} \prod_{i=1}^n \lambda(t_i) \quad (2)$

4  $(0 = t_0 < t_1 < \dots < t_n < T)$

5 with the corresponding log-likelihood

6  $\log L(t_1, \dots, t_n) \equiv \log f(t_1, \dots, t_n) = \sum_{i=1}^n \log \lambda(t_i) - \Lambda(T) \quad (0 < t_1 < \dots < t_n < T). \quad (3)$

7 Formulas (1), (2) imply that the joint probability density of arrivals conditioned on their  
8 total number  $N(T) = n$  is given by

9  $f_n(t_1, \dots, t_n) = \frac{n!}{\Lambda(T)^n} \prod_{i=1}^n \lambda(t_i) \quad (0 < t_1 < \dots < t_n < T), \quad (4)$

10 which coincides with the density of the increasing order statistics for  $n$  independent  
11 points thrown in the interval  $[0, T]$  with the density  $\lambda(t)/\Lambda(T)$  each (e.g., for a con-  
12 stant rate this amounts to the uniform distribution on  $[0, T]$ ). In turn, formula (4) im-  
13 plies that the values  $t'_i = \Lambda(t_i)/\Lambda(T)$  have the same distribution as the uniform order  
14 statistics on  $[0, T]$  (Lewis and Shedler, 1976), which allows one to construct approxi-  
15 mate goodness-of-fit tests for particular classes of the rate functions using estimates  
16 of the model parameters (see Section 2.2 below).

17 In the present paper, the rate  $\lambda(t)$  is assumed to be a step function,

18  $\lambda(t) = h_j \quad \text{if} \quad s_j \leq t < s_{j+1} \quad j = 0, 1, \dots, k \quad (0 = s_0 < s_1 < \dots < s_k < s_{k+1} = T), \quad (5)$

19 where the number  $k$  of change-points, as well as their positions  $s_j$  and heights  $h_j$  of  
20 the corresponding steps are not predetermined and therefore have to be statistically  
21 estimated.

## 22 2.2. Statistical tests for the Poisson model

23 To assess the utility of the change-point Poisson process as a fitting model, one  
24 can carry out statistical tests to compare three different models:

- 25 •  $M_0$ : homogeneous Poisson process,  $\lambda(t) \equiv \text{const}$ ;
- 26 •  $M_1$ : Poisson process with a log-linear rate,  $\lambda(t) = \alpha \exp(-\beta t)$ ;
- 27 •  $M_2$ : Poisson step rate model with one change-point.

28 For observed threshold exceedance data  $X : 0 < t_1 < t_2 < \dots < t_n < T$ , denote

$$S_n = \sum_{i=1}^n u_i, \quad u_i = t_i/T \quad (i = 1, \dots, n). \quad (6)$$

A classical goodness-of-fit test for  $M_0$  is based on the uniform distribution of arrivals (Section 2.1), implying that  $S_n$  is approximately normal with mean  $n/2$  and variance  $n/12$  (Cox and Lewis, 1966). A test for  $M_0$  against  $M_2$  (Raftery, 1989) is based on the statistic (Raftery, 1989)

$$\Delta_n = \frac{1}{\sqrt{n}} \max\{|g(i-1, u_i)|, |g(i, u_i)|\}, \quad (7)$$

where the maximum is taken over all  $u_i \in [0.01, 0.99]$ , and

$$g(i, u) = i\sqrt{(1-u)/u} - (n-i)\sqrt{u/(1-u)} \quad (i = 1, \dots, n, \quad 0 < u < 1).$$

Approximate critical values can be computed from the formula (Akman and Raftery, 1986)

$$P(\Delta_n > z) \approx \sqrt{\frac{2}{\pi}} e^{-z^2/2} (\xi z - \xi z^{-1} + z^{-1}), \quad \xi = \log \frac{0.99}{0.01} = 4.595 \dots \quad (8)$$

For the model  $M_1$ , the coefficient  $\alpha$  is a nuisance parameter eliminated by switching to the conditional distribution with  $n > 0$  observed arrivals in  $[0, T]$ . The slope parameter  $\beta$  is estimated from the maximum likelihood equation  $\frac{\partial \log L_n}{\partial \beta}(\hat{\beta}) = 0$ , in view of (3) reduced to (Cox and Lewis, 1966)

$$\frac{1}{\hat{\beta}T} - \frac{e^{-\hat{\beta}T}}{1 - e^{-\hat{\beta}T}} = \frac{S_n}{n}. \quad (9)$$

As was explained in Section 2.1, the modified values

$$t'_1 = \frac{\hat{\Lambda}(t_1)}{\hat{\Lambda}(T)}, \quad \dots, \quad t'_n = \frac{\hat{\Lambda}(t_n)}{\hat{\Lambda}(T)}, \quad (10)$$

with  $\hat{\Lambda}(t) = \int_0^t \hat{\lambda}(u) du = \alpha \hat{\beta}^{-1} (1 - \exp(-\hat{\beta}t))$ , are approximately distributed as the uniform order statistics, so the uniformity test described above may be applied (Lewis and Shedler, 1976).

Models can also be compared using the Bayes factor  $B_{ij}$  (Raftery, 1996), defined as the ratio of posterior to prior odds for  $M_i$  against  $M_j$ :

$$B_{ij} = \frac{P(X | M_i)}{P(X | M_j)}, \quad X : t_1 < \dots < t_n. \quad (11)$$

1 An indicative scale for interpreting  $B_{ij}$  is given in Raftery (1996, p. 165); e.g., values  
 2 from 12 to 150 are classed as strong evidence in favour of the null hypothesis  $M_j$ .  
 3 The Bayes factors for  $M_0$  against  $M_1$  and  $M_2$  (Raftery, 1989) are given by

$$4 \quad B_{01} = 0.6449(n-1) \left[ \int_0^\infty e^{-S_n y} \left( \frac{y}{1-e^{-y}} \right)^{n-1} dy \right]^{-1}, \quad (12)$$

$$5 \quad B_{02} = \frac{4\sqrt{\pi} \Gamma(n + \frac{1}{2})}{\sum_{i=0}^n J_i \Gamma(i + \frac{1}{2}) \Gamma(n - i + \frac{1}{2})}, \quad (13)$$

6 where  $\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$  and  $J_i = \int_{u_i}^{u_{i+1}} x^{-(i+1/2)} (1-x)^{-(n-i+1/2)} dx$ . Using equation  
 7 (11), one also obtains

$$8 \quad B_{12} = \frac{B_{02}}{B_{01}}. \quad (14)$$

### 9 **3. Model preprocessing**

#### 10 *3.1. Thresholding*

11 To set a specific threshold in order to extract the statistics of exceedances, it  
 12 might seem natural to use the air quality standards (AQS, 2007). However, there are  
 13 several difficulties: (i) the standards are set out in terms of hourly or maximum daily  
 14 rolling 8-hour averages, while we work with the daily data; (ii) the standards limit  
 15 permitted frequencies of threshold violations, which causes the “moving window” type  
 16 dependence in the exceedance data; (iii) the standard threshold values produce  
 17 scarce statistics of the resulting exceedances (e.g., 0.11% for NO<sub>2</sub>).

18 A simple alternative is to use the quantile thresholding at, say, a 90th empirical  
 19 percentile, leading to the specific threshold values of 96  $\mu\text{g m}^{-3}$  (NO<sub>2</sub>), 185  $\mu\text{g m}^{-3}$   
 20 (NO) and 2.1  $\text{mg m}^{-3}$  (CO). Cumulative counts of the resulting exceedance data for  
 21 all three pollutants are shown in Fig. 1.

#### 22 *3.2. Missing values*

23 Our concentration level data contain a noticeable proportion of missing values  
 24 (6.6%, 5.5% and 9.9% in the NO<sub>2</sub>, NO and CO data, respectively), which considera-  
 25 bly decreases the exceedance statistics available and therefore may adversely affect  
 26 the estimation of the unknown rate function  $\lambda(t)$ .

27 One can improve estimation by compensating for the omissions (Cox, 1981). In  
 28 this work, we impute the missing values using a (two-sided) moving average estima-  
 29 tor with the window size  $\pm 65$  days. We first impute raw concentration values drawn  
 30 independently according to their estimated (univariate) distribution, and then apply  
 31 thresholding. Even though such interpolation neglects possible short-range correla-

1 tions in the concentration time series, it may be expected to work well due to a sub-  
2 sequent high-level thresholding (Leadbetter et al., 1983).

3  
4 **<Figure 1>**

5  
6  
7 **<Table 1>**

8  
9  
10 The imputed exceedances for each pollutant are indicated in Fig. 1. Table 1  
11 shows the MCMC results for NO<sub>2</sub> obtained before and after imputation. In both  
12 cases, the number of change-points identified by MCMC was found to be  $k = 1$ , with  
13 almost the same locations; however, the step heights of the rate function noticeably  
14 increased (by about 5%) as a result of more exceedances.

### 15 3.3. Deseasonalisation

16 It has to be taken into account that the concentration levels of all three pollutants  
17 under study (NO<sub>2</sub>, NO and CO) are affected by meteorological conditions, e.g. induc-  
18 ing yearly oscillations (seasonality) in time series. For nitrogen oxides NO<sub>2</sub> and NO,  
19 this may be explained by photochemical reactions between NO<sub>2</sub> and NO involving  
20 ozone O<sub>3</sub> (Clapp and Jenkin, 2001), leading to decreased levels of NO<sub>x</sub> during the  
21 summer due to higher solar radiation. Carbon monoxide CO is relatively nonreactive,  
22 but low temperatures in winter contribute to higher CO concentrations due to incom-  
23 plete combustion in vehicle engines (CCME, 2002). Yearly cycles in the concentra-  
24 tion time series are clearly visible in the daily data plots (Fig. 2, top), confirmed by  
25 computing the autocorrelation at lag 365 days (0.17 for NO<sub>2</sub> and NO, 0.33 for CO).  
26 Furthermore, analyzing the estimated spectral density of the observed time series  
27 (Venables, 2002), seasonality was found to be significant.

28 Deseasonalisation may be based on fitting a regression model for the log-  
29 transformed data (Cox and Lewis, 1966)

$$30 \quad \log \tilde{X}(t) = \log X(t) - a \cos \omega t - b \sin \omega t - c, \quad \omega = \frac{2\pi}{365} (\text{days}^{-1}), \quad (15)$$

31 Possible deviations from stationarity in the concentration data are likely to occur on a  
32 scale much slower than the annual variability, hence the regression estimation may  
33 be expected to give satisfactory results (cf. Cox and Lewis, 1966). Periodic trends  
34 filtered out from the raw data are shown in Fig. 2 (top). Deseasonalised data  $\tilde{X}(t)$   
35 are plotted in Fig. 2 (bottom) showing a much improved scattering, with no visible pe-  
36 riodic pattern

37  
38 **<Figure 2>**

1 Alternative methods of deseasonalisation may be based on moving-average, ker-  
2 nel type local smoothing or spectral techniques (Cox and Lewis, 1966; Lewis and  
3 Shedler, 1976). One could also take into account weekly periodicities due to traffic  
4 patterns caused by peak hours and weekday–weekend oscillations. However, weekly  
5 oscillations occur much more frequently as compared to the annual scale, so are bet-  
6 ter averaged out by fitting a step rate Poisson model; in addition, the weekly patterns  
7 are likely to be rather stable across decades.

### 8 *3.4. Declusterisation*

9 The time series of threshold exceedances by air pollution concentrations is likely  
10 to involve dependence (autocorrelation) manifested in clustering of consecutive ex-  
11 ceedances, which may render Poisson model unsuitable. However, thanks to a rela-  
12 tively short range of such correlations, the dependence problem can be rectified by a  
13 suitable thinning, or “declusterisation” (Leadbetter et al., 1983; Shively, 1991), usually  
14 achieved by retaining one point from each cluster (merged with the subsequent point  
15 below the threshold to remove the tie).

16 Significance of dependence before and after declusterisation may be tested using  
17 the sample correlation coefficient for gaps between exceedances or by estimating  
18 serial correlation at various lags (Venables, 2002). We deployed the nonparametric  
19 Wald–Wolfowitz runs test of independence (Gibbons and Chakraborti, 2003) applied  
20 to a two-valued data sequence representing threshold exceedances versus non-  
21 exceedances. Note that the runs test applied to the entire declustered sequence may  
22 still reject  $H_0$  due to nonstationarity, so it should be used on separate intervals be-  
23 tween the change-points, which need to be estimated in advance.

24 In this study, declusterisation resulted in a reduction of the data size by about 30–  
25 35%. The runs test accepted independence on all intervals between the estimated  
26 change-points, with the p-values above 0.10.

## 27 **4. Results for the air pollution data**

28 In this section, we report the results of estimation of the unknown step rate func-  
29 tion  $\lambda(t)$  to fit a non-homogeneous Poisson process to threshold exceedance data  
30 obtained from the NO<sub>2</sub>, NO and CO concentration level measurements. The data  
31 were collected in the Leeds city centre by a permanent monitoring station located  
32 about 30 metres from a busy 4-lane inner-city road and 150 metres from an urban  
33 motorway, with traffic flow of approximately 21,500 and 93,500 vehicles per day, re-  
34 spectively (AURN, 2010). The data used in the study (AQA, 2010) correspond to the  
35 daily maxima for each of the three pollutants, observed from 4 January 1993 to 31  
36 December 2009 (i.e., for  $T = 6206$  days).

### 37 *4.1. Posterior estimates of the rate function*

38 From the cumulative plot of exceedance frequencies for NO<sub>2</sub> (Fig. 1), it is evident  
39 that the initial part of the graph up to  $t \approx 2500$  is close to a linear function (with the  
40 slope about  $280/2500 = 0.112$ ), followed by a noticeable drop of the slope to ap-

1 proximately  $(300 - 280)/(3600 - 2500) = 0.018$  until  $t \approx 3600$ . One may also suspect  
2 one or two less significant change-points further on (e.g., near  $t \approx 4500$ ). Similarly,  
3 the cumulative graphs for NO and CO in Fig. 1 suggest about four to five hypothetical  
4 change-points each. Hence, as a prior distribution for the unknown number  $k$  of  
5 change-points it is reasonable to choose a Poisson law with mean 4.5.

6 The programming code for the Markov chain Monte Carlo (MCMC) algorithm  
7 adapted from Green (1995) was implemented in R using a standard laptop (AMD  
8 processor DualCore 2.2 GB, RAM 4GB). The sampler takes 1000 steps within 20  
9 seconds. Manual programming was necessary since the standard MCMC engines  
10 (e.g., WinBUGS, <http://www.mrc-bsu.cam.ac.uk/bugs/>) are not yet suitable for models  
11 with variable dimension.

12 After a burn-in period of length 15,000–20,000, pilot runs established that the dy-  
13 namics of the Markov chain have reached the stationary regime. As usual, conver-  
14 gence of the sampling algorithm was monitored by visual inspection of the output  
15 plots and also by simple convergence diagnostic tools based on the MCMC output  
16 for a single Markov chain (Gilks et al., 1996). In addition, reliability of the burn-in runs  
17 was verified by trying different initial values, also utilised to control the termination  
18 time of the MCMC algorithm. In the stationary regime, the MCMC was run for  
19 500,000 updates. Examining the autocorrelation plots, it was determined that inde-  
20 pendence within the sample is ensured by selecting every 40th generated value, thus  
21 providing a posterior sample of size  $500,000/40 = 12,500$ .

22  
23 **<Figure 3>**  
24  
25

26 As a result of the MCMC performance, it was concluded that the MCMC sampler  
27 spent most of its time in parametric states corresponding to  $k = 1$ ,  $k = 3$  and  $k = 4$   
28 change-points for NO<sub>2</sub>, NO and CO, respectively (Fig. 3). The change-point locations  
29  $s_j$ , estimated via the modes of the posterior marginal densities, are shown in Fig. 4.  
30 Specific time values of the change-points, along with 50% credible intervals, are  
31 given in Table 2. The step heights  $h_j$  of the rate function  $\lambda(t)$  were estimated via the  
32 median of the posterior marginal distributions (see Fig. 5 and Table 3). Note that the  
33 parametric dimension of the model is bounded by

34 
$$\dim(k) + \dim(s_j) + \dim(h_j) \leq 1 + 4 + 5 = 10,$$

35 furnishing  $12,500/10 = 1250$  sampled values per parameter, which ensures good ac-  
36 curacy of the posterior estimation.

37  
38 **<Figure 4>**  
39

40 **<Table 2>**  
41

1  
2 <Figure 5>

3  
4  
5 <Table 3>

6  
7 <Figure 6>

8  
9  
10 The graph of the step function  $\hat{\lambda}(t)$  (with estimated change-points  $s_j$  and step  
11 heights  $h_j$ ) is shown in Fig. 6 along with an “integral” estimate  $\bar{\lambda}(t)$  of the posterior  
12 mean rate function calculated by averaging 5000 sample step rate functions drawn  
13 from the MCMC posterior distribution. The transition areas for  $\bar{\lambda}(t)$  near the change-  
14 points are typically quite narrow, confirming a stepwise nature of the rate function  
15  $\lambda(t)$ . Fig. 6 also shows the cumulative graphs simulated via Poisson processes with  
16 rates  $\hat{\lambda}(t)$  and  $\bar{\lambda}(t)$ , contrasted with the observed cumulative plot (cf. Fig. 1). As  
17 could be anticipated, the integrated (cumulative) estimates demonstrate a much bet-  
18 ter fit to the data.

#### 19 4.2. Verification of the Poisson change-point model

20 The fitted step rate Poisson model can be verified by checking uniform distribution  
21 of exceedance occurrences between the consecutive change-points (Section 2.1).  
22 Specifically, for NO<sub>2</sub> there are two step intervals separated by the estimated change-  
23 point  $s_1 = 2490$  and containing  $n_0 = 257$  and  $n_1 = 133$  exceedances. Using the uni-  
24 formity test based on (6), we obtain the approximate p-values 0.6326 and 0.4907, re-  
25 spectively, confirming good quality of the fitted model. For NO and CO, the results  
26 were similar.

27 Further verification is motivated by challenging the lack of change-points on some  
28 intervals. For instance, to compare two competing models  $M_0$  and  $M_2$  for NO<sub>2</sub> on the  
29 interval  $[2490, 6206]$ , we evaluate the test statistic (7) (with  $n = n_1 = 133$ ) and use  
30 formula (8) to find the approximate p-value 0.1457, indicating that model  $M_0$  is ac-  
31 ceptable. Moreover, formula (13) gives the Bayes factor  $B_{02} = 33.39$ , which consti-  
32 tutes strong evidence in support of  $M_0$ . Similar calculations for NO and CO also con-  
33 firmed the absence of “extra” change-points.

34 It is interesting to assess the stepwise variation of the rate function  $\lambda(t)$  as op-  
35 posed to a more gradual change, especially for NO and CO (Fig. 1), by testing model  
36  $M_1$  against the alternative  $M_2$ . Starting with the NO data, from equation (9) we esti-  
37 mated the slope  $\hat{\beta} = 0.000313$  and tested the uniformity of the modified ex-  
38 ceedances (10). The statistic (6) gave the p-value 0.025, suggesting that  $M_1$  should  
39 be rejected. This was confirmed by calculating the Bayes factor  $B_{21} = 53.17$  using

1 equations (12), (13) and (14). Similar results were obtained for CO, confirming that  
2 model  $M_1$  is less suitable than the fitted step rate model.

### 3 4.3. *Validation of the fitted Poisson model*

4 To validate the employed MCMC sampler, we simulated data sets from non-  
5 homogeneous Poisson processes with the estimated step function  $\hat{\lambda}(t)$  as a rate  
6 function and again applied MCMC to see if there were any significant differences in  
7 the estimates. The results (Fig. 7) show an excellent match of the posterior estimates  
8 for the simulated data sets with those obtained earlier for the observed data.

9

10 <Figure 7>

11

12 <Figure 8>

13

14

15 Following the method of *posterior predictive simulation* (Gilks et al., 1996), repli-  
16 cated the counting processes and compared them to the corresponding cumulative  
17 plots (Fig. 8). Due to random fluctuations, deviations of the replicated cumulative  
18 plots increase over time. The accuracy may be improved by simulating Poisson  
19 processes conditioned to have the same number of occurrences as in the real data.  
20 Conditional replicates are presented in Fig. 8, showing an excellent fit with the ob-  
21 served plot and also consistent with the sample mean plot for data sets simulated  
22 with a posterior estimate  $\hat{\lambda}(t)$ .

## 23 5. Discussion

### 24 5.1. *General comments*

25 Statistical analysis of observed pollutant concentrations developed in this paper is  
26 essentially based on the extreme value approach by utilising high threshold ex-  
27 ceedances for statistical inference. We focused on the corresponding point process  
28 of exceedance occurrences in time, whereas the actual values of exceedances were  
29 not taken into consideration. This simplified model has an advantage of being less  
30 sensitive to deviations from the model assumptions, and hence more useful in real  
31 data applications. The algorithmic tractability of the model makes it suitable for an ef-  
32 ficient MCMC multi-parametric estimation. As has been demonstrated in this work,  
33 MCMC provides a reliable identification of the unknown change-points in the rate  
34 function, including their number, location and the corresponding step heights.

35 A global downward trend in the behaviour of the rate functions for each pollutant  
36 (Fig. 6) indicates a general improvement in the air quality, arguably promoted by en-  
37 vironmental actions (NAEI, 2010). It may also be attributable to the improvement of  
38 car engines (in particular, using the three-way catalytic converters as standard) and  
39 to the developing change in drivers' attitudes (Glaister, 2002). Since a step rate

1 model is unsuitable for capturing this kind of dynamics, it is interesting to consider a  
2 more general class

$$3 \quad \lambda_{\beta}(t) = e^{-\beta t} \lambda(t), \quad (16)$$

4 where  $\lambda(t)$  is a pure step function (5). Model (16) can be reduced to (5) by filtering  
5 out the log-linear term together with annual periodicity, using the regression equation  
6 (cf. (15))

$$7 \quad \log \tilde{X}(t) = \log X(t) + \beta t - a \cos \omega t - b \sin \omega t - c. \quad (17)$$

8 Fig. 9 illustrates the results of fitting the combined trend (17) to the raw data; note a  
9 better fit as compared to a pure periodic trend (cf. Fig. 2).

10 Due to the computational complexity of the model and the distribution-oriented na-  
11 ture of Bayesian inference, specific conclusions of the MCMC-based analysis should  
12 be treated with caution. There is a need for an additional cross-examination and veri-  
13 fication of these findings in co-operation with experts from other subject areas, in-  
14 cluding policy makers as well as environmental and transport modellers and practi-  
15 tioners. Special care is required when relating the conclusions of the threshold ex-  
16 ceedance analysis with official statistics (see, e.g., NAEI, 2010; Faulkner and Rus-  
17 sell, 2010) usually based on averaged observations rather than on extreme values.  
18 This underpins possible significant discrepancies in the interpretation, since extreme  
19 value statistics may appear more sensitive to traffic pollution changes than the typi-  
20 cal, bulk values (Smith, 2004).

21  
22

23 <Figure 9>

24

## 25 5.2. Mapping and interpretation

26 The following is a preliminary attempt of “mapping” of our results, i.e., relating the  
27 identified change-points with real transport and environment related circumstances.

28 5.2.1. *High Occupancy Vehicle lane.* Two change-points in Table 2 (29 October  
29 1999, NO<sub>2</sub>, and 11 November 1999, NO) may be associated with the High Occu-  
30 pancy Vehicle (HOV) lane scheme, designed to give priority to vehicles with more  
31 than one commuter (HOVL, 2004). The HOV lane in Leeds was opened on 11 May  
32 1998 on the A647 Stanningley Road and Stanningley By-Pass, which form the prin-  
33 cipal radial route to the west of the city centre. By the official assessment (HOVL,  
34 2004), a relatively small improvement of air quality was attributed “to reduced vehicle  
35 emissions rather than to the impact of the HOV lane”. However, the reduced conges-  
36 tion on A647 could also affect the traffic in the city centre where our data were col-  
37 lected (AQA, 2010). In addition, due to predominantly westerly winds in North East  
38 England (MET-O, 2010), an improved air quality near A647 might positively contrib-  
39 ute to the city centre (5 miles east of the HOV lane).

1       5.2.2. *East Leeds Link Road*. Leeds City Council initiated a number of projects to  
2 improve the highway network, including the East Leeds Link Road (ELLR) designed  
3 to reduce through-traffic in the city centre (ELLR, 2010). The new 4 km dual car-  
4 riageway link between junction 45 of the M1 and the Leeds Inner Ring Road was  
5 opened on 10 February 2009, expected to provide “indirect” improvements in the air  
6 quality, with low to moderate impact (AQAP, 2004). Our model has not detected any  
7 change-points around February 2009 or after, but using more data beyond December  
8 2009 might make it different.

9       5.2.3. *Leeds Free City Bus*. This zero-fare service (3.4 miles long loop in the city  
10 centre) began on 30 January 2006. It is estimated that the scheme has reduced the  
11 numbers of cars entering the city centre (LCCA, 2010) and hence contributed to-  
12 wards meeting a target of 10% reduction by 2011 in NO<sub>2</sub> (LCCA, 2010). The  
13 scheme’s impact was not strong enough to be picked up by our model, perhaps be-  
14 ing masked by the general downward trend.

15       5.2.4. *Christmas 2004*. Two change-points in Table 2 (1 December 2004, NO,  
16 and 8 January 2005, CO), with overlapping 50%-credible intervals, are both close to  
17 the Christmas and New Year holidays when the traffic drops dramatically. It is un-  
18 clear if this alone has affected the threshold exceedance data so much as to lead to  
19 a statistically significant change-point for NO and CO. Incidentally, note a strong drop  
20 in the raw time series for NO<sub>2</sub> around 16 December 2001 (Fig. 2); however, no  
21 change-points were detected anywhere near.

22       5.2.5. *Construction works*. There were construction works near the monitoring  
23 station in the city centre, from mid-January 2005 until end of June 2009 (AURN,  
24 2010). There were some concerns that the associated transportation and engineering  
25 activities might be responsible for sporadic elevated pollution levels in the area. From  
26 Table 2, the latest identified change-point was 8 December 2005 (for CO), with the  
27 50% credible interval covering the commencement date of the construction. How-  
28 ever, causal association seems unlikely due to a decrease of the estimated rate,  
29 suggesting that any contribution to the air pollution in NO<sub>2</sub>, NO and CO was relatively  
30 mild as compared to a presumably dominant input from intensive traffic on the nearby  
31 roads (AURN, 2010). On the other hand, contributions to other pollutants might have  
32 been more prominent (e.g., particulate matter due to dust) and hence could be  
33 picked up by our methods; indeed, PM<sub>10</sub> concentration data (not reported here) re-  
34 veal substantially elevated levels from April to December 2005.

35       5.2.6. *Fuel crisis 2000*. The UK “fuel crisis”, lasting one week in September 2000,  
36 was caused by the increase of the fuel prices and subsequent public protests includ-  
37 ing blockade of oil facilities, which led to disruption of fuel deliveries and hence to  
38 significant reductions in traffic flows (Hathaway, 2001). According to AQAP (2004),  
39 this resulted in a dramatic improvement of the air quality in the city (e.g., by 30–40%  
40 for NO<sub>2</sub>). Our model has not picked up any noticeable change near this point in time,  
41 which is not surprising in view of a very short-lived impact of this episode.

1       5.2.7. *European emission standards.* These standards determined acceptable  
2 exhaust emission limits for new vehicles sold within the European Union (EUR,  
3 2003); e.g., for passenger cars the staged inception dates were: July 1992 (Euro 1),  
4 January 1996 (Euro 2), January 2000 (Euro 3), January 2005 (Euro 4) and Septem-  
5 ber 2009 (Euro 5). Some of these dates may be related to the identified change-  
6 points (Table 2). In particular, 29 October 1999 (NO<sub>2</sub>) and 11 November 1999 (NO)  
7 are close to Euro 3 (covered by the 50% credible intervals), while the Christmas 2004  
8 change-points (Section 5.2.4) may be linked with Euro 4. Furthermore, 18 May 1995  
9 (NO) and 5 April 1995 (CO) may be attributed (even though with somewhat lower  
10 credibility) to Euro 2 which significantly tightened up the CO standard for diesel cars  
11 and reduced the NO<sub>x</sub> permitted levels by 25–45% (EUR, 2003).

## 12   **6. Conclusions**

13       In this paper, we have used a non-homogeneous Poisson model to fit the ob-  
14 served series of threshold exceedance occurrences in time, extracted from the daily  
15 pollution concentration data obtained in Leeds (UK) in 1993–2009. The time depend-  
16 ence of the unknown Poisson rate was assumed to be a step function, with an arbi-  
17 trary number of possible change-points and no restrictions on its values (in particular,  
18 with no *a priori* requirement of monotonicity). The principal aim of the modelling was  
19 to identify statistically significant change-points in the rate function, signalling notable  
20 changes in the pollution dynamics patterns. Statistical estimation of the unknown rate  
21 function was carried out using Bayesian posterior samples obtained by adapting an  
22 MCMC algorithm proposed by Green (1995). Our results have demonstrated the  
23 computational and statistical efficiency of the MCMC method.

24       Once the statistical identification is carried out and validated (e.g., by simulations),  
25 it is important to seek a meaningful interpretation of the results, bearing in mind up-  
26 dates in the environmental policies by local and central authorities as one possible  
27 source of the emerging change-points. Certain care is needed in practical applica-  
28 tions of the Poisson model, since “false” change-points may be induced by meteoro-  
29 logical conditions, seasonality, missing values and cluster correlations in threshold  
30 exceedances. We have addressed these issues by the data preprocessing (including  
31 deseasonalisation and declusterisation), supported by a graphical and statistical as-  
32 sessment of the results.

33       The application of non-homogeneous Poisson processes with multiple change-  
34 points may provide a suitable modelling framework in the air quality management,  
35 which can be used to analyse pattern changes in the violation of air quality standards  
36 by various pollutants. In particular, these methods may be instrumental in obtaining  
37 important feedback about environmental actions and assessing their impact on the  
38 dynamical patterns of potentially hazardous pollutants.

## 39   **Acknowledgements**

40       J. Gyarmati-Szabó was supported by an EPSRC Doctoral Training Grant, by the  
41 Strategic Fund of the Institute for Transport Studies and by a Postgraduate Research

1 Scholarship of the School of Mathematics (University of Leeds). L.V. Bogachev was  
2 partially supported by a Leverhulme Research Fellowship.

### 3 **References**

- 4 Achcar, J.A., Rodrigues, E.R., Tzintzun, G., 2009. Using non-homogeneous Poisson  
5 models with multiple change-points to estimate the number of ozone exceedances  
6 in Mexico City. *Environmetrics*, published online, doi: 10.1002/env.1029 (ac-  
7 cessed 30 March 2010).
- 8 Akman, V.E., Raftery, A.E., 1986. Asymptotic inference for a change-point Poisson  
9 process. *Annals of Statistics* 14, 1583–1590.
- 10 AQA, 2010. The UK Air Quality Archive,  
11 [http://www.airquality.co.uk/data\\_and\\_statistics.php](http://www.airquality.co.uk/data_and_statistics.php) (accessed 6 April 2010).
- 12 AQAP, 2004. Air Quality Action Plan, Leeds City Council,  
13 <http://www.leeds.gov.uk/files/Internet2007/2010/15/doc%2012.pdf> (accessed 16  
14 October 2010).
- 15 AQMA, 1997. Air Quality Management Areas. Defra <http://aqma.defra.gov.uk/>. Infor-  
16 mation about Leeds:  
17 [http://aqma.defra.gov.uk/maps.php?map\\_name=fulluk&la\\_id=143](http://aqma.defra.gov.uk/maps.php?map_name=fulluk&la_id=143) (accessed 26  
18 July 2010).
- 19 AQS, 2007. The Air Quality Strategy for England, Scotland, Wales and Northern Ire-  
20 land, Vols. 1, 2. Cm 7169 NIA 61/06-07, Defra,  
21 <http://www.defra.gov.uk/environment/quality/air/airquality/strategy/> (accessed 5  
22 May 2008).
- 23 AURN, 2010. UK Automatic Urban and Rural Network, Defra, Leeds Centre (station  
24 #54), <http://aurn.defra.gov.uk/stations/viewStation.php?id=54> (accessed 21 Octo-  
25 ber 2010).
- 26 Ayres, J.G., 2002. Chronic effects of air pollution. *Occupational and Environmental*  
27 *Medicine* 59, 147–148.
- 28 Brimblecombe, P., 1998. History of urban air pollution. In: Fenger, J. et al. (Eds.), *Ur-  
29 ban Air Pollution — European Aspects*, Kluwer, Dordrecht, pp. 7–20.
- 30 CCME, 2002. The Ongoing Challenge of Managing Carbon Monoxide Pollution in  
31 Fairbanks, Alaska. National Academy Press, Washington;  
32 <http://www.nap.edu/catalog/10378.html> (accessed 6 October 2010).
- 33 Clapp, L.J., Jenkin, M.E., 2001. Analysis of the relationship between ambient levels  
34 of O<sub>3</sub>, NO<sub>2</sub> and NO as a function of NO<sub>x</sub> in the UK. *Atmospheric Environment* 35,  
35 6391–6405.
- 36 Cox, D.R., 1981. Statistical analysis of time series: Some recent developments.  
37 *Scandinavian Journal of Statistics* 8, 93–108.

- 1 Cox, D., Lewis, P.A.W., 1966. The Statistical Analysis of Series of Events. Methuen,  
2 London.
- 3 ELLR, 2010. East Leeds Link Road. Leeds City Council, [www.leeds.gov.uk/ellr](http://www.leeds.gov.uk/ellr) (ac-  
4 cessed 1 August 2010).
- 5 EUR, 2003. Trends in Vehicle and Fuel Technologie. Report EUR 20746 EN, EC  
6 Joint Research Centre, <http://ftp.jrc.es/EURdoc/eur20746en.pdf> (accessed 16 Oc-  
7 tober 2010).
- 8 Faulkner, M., Russell, P., 2010. Review of Local Air Quality Management. Report,  
9 IHPC,  
10 [http://www.defra.gov.uk/environment/quality/air/airquality/local/documents/laqm-  
11 report.pdf](http://www.defra.gov.uk/environment/quality/air/airquality/local/documents/laqm-report.pdf) (accessed 16 October 2010).
- 12 Gibbons, J.D., Chakraborti, S., 2003. Nonparametric Statistical Inference, 4th edition.  
13 Marcel Dekker, New York.
- 14 Gilks, W.R., Richardson, S., Spiegelhalter, D.J., 1996. Introducing Markov chain  
15 Monte Carlo. In: Gilks, W.R. et al. (Eds.), Markov Chain Monte Carlo in Practice,  
16 Chapman & Hall/CRC, London, pp. 1–19.
- 17 Glaister, S., 2002. UK transport policy 1997–2001. Oxford Review of Economic Pol-  
18 icy 18, 154–186.
- 19 Green, P.J., 1995. Reversible jump Markov chain Monte Carlo computation and  
20 Bayesian model determination. Biometrika 82, 711–732.
- 21 Hathaway, P., 2001. The effect of the fuel ‘protest’ on road traffic. In: Transport  
22 Trends: 2001 Edition. TSO, London;  
23 [http://web.archive.org/web/20040224182437/http://www.dft.gov.uk/stellent/groups/  
24 dft\\_transstats/documents/page/dft\\_transstats\\_505952.pdf](http://web.archive.org/web/20040224182437/http://www.dft.gov.uk/stellent/groups/dft_transstats/documents/page/dft_transstats_505952.pdf) (accessed 1 August  
25 2010).
- 26 HOVL, 2004. Case study: High occupancy vehicle lanes, A647 Stanningley Road,  
27 Leeds. In: Bus Priority: The Way Ahead, 2nd edition,  
28 [http://www.dft.gov.uk/pgr/regional/buses/bpf/busprioritythewayahead12/  
29](http://www.dft.gov.uk/pgr/regional/buses/bpf/busprioritythewayahead12/) (accessed 8 October 2010).
- 30 LCCA, 2010. 9th Leeds City Centre Audit 2010. Leeds City Council,  
31 [http://www.leedsinitiative.org/WorkArea/DownloadAsset.aspx?id=18786  
32](http://www.leedsinitiative.org/WorkArea/DownloadAsset.aspx?id=18786) (accessed 10 October 2010).
- 33 Leadbetter, M.R., Lindgren, G., Rootzén, H., 1983. Extremes and Related Properties  
34 of Random Sequences and Processes. Springer, New York.
- 35 Lewis, P.A.W., Shedler, G.S., 1976. Statistical analysis of non-stationary series of  
36 events in a data base system. IBM Journal of Research and Development 20,  
37 465–482.
- 38 MET-O, 2010. Met Office UK. Climate: North East England,  
39 <http://www.metoffice.gov.uk/climate/uk/ne/> (accessed 8 October 2010).

1 NAEI, 2010. UK Emissions of Air Pollutants 1970 to 2008. Report, National Atmos-  
2 pheric Emission Inventory (NAEI),  
3 [http://www.airquality.co.uk/reports/cat07/1009030925\\_2008\\_Report\\_final270805.](http://www.airquality.co.uk/reports/cat07/1009030925_2008_Report_final270805.pdf)  
4 [pdf](http://www.airquality.co.uk/reports/cat07/1009030925_2008_Report_final270805.pdf) (accessed 16 October 2010).

5 Raftery, A.E., 1989. Comment: Are ozone exceedance rates decreasing?. *Statistical*  
6 *Science* 4, 378–381.

7 Raftery, A.E., 1996. Hypothesis testing and model selection. In: Gilks, W.R. et al.  
8 (Eds.), *Markov Chain Monte Carlo in Practice*, Chapman & Hall/CRC, London, pp.  
9 163–187.

10 Shively, T.S., 1991. An analysis of the trend in ground-level ozone using non-  
11 homogeneous Poisson processes. *Atmospheric Environment* 25, 387–395.

12 Smith, R.L., 1989. Extreme value analysis of environmental time series: An applica-  
13 tion to trend detection in ground-level ozone. *Statistical Science* 4, 367–377.

14 Smith, R.L., 2004. Statistics of extremes, with applications in environment, insurance,  
15 and finance. In: Finkenstädt, B., Rootzén, H. (Eds.), *Extreme Values in Finance,*  
16 *Telecommunications, and the Environment*, Chapman & Hall/CRC, Boca Raton,  
17 pp. 1–78.

18 Smith, R.L., Shively, T.S., 1995. Point process approach to modeling trends in tropo-  
19 spheric ozone based on exceedances of a high threshold. *Atmospheric Environ-*  
20 *ment* 29, 3489–3499.

21 Thompson, M.L., Reynolds, J., Cox, L.H., Guttorp, P., Sampson, P.D., 2001. A re-  
22 view of statistical methods for the meteorological adjustment of tropospheric  
23 ozone. *Atmospheric Environment* 35, 617–630.

24 TMA, 2004. *Traffic Management Act 2004 (c. 18)*. TSO, London;  
25 <http://www.statutelaw.gov.uk/content.aspx?activeTextDocId=1606563> (accessed  
26 26 July 2010).

27 Venables, W., 2002. *Modern Applied Statistics with S*, 4th edition. Springer, New  
28 York.

29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

1  
2  
3

## Figure captions

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37

**Fig. 1.** Cumulative count plots of threshold exceedances for the observed raw data (*black*) and for the data with the missing values estimated via the moving-average with a two-sided window of size  $\pm 65$  days (*grey*). Black ticks on the time axis indicate the occurrence of missing values; red points on the grey curve show the imputed exceedance values.

**Fig. 2.** *Top:* observed daily maximum concentrations of  $\text{NO}_2$ ,  $\text{NO}$  and  $\text{CO}$  (*black*) and the estimated yearly periodic trend for the log-transformed data (*red*). *Bottom:* de-seasonalised data, with the yearly periodic trend removed.

**Fig. 3.** Estimated MCMC posterior distribution of the number  $k$  of change-points in the step rate function  $\lambda(t)$ , obtained for the observed threshold exceedance data (*red*) and for data simulated via non-homogeneous Poisson process with the estimated step function  $\hat{\lambda}(t)$  as the rate (*blue*).

**Fig. 4.** Estimated MCMC posterior density functions of the marginal distributions for the locations  $s_j$  of the identified change-points (cf. Table 2), obtained for the observed threshold exceedance data (*red*) and for data simulated via non-homogeneous Poisson process with the estimated step function  $\hat{\lambda}(t)$  as the rate (*blue*). The estimates were computed using a Gaussian kernel with bandwidth 95 days.

**Fig. 5.** Estimated MCMC posterior density functions of the univariate marginal distributions for the heights  $h_j$  of the estimated change-points (cf. Table 3), obtained for the observed threshold exceedance data (*red*) and for data simulated via the estimated step rate function  $\hat{\lambda}(t)$  (*blue*). The estimates were computed using a Gaussian kernel with bandwidth  $0.003 \text{ days}^{-1}$ .

**Fig. 6.** MCMC results for each pollutant: estimated rate functions (*left axes*) and the corresponding cumulative plots for data sets simulated via non-homogeneous Poisson process (*right axes*). Colour-coded plots correspond to the posterior mean rate  $\bar{\lambda}(t)$  (*red*) and a step rate function  $\hat{\lambda}(t)$  estimated from the posterior distribution (*blue*). *Black* dotted graphs show the cumulative plots of observed threshold exceedances. Dashed vertical lines indicate the locations of estimated change-points.

1 **Fig. 7.** Validation of the MCMC performance in the step rate Poisson model estima-  
2 tion. The rate graphs (*left axes*) show the step rate function  $\hat{\lambda}(t)$  (*blue*) estimated for  
3 the observed threshold exceedances from the MCMC posterior distribution, and the  
4 posterior mean rate  $\bar{\lambda}(t)$  (*red*) obtained for a data set simulated via a non-  
5 homogeneous Poisson process with rate  $\hat{\lambda}(t)$ . Cumulative plots (*right axes*) corre-  
6 spond to data sets simulated via non-homogeneous Poisson processes with rate  
7  $\hat{\lambda}(t)$  (*blue*) or  $\bar{\lambda}(t)$  (*red*).

8

9 **Fig. 8.** Bayesian validation of the MCMC results using replicated counting processes.  
10 The graphs show cumulative plots for simulated data, non-conditional (*left*) and con-  
11 ditional (*right*) on the total number of observed threshold exceedances. *Black*: cumula-  
12 tive plots for the original exceedance data. *Grey*: 1000 replicated trajectories simu-  
13 lated each with a rate function independently sampled from the MCMC posterior dis-  
14 tribution. *Red*: pointwise average of these trajectories. *Blue*: pointwise average of  
15 another 1000 simulated trajectories of a Poisson counting process with the step rate  
16 function  $\hat{\lambda}(t)$  estimated from the posterior distribution.

17

18 **Fig. 9.** Combined log-linear/periodic trend of the form (17) (*red*) fitted to the the log-  
19 transformed raw data (*black*).

20

21

22

23

24

25

26

27

28

29

30

1

2

## Table captions

3

### 4 **Table 1**

5 Comparison of the MCMC estimation results for NO<sub>2</sub> before and after imputation of  
6 missing values. Here  $n$  = number of exceedances in  $[0, T]$ ;  $k$  = number of change-  
7 points;  $s_1$  = location of the (first) change-point;  $h_0, h_1$  = step heights of the rate  
8 function on the intervals  $(0, s_1)$ ,  $(s_1, T)$ , respectively.

9

### 10 **Table 2**

11 Posterior modes of the estimated change-points  $s_j$  and the corresponding lower and  
12 upper quartiles, expressed as day numbers (ranging from 1 to 6206) as well as  
13 calendar dates (dd-mm-yy).

14

### 15 **Table 3**

16 Posterior modes of the estimated heights  $h_j$  of the unknown step rate function  $\lambda(t)$ ,  
17 with the corresponding lower and upper quartiles.

18