

NISS

Statistical Analyses of Freeway Traffic Flows

Claudi Tebaldi, Mike West, and Alan F. Karr

Technical Report Number 111
November, 2000

National Institute of Statistical Sciences
19 T. W. Alexander Drive
PO Box 14006
Research Triangle Park, NC 27709-4006
www.niss.org

STATISTICAL ANALYSES OF FREEWAY TRAFFIC FLOWS

CLAUDIA TEBALDI¹, MIKE WEST² & ALAN F. KARR³

NOVEMBER 9, 2000

This paper concerns the exploration of statistical models for the analysis of observational freeway flow data, and the development of empirical models to capture and predict short-term changes in traffic flow characteristics on sequences of links in a partially detectorised freeway network. A first set of analyses explores regression models for minute-by-minute traffic flows, taking into account time of day, day of the week, and recent upstream detector-based flows. Day- and link-specific random effects are used in a hierarchical statistical modelling framework. A second set of analyses captures day-specific idiosyncrasies in traffic patterns by including parameters that may vary throughout the day. Model fit and short-term predictions of flows are thus improved significantly. A third set of analyses includes recent downstream flows as additional predictors. These further improvements, though marginal in most cases, can be quite radically useful in cases of very marked breakdown of freeway flows on some links. These three modelling stages are described and developed in analyses of observational flow data from a set of links on Interstate Highway 5 (I-5) near Seattle.

Keywords: *Bayesian inference; Hierarchical regression models; Short-term prediction; Dynamic hierarchical regression models; Traffic flow.*

¹Correspondence to Claudia Tebaldi, Athene Software inc., 2060 Broadway, Suite 300, Boulder, CO, 80302. claudiat@athenesoft.com

²Institute of Statistics and Decision Sciences, Duke University, Durham, NC.

³National Institute of Statistical Sciences and Department of Statistics, University of North Carolina, Chapel Hill, NC.

1. INTRODUCTION

This article describes and illustrates statistical models of freeway traffic flows and their uses in evaluating and predicting short-term changes in flows on sequences of freeway segments. This investigation arises from a study undertaken by the National Institute of Statistical Sciences (NISS) to explore the utility of traffic detector data in presaging changes in traffic flow patterns. As part of this project, detector records of vehicle counts were collected along a section of Interstate Highway 5 (I-5) near Seattle. Issues motivating the data collection exercise include questions about detector placement design in connection with the utility of detector-based data in short-term forecasting of changes, and particularly breakdowns, in flow rates over identified stretches of the highway. Short-term forecasting of traffic flows is an area of growing interest. Omnipresent road congestion from one side and development in electronics from the other have made this a global issue and a fertile field for prediction technology. Recently the International Journal of Forecasting recognized it by devoting a section to contributions specific to it. (IJF, 1997; vol.13, n.1). But prior to addressing such issues, the basic structure of traffic flow patterns, as measured by detector count data, needs illuminating; the study reported here represents aspects of our statistical investigation of this basic structure.

Our primary objectives are to investigate models capable of tracking the development of traffic flows throughout the day on highway links, identified as sections of highway between consecutive detectors. For each link, traffic flows are modelled as functions of sets of explanatory variables that may include

- detector-based measures of traffic flows at preceding time points at upstream and downstream locations;
- topography of the locations involved (distances between detectors, presence of on- and off-ramps between detectors);
- time of the day, and
- day of the week.

Our study is based on a southbound section of I-5, just north of Seattle, fitted with a series of single-loop detectors. Further details on the area and data collection exercises appear in Graves et al (1997). From the data collected by these authors, we identified and extracted data from a 6-mile stretch of highway with 15 detector locations. Each location has one single-loop detector per lane, with either three or four lanes per location. The data explored represent a total of fifty days of recordings (from the end of May to the middle of August 1996), with recorded flows for 4 hours (7-11am) on each weekday morning in this period. The detector flows record vehicles traveling in all lanes, i.e., are aggregated across lanes at each location. On and off ramps are not detectorised, so contributing the major components of uncertainty to the modelling problem.

All detectors are single loop, presence-type detectors, recording 20 second flows which were aggregated to the 1 minute level for our analysis. Hence the raw response data are reported numbers of vehicles per minute at each of the detector locations. In addition to natural variability in flows based on traffic patterns and local conditions, unexplained variability in predicted flows will be partially contributed by detector errors and inaccuracies. At the 1 minute level, we are in a position to assess relationships between flows and predict one or two minutes ahead, section by section of the highway.

This report discusses models fitted to 9 days of data, beginning on 6/17/96 and we present summary inferences for a selection of 6 detector locations. Figure 1 is a schematic of the section of I5 in question, with the detector locations labelled 1-6 and the lengths of the highway links between detectors marked in miles.

We begin with basic regression models suggested by the physical highway layout and detector placement, and supported by preliminary exploratory data analyses. These models take into account time of day, day of week and recent upstream detector-based flows. Relationships between days are captured through the use of day- and link-specific random effects in an hierarchical modelling framework. Refinements of these models to include parameters that may vary moderately throughout the course of the day are then investigated, and found to improve model fit and short-term predictions of flows significantly. Further model extensions to include recent downstream flows as additional covariates provide further improvements, marginal in most cases but very radically useful on one or two specific highway

links. Each of these modelling stages is described below, with illustrations of model fit and short-term predictive performance across a selected set of links on I5. Appendices include supporting technical material on models and analysis, and some summary comments and discussion of potential further developments appear in the final section.

2. HIERARCHICAL REGRESSION MODELS

2.1. DETECTOR SPECIFIC REGRESSION MODELS

We begin with initial regression models for individual link flows. The highway section has six links identified by their terminal detector nodes labelled $i = 1, \dots, I = 6$. Our data represent an initial nine days, indexed by $j = 1, 2, \dots, J = 9$ and, at the 1 minute resolution level, we observe link flows at minutes $t = 1, \dots, T = 228$ during the morning period of each day. Write y_{ijt} for the vehicle flow observed on link i during day j and at minute t ; the y_{ijt} are in units of vehicles per minute based on the raw detector counts aggregated to the 1 minute level.

The conceptual basis for the initial regression framework has the following components.

- Between each pair of detectors but numbers 2 and 3 there is at least one entry or exit ramp. Much of the variability in link flows as measured at terminal detectors will be driven by the on/off traffic on ramps. As the ramps are not detectorised, we have no basis for modelling or forecasting the ramp activity and so simply describe it through parameters to be estimated based on past data. Flows at link 1 will be wholly represented this way.
- Traffic flows at any detector at time t are driven partly by flows on upstream links at times $t - 1, t - 2$, and so forth. The relationship is, of course, speed and congestion dependent, and moderated by the possibly dominating in/out flows on intervening ramps. At a first step, our regressions model flows directly in terms of proportions of vehicles passing through in intervals of 1 minute, 2 minutes and so forth.
- Days are treated exchangeably, initially. That is, parameters defining the daily pattern of flows and the regression responses to upstream traffic are assumedly drawn from

a population of “daily flow” parameters, one for each link in the highway network. Thus, while allowing for close or common parameters, the model will allow for daily effects due to prevailing weather patterns, road conditions, driver behaviour and so forth. This is a standard random effects or hierarchical modelling concept; the days differ as their parameters are different, but the differences are not predicted, being treated as purely random and to be estimated. Predictions of future days with no observed data is therefore immediately possible as day-specific parameters are drawn from these underlying population models, and the characteristics of the population distributions are estimated based on the observed days data. Also, differences between days may be investigated through inferences on day-specific parameters.

- Beyond a natural daily pattern and the regression responses to upstream traffic, residual variation in link flows is assumed to be purely random and independent across links, in a standard normal linear regression model. Note that we do not transform the data, maintaining the original vehicle flow scale so as to be able to interpret regression parameters on upstream flows as proportionate contributions to predicted flows. Model validation and residual analyses confirm model adequacy.

As an example, consider link 5, terminated by detector 5. This link has one exit and three entry ramps, so that traffic recorded at detector 5 will likely be heavily influenced by incoming and exiting traffic. Given the distance of almost two miles from detector 4 to 5, the relationship between flows at detector 5 and those in the past at detector 4 will be of limited predictive value unless the ramp activity is minor.

The mathematical form of the regression models is as follows. For link i on day j at minute t we have flows

$$y_{ijt} = s_{ij}(t) + b'_{ijt}\beta_{ij} + \epsilon_{ijt} \quad (1)$$

where

- $s_{ij}(t)$ represents the smooth trend over day j driven by on/off ramp flows,
- b_{ijt} is a vector of known flows from preceding upstream links $i - 1, i - 2, \dots$ at recent times $t - 1, t - 2, \dots$, whose size depends on the number of upstream links and maximum time lag chosen,

- β_{ij} is the regression parameter defining the expected response to these predictor flows, and
- ϵ_{ijt} represents residual, unexplained variation in flows.

The daily trend patterns are modelled via cubic splines. Following exploratory data analyses we adopt splines with terminal knots at $t = 0$ and $t = T = 228$, and two interior knots at minutes $t = 60$ and $t = 150$. Then $s_{ij}(t)$ is a day- and link- specific constant plus a cubic spline interpolating these knots. This requires a total of six parameters denoted by the 6-vector α_{ij} , and can be represented as the linear model

$$s_{ij}(t) = a'_{it}\alpha_{ij}$$

where a_{it} is the 6-vector taken as the relevant row of the basis matrix for the spline function over time t . We therefore have

$$y_{ijt} = x'_{ijt}\theta_{ij} + \epsilon_{ijt}$$

where $x'_{ijt} = (a'_{it}, b'_{ijt})$ and $\theta'_{ij} = (\alpha'_{ij}, \beta'_{ij})$. Combining across the day $t = 1, \dots, T$, write Y_{ij} for the column T -vector of flows y_{ijt} , X_{ij} for corresponding matrix whose rows are the minute-specific, known row-vectors x'_{ijt} and e_{ij} for the T -vector of elements ϵ_{ijt} . Then the link i , day j model can be expressed as the linear model

$$Y_{ij} = X_{ij}\theta_{ij} + e_{ij}. \quad (2)$$

Additionally, we assume independence and normality of the ϵ_{ijt} , namely $N(\epsilon_{ijt}|0, \sigma_i^2)$ for each i, j, t where $N(e|m, v)$ denotes a normal distribution, of mean m and variance v , for the random quantity e . This implies an essentially standard normal linear model for the flows.

Link by link, the regression components b_{ijt} are detailed as follows. The selection of link-specific “upstream” predictors is based on our initial exploratory data analysis and modelling (not described here) and physical considerations.

1. Detector 1 has no upstream links so the regression term is absent at $i = 1$. Thus X_{1j} is simply the $T \times 6$ spline basis matrix and $\theta_{1j} = \beta_{1j}$ are the 6-vectors of day-specific trend parameters.

2. Detector 2 flows are partly predicted by flows at detector 1 at times $t - 1$ and $t - 2$. There is a distance of 1.28 miles between the two detectors, and so, supposing normal conditions of traffic, a proportion of cars recorded 1 and 2 minutes before at detector 1 can be expected to be recorded at detector 2. So β_{2j} is a 2-vector, X_{2j} is $T \times 8$, and θ_{2j} is an 8-vector.
3. Detector 3 flows use the same regressor information as detector 2, as the distance between the two is too short to allow for any correlation between the traffic at detector 2 one minute before and the traffic at detector 3 at present. So β_{3j} is a 2-vector, X_{3j} is $T \times 8$, and θ_{3j} is an 8-vector.
4. Detector 4 flows are regressed on past flows from all three of the upstream detectors: the flows at detectors 2 and 3 at time $t - 1$, and the flows at detector 1 at times $t - 2$ and $t - 3$. So β_{4j} is a 4-vector, X_{4j} is $T \times 10$, and θ_{4j} is a 10-vector.
5. Detector 5 flows are predicted by flows at detector 4 just 2 and 3 minutes before, the distance between the two detectors being almost 2 miles. So β_{5j} is a 2-vector, X_{5j} is $T \times 8$, and θ_{5j} is an 8-vector.
6. Detector 6 is predicted only by flows at detector 5 in the previous minute. So β_{6j} is a 1-vector, X_{6j} is $T \times 7$, and θ_{6j} is a 7-vector.

The hierarchical modelling components arise as descriptors of the variation in both spline-trend and regression parameters across days. At link i , we assume that the θ_{ij} are drawn from a multivariate normal distribution of parameters, representing day-by-day variability at this link. Thus

$$\theta_{ij} \sim N(\theta_{ij} | \mu_i, \Sigma_i) \tag{3}$$

independently across all days j and links i . Here μ_i represents the expected parameters (spline and regression) at link i about which individual days are distributed. The extent and nature of the day-specific variation is moderated by link-specific variance matrices Σ_i .

The model specification is completed by defining prior distributions for the population parameters, or hyperparameters, μ_i , Σ_i and σ_i^2 for each link i . We use standard, conditionally

conjugate priors here, namely

$$\mu_i \sim N(\mu_i | m_i, M_i),$$

independently,

$$\Sigma_i^{-1} \sim W(\Sigma_i^{-1} | (\rho_i R_i)^{-1}, \rho_i)$$

and

$$\sigma_i^{-2} \sim G(\sigma_i^{-2} | \nu, \nu \tau^2 / 2)$$

where W and G denote Wishart and Gamma distributions, respectively. The defining quantities ν, τ and the m_i, M_i, ρ_i and R_i for each i , are fixed and specified. The analysis reported below assumed essentially uninformative hyperpriors, with very large elements on the diagonals of the prior variance matrices M_i for the link-specific parameters μ_i . Useful background to normal, linear hierarchical models and prior forms can be found in, for example, Lindley and Smith (1972), Smith (1973), and more recently and specifically in Gelfand et al (1990), and in several of the contributed chapters in Gilks et al (1996).

Analysis of the nine days of observed data involves computing posterior distributions for all model parameters and the population hyperparameters. Analysis is via standard methods of Bayesian simulation using Markov chain Monte Carlo techniques, and basically follows developments in the final two references above. Basic mathematical details appear in the appendix to this paper. With samples from posterior distributions available, we are in a position to compute posterior estimates and associated probability intervals, predictions, and so forth. Some such summaries are discussed in the next section.

2.2. SOME SUMMARY INFERENCES AND MODEL FIT ASSESSMENT

A selection of posterior summaries is reviewed here, all based on standard graphical and numerical representations of the outputs of the simulation based analysis of the full posterior distribution for all model parameters and hyperparameters.

Begin with Figure 2. This summarises marginal posteriors for the elements of μ_2 , the link 2 hyperparameters. As noted above the first six elements represent the spline parameters of the population of days for this link, and the final two are the population regression coefficients on past flows on link 1. In each frame we display density estimators of posterior

samples for these hyperparameters. Also, the “flat” density curve in each frame represents the diffuse marginal prior for each variable; this is displayed to illustrate just how relatively diffuse the priors are in connection with the claim that the priors are uninformative. These graphs are representative of marginal posteriors for hyperparameters across all links.

The sequence of Figures 3 to 8 display summary posterior inferences for link and day specific parameters θ_{ij} . The spline function parameters are summarised in terms of posterior means of the daily patterns $s_{ij}(t)$ plotted, for each link i and day j , as functions over time t during the day. Thus, for $j = 1, \dots, 9$, Figure 3 plots the posterior mean $E[s_{1j}(t)|Y] = a'_{1t}E[\alpha_{1j}|Y]$ over $t = 0, \dots, 228$. In addition, the spline function of the underlying population trends, link by link, are graphed; thus Figure 3 includes the graph over t of $a'_{1t}E[m_1|Y]$. Notice the similarities of the trends across days on this link 1. The variability in daily patterns exhibited indicates meaningful random effects day by day, consistent with non-negligible values of the scale factors in Σ_1 . Note one peculiar Wednesday that departs substantially from the underlying population trend and has quite different characteristics than the other eight days, the basic flows being much lower than typical.

For the remaining links, 2-6 inclusive, similar graphs of day-specific and underlying daily trend patterns appear in Figures 4 to 8. Link 1 is described entirely in terms of the spline patterns, whereas links 2-6 have additional parameters related to the regressions on past upstream flows. Hence these figures include additional graphs summarising posterior inferences for the day-specific regression parameters, and the underlying regression hyperparameters, for each of these links. For example, consider link 2 summarised in Figure 4. This link has just two regression parameters for past upstream flows, and the upper frames in the figure provide boxplots of the samples from the corresponding posteriors. Those labelled by names of days represent the day-specific regression parameters, and those labelled “mu” represent inferences on the underlying hyperparameters. Here again we note similarities across days, though with some evident variability day-to-day consistent with the random effects, hierarchical structure. Note also the significance of the regression parameters indicated by the general lack of coverage of zero values by the posterior intervals. The remaining figures provide similar displays for the other links. We note that the one peculiar Wednesday stands out across all six links in terms of the departure of the daily spline trend

pattern from the norm, a feature that is currently unexplained but which is naturally and adequately assumed by the hierarchical model without overly distorting inferences for other days.

Though the model assumes that the link specific patterns θ_{ij} are exchangeable across days, i.e., that their prior distribution depends only on link index i , there is some suggestion of weekday similarities from the figures of daily trends. Thus refined models might incorporate partially exchangeable priors in which, for example, the parameters for any Tuesday are a random sample from a population of Tuesdays, rather than the broader population of all days. We do not do this here formally, but do use the idea in developing out-of-sample predictive distributions for model assessment in the following subsection.

Initial examination of aspects of model fit and adequacy are explored through fitted values and estimated residuals. The two Figures 9 and 10 provide some graphical displays of data and fitted values for two days, the first Tuesday and the second Thursday of the period under study. The reader can consider them representative of the general results. In each case the flow data on all six links are plotted, and the fitted values $\hat{y}_{ijt} = x'_{ijt}E[\theta_{ij}|Y]$ are overlaid. On link 1 the fitted values simply represent the trend, but across the remaining links also include the estimated regression effects from past upstream flows. From these graphs it is clear that the regression components are very effective in tracking the short-term variability in link-to-link flows on shorter links. Link 5 is the least well modelled as it is relatively long, so that recent downstream flows are less indicative of current activity at the fifth detector, and more heavily influenced by the in/out flows on the several intervening ramps. The fitted flows on the remaining links are remarkably accurate, capable of tracking both minor and more significant fluctuations. Note in particular the marked response to the major slowdown on links 5 and 6, and the ability of the link 6 model to rapidly respond and adjust.

The model explains a significant degree of observed variability across all detectors and days. Daily sample variances for each detector are quite variable across the full 52 days of data. Average sample variances by detector are roughly 228, 211, 214, 228, 270 and 175, but the values range quite widely, from minima of roughly 69, 94, 95, 117, 108 and 71, to maxima of roughly 1475, 458, 456, 483, 787 and 709. The average sample variances for the

specific nine days in our study are roughly 183, 217, 223, 240, 271 and 176. By comparison, the posterior means of the model residual variances σ_i^2 for the six links are roughly 85, 100, 110, 111, 122 and 71, respectively, so indicating a very meaningful reduction in observed variability, though there is substantial unexplained variation in the data. Further residual analyses confirm, across all detectors, no significant evidence of departures from the assumed normality, and no evidence of residual correlation structure.

2.3. PREDICTION OF NEW DAYS

More incisive and relevant assessments of model adequacy are based on out-of-sample predictive performance, as opposed to the above in-sample assessments using fitted values. To explore this we examined conditional flow forecasts across a representative selection of additional days whose data are available. Figures 11 and 12 provide plots of the data for two randomly chosen future days, and overlaid are point forecasts derived as forecast means from approximate conditional predictive distributions. Specifically, these point forecasts are simply the conditional means in (2) with each θ_{ij} replaced by posterior means from the data analysis of the original nine days. For the Tuesday 2nd July forecasts, the parameters are estimated by posterior means from the first Tuesday of the original nine days, and those for Thursday 8th August by posterior means from the second Thursday of the original nine days. This represents a minor departure from the model in that we are now implicitly identifying and, in an ad-hoc manner, incorporating a possible day-of-week effect, as was earlier alluded to. Further, this does not fully exploit the predictive ability of the models as more formal exploration of predictive distributions might. Nevertheless, performance in terms of the match of point forecasts to outcomes, is excellent, and really comparable to the in-sample fitted values of the earlier figures.

Real-time use of these models will involve sequential analysis, proceeding through the day and successively updating posterior distributions across links as flow data is received. This is the context in which models should be developed and assessed. Before exploring this, we embark on relatively minor model generalisation that moves us naturally into a sequential processing format while adding potential for improvements in short-term predictive ability.

3. DYNAMIC HIERARCHICAL REGRESSION MODELS

Improvements in day-specific, short-term forecasts can be expected to arise from refined models that adapt to the observed within-day variability in flows using dynamic modelling methods from time series (West and Harrison 1997). As it stands the basic regression model of equations (2) and (3) provides day-specific effects that are constant over the course of the day. This constancy is certainly most appropriate for the basic spline parameters representing the in/out flows from ramps, but perhaps a little rigid in connection with the regression parameters β_{ij} that represent the transfer effects of historical upstream flows. Changing road and weather conditions, mix of vehicle type, idiosyncratic driver behaviour, and so forth all impact randomly on within-day flows in the very short-term. Such unpredictable influences on link flow rates are captured in the regression model wholly by the residual terms ϵ_{ijt} , whereas their effects are, in part, made physically evident through the transfer responses to downstream flows and hence in a modelling framework through changes in regression parameters. This is a traditional concept that underlies the entire field of dynamic modelling and Bayesian forecasting (West and Harrison 1997). A typical immediate benefit of recognising and modelling time-varying regression parameters is increased short-term forecasting accuracy and reduced forecast uncertainty, as the examples in West and Harrison (1997, chapters 2 and 3, particularly) vividly demonstrate. Also, the use of Bayesian dynamic modelling for short term forecasting has been proven effective by Whittaker and al. (1997), though under a different perspective, which models the physical process underlying a traffic network from a theoretical point of view. Here we remain in a strictly empirical domain, letting the incoming data modify the parameters' estimate.

We detail and explore the simplest dynamic model extension of the hierarchical regressions above. The only change to the basic model form in equation (1) is that the regression parameters β_{ij} are now considered as possibly varying over the day at the minute-by-minute level $t = 1, \dots, T$. Thus, for link i on day j at minute t we now model detector flows as

$$y_{ijt} = s_{ij}(t) + b'_{ijt}\beta_{ijt} + \epsilon_{ijt} \quad (4)$$

where $s_{ij}(t)$ represents the smooth trend over day j as earlier, b_{ijt} is a vector of known flows from preceding upstream links as earlier, and $\epsilon_{ijt} \sim N(\epsilon_{ijt}|0, \sigma_i^2)$ represents residual variation in flows, also as in the original formulation. The only change from the original model is that the regression parameter defining the expected response to the predictor flows

b_{ijt} is now time-dependent, denoted by β_{ijt} at minute t . The simplest exploratory model for time-variation in these parameters is a dynamic regression model that simply allows for small changes through the day, but which does not anticipate directions of change (West and Harrison 1997, chapter 3). In this normal linear model, we adopt the random walk

$$\beta_{ijt} = \beta_{ij,t-1} + \omega_{ijt} \quad (5)$$

where ω_{ijt} represents a sequence of zero-mean, normally distributed random terms, referred to as evolution errors or innovations, $N(\omega_{ijt}|0, W_{ijt})$. Between minutes $t - 1$ and t the regression parameter vector changes randomly via the addition of the innovation ω_{ijt} , and such changes can therefore be estimated based on the sequentially received and processed flow data. It is assumed that the innovations are independent through time t within days j , and also across links i . The model requires that we specify or estimate the innovation variances W_{ijt} , and this is done using standard methods of discounting (West and Harrison 1997, chapter 6). Some technical details of this are given in the appendix here, together with a summary of the model formulæ and resulting computations. The essential detail is that this approach defines innovations variances to depend on a single scalar *discount factor* $\delta \in (0, 1]$, i.e., $W_{ijt} = W_{ijt}(\delta)$. The structure is such that $W_{ijt}(1) = 0$, so that $\delta = 1$ reduces the dynamic model to the original constant regression $\beta_{ijt} = \beta_{ij}$ constant through the day. This provides opportunity to assess the original constant model and compare its predictive performance and fit to alternative dynamic models with $\delta < 1$. Further, each variance $W_{ijt}(\delta)$ is a decreasing function of δ , so that smaller discount factors lead to larger variances and hence greater potential variability in the regression parameters. As in other areas of dynamic modelling, discount factors close to but less than unity are expected to be appropriate; this allows for very small degrees of change in the regression relationship to adapt to random variations throughout the course of the day. This is confirmed by several re-analyses of the original nine days of data, as we shall see below.

3.1. PERSPECTIVE AND MODEL FITTING

The sequential dynamic model analysis proceeds in parallel across links, so consider any link i at the start of the “current” day j . Based on the past data and analysis, we have current

initial information D_{j0} which provides the basis for initial distributions (across all links) for the initial values of the current days parameters θ_{ij0} and σ_i^2 . The joint distribution has the conjugate form, namely $N(\theta_{ij0}|m_{ij0}, \sigma_i^2 C_{ij0})$ multiplied by $G(\sigma_i^{-2}|n_{ij0}/2, n_{ij0}S_{ij0}/2)$, as noted in appendix. Our analyses of new days of data reported below bases these initial distributions on the results of the static analysis of the first nine days of data, exactly as described in the out-of-sample prediction analysis under the static model. Thus the variance estimate S_{ij0} and degree of freedom n_{ij0} are taken as the posterior values from the nine days analysis on link i , the estimate m_{ij0} is the posterior value for the θ vector on the last weekday of the same name, and C_{ij0} is the posterior mean of Σ_i from the static analysis. This “initialisation” again recognises the minor departure from the standard exchangeable model in adapting to the perceived day-of-week similarities through the parameter means m_{ij0} ; otherwise, we would simply set m_{ij0} equal to the posterior mean of μ_i from the static analysis of the initial nine days of data.

Proceeding sequentially through the day, $t = 1, 2, \dots$, we sequentially receive and analyse the observed flows across all links, and the analysis adaptively updates the posterior distributions for model parameters. At time t during the day, write D_{jt} for all data and information available at the time, so that $D_{jt} = \{D_{j,t-1}, Y_{jt}\}$ where Y_{jt} is the full set of observed flows on all links in time interval (here minute) t . Based on this data, the posterior distributions for link i parameters are

$$N(\theta_{ijt}|m_{ijt}, \sigma_i^2 C_{ijt}) \times G(\sigma_i^{-2}|n_{ijt}/2, n_{ijt}S_{ijt}/2)$$

where the defining quantities are updated sequentially. The values m_{ijt} and S_{ijt} provide “on-line” point estimates of the model parameters θ_{ijt} and σ_i^2 at this point t during day j .

Key to real-time implementation of these models are step ahead forecast distributions, i.e., forecasts for near-term future flows based on past data. These are easily computable using standard dynamic model theory. Of primary interest are the one-step ahead forecast distributions, namely $p(y_{ij,t+1}|D_{jt})$ on link i . At time point t , this provides the one-step ahead forecast or prediction of flows in the next time interval. In addition to the interest in real-time application, these distributions are the cornerstones of methods for assessing model fit and adequacy; consistency of observed flows with these forecast distributions is a more critical and meaningful benchmark of model quality than the “retrospective” or fitted

values earlier explored (West and Harrison 1997, especially chapters 10 and 11), which we nonetheless document in Figure 15 and 16. Under the current model structures, these are Student T distributions with trivially computed moments (see appendix).

Figures 13 and 14 display graphs of observed flows for all links under a model as described above and with a discount factor of $\delta = 0.95$ determining the variance matrices $W_{ijt}(\delta)$ that describe the time-variation in the link-specific regression parameters β_{ijt} . The sequential adaptation to new data is clear and most apparent on link 1 where we see the estimated spline-based trend modified throughout the day. By comparison with the earlier figures of fitted values note that one-step forecasts inevitably appear less “on target” – they are forecasts not retrospective fits. In fact the dynamic models are indeed preferred in terms of model fit, some aspects of which are summarised through comparisons of model likelihood measures (West and Harrison 1997, especially sections 2.6, 3.4 and 11.4). In this setting the model likelihood is simply the product of observed one-step ahead forecast densities, i.e., the product of all terms $p(y_{ij,t+1}|D_{jt})$ over $t = 0, \dots, T - 1$. Running separate analyses that differ only through the value of the discount factor δ allows comparison of the values and so an assessment of data-based support for differing degrees of variability in the regression parameters. We summarise the results here comparing two discount factors $\delta = 1$ (the static model, no variations in parameters) and $\delta = 0.95$. Table 1 quotes the log-likelihoods for each link separately and for two representative days, July 2nd and August 8th. Evidently, this supports the dynamic model across all links.

We also show in Figure 15 and 16 the smoothed (i.e. retrospective) means of the distribution of the flows. This type of representation is more directly comparable to the fitted series in Figure 9 and 10. More so than the one-step-ahead forecast, necessarily less “on-target” as mentioned above. Clearly these fits appear to follow the observed counts very satisfactorily.

Further illustration of the benefits of dynamic modelling in this context arises in cases of marked changes in traffic flows. Figure 17 plots flows and out-of-sample fitted values on Wednesday, 31st July at detector 4. Evidently the standard model does very poorly in predicting the major change in patterns of flows in the mid-morning. By comparison, the dynamic model is capable of rapidly adapting to pick-up the changes and thereafter

adequately track and predict near-term flows.

4. INCORPORATING DOWNSTREAM DATA

Additional model extensions to explore the incorporation of downstream flow data have been developed and analysed in a similar fashion. Particularly in the cases of significant slowdowns of flows at a given link, it is evident that the near-term development of flows on the immediate upstream link will likely be heavily impacted. Hence it is of interest to extend the regression components of existing models to include observed downstream flows. This has been done, link-by-link, including past downstream flows one or two minutes back in time, using only flows at contiguous links. Over most of the links and on the basis of analyses across the same nine days of flow data, the improvements in model fit and short-term forecast performance are almost negligible. Figures 18 and 19 provide graphs of the observed flows at detectors 1 and 3 over a selection of days. Superimposed are fitted values from both the original model, based only on upstream flows, and the extended model that incorporates recent downstream data. Evidently, the differences are very small.

The story is not quite so negative, however. Recall that the original model is relatively poor in tracking flows on one of the Wednesdays of the initial set of days. It turns out that adding in regression terms on downstream data does have a more apparent impact, in terms of slightly improved fit, for data on one link, link 3, on Wednesdays. On one specific Wednesday (7th August), this link experienced a radical slowdown in the early rush hour, and this is very well predicted on the basis of the feedback from the downstream link, in this case representing a very substantial improvement over the “upstream only” model; see Figure 20. It is unclear just why this link:day combination stands out in this respect, but the results here do indicate the potential for real-time presaging of major slowdowns in this model extension.

5. SUMMARY COMMENTS

The link-specific regression and dynamic regression models identified in our analyses account for a good deal of the observed variability in minute-by-minute traffic flows as recorded by existing detectors. At this time resolution, basic spline functions for the day-specific

profiles combined with selected upstream flows from very recent time intervals reasonably adequately characterise changes in traffic flow patterns for current time and one-step ahead predictions. The daily patterns reproduce across common days of the week and so historical data may be used to refine the estimation of the day-of-week specific spline parameters. Modified models that permit day- and minute- specific changes in the regressions on recent upstream flows quite significantly improve the accuracy of short-term forecasts of flows through adaptation to more local conditions. These dynamic regressions are recommended for empirical modelling and short-term prediction and assessment of changes in flow patterns. Our analyses have identified a specific link-day combination that departs significantly from this otherwise adequate model across other links and days, and for which additional explanations are needed. In general, the incorporation of recent flows on downstream links does not markedly improve model fit and predictive ability, with the very notable exception of a day on which flows break down quite substantially and the feedback of the slowdown to upstream links is evidently strongly predicted by the use of recent downstream data. Hence the extended models incorporating downstream data should generally be entertained in anticipation of their relevance in cases of extreme slowdowns.

Open questions and areas for further development include considerations of flow predictions several minutes ahead. In the current I5 study, the models do not have this capacity, and additional detectors for on-ramp traffic would be needed to refine the models usefully in this direction. They do, however, offer potential for learning through “What if?” analyses that explore predictions under assumed patterns of on-ramp flows. A second area for further study concerns more formal modelling of the day-of-week effect. This may involve developing and estimating hierarchical models that describe correlations between the spline parameters across days of the week, and that permit marked departures for idiosyncratic days to respond, for example, to significant changes in weather patterns. Further areas of research that should prove fruitful include extended models that correlate changes in the dynamic regression parameters across links, and incorporation of additional predictors such as occupancy measures (if available from detector data) to proxy local average speed of traffic. Finally, the potential of such models as management tools will be enhanced by the overlaying of feed-forward intervention methods (West and Harrison, 1997) to mod-

ify model estimates and forecasts in the light of external information, such as changes in weather conditions, breakdowns, lane or ramp closures, and so forth.

6. APPENDIX

6.1. FITTING HIERARCHICAL REGRESSION MODELS

The structure of conditional posterior distributions in the basic hierarchical regression model, and the resulting implementation of Markov chain Monte Carlo methods of analysis, is summarised here. We implement a direct Gibbs sampling analysis in which each of several sets of parameters are simulated at each stage of an overall iterative framework. The parameters are simulated from conditional distributions that are determined by currently fixed values of other parameters, and which changes as these other parameters are resampled from their conditional distributions.

The specific set of conditional distributions used here are as follows. Recall that we have $I = 6$ links, $J = 9$ days and $T = 228$ minutes per day. The distributions are given for each detector $i = 1, \dots, I$, with the understanding that the corresponding link-specific parameters are conditionally independent with the stated distributions.

- For each of the θ_{ij} , we have the multivariate normal distribution $N(\theta_{ij}|t_{ij}, T_{ij})$ where

$$t_{ij} = T_{ij}(\Sigma_i^{-1}\mu_i + X_{ij}'Y_{ij}\sigma_i^{-2})$$

and

$$T_{ij} = (\Sigma_i^{-1} + X_{ij}'X_{ij}\sigma_i^{-2})^{-1}.$$

- For each of the μ_i , we have the multivariate normal distribution $N(\mu_i|m_i^*, M_i^*)$ where

$$m_i^* = M_i^*(M_i^{-1}m_i + J\Sigma_i^{-1}\bar{\theta}_i)$$

and

$$M_i^* = (M_i^{-1} + J\Sigma_i^{-1})^{-1}.$$

Here $\bar{\theta}_i = \sum_{j=1}^J \theta_{ij}/J$.

- For each of the Σ_i , we have the inverse Wishart distribution

$$W(\Sigma_i^{-1}|((\rho + J)R_i)^{-1}, (\rho + J))$$

where

$$R_i = (\rho R + \sum_{i=1}^d w_{ij} w'_{ij}) / (\rho + J)$$

with $w_{ij} = \theta_{ij} - \mu_i$.

- For each σ_i^2 , we have the conditional inverse gamma distribution

$$G(\sigma_i^{-2} | (\nu + JT)/2, (\nu + JT)\tau_i^2/2)$$

where

$$\tau_i^2 = (\nu\tau^2 + \sum_{j=1}^J e_{ij}e'_{ij}) / (\nu + JT)$$

with $e_{ij} = Y_{ij} - X'_{ij}\theta_{ij}$.

The Gibbs sampler cycles repeatedly through these four items in turn, at each step simulating a new value of the corresponding parameters to be used in future stages. The result of the iterative sampling, after eliminating a number of initial iterations devoted to the initial burn-in process, are values for the variables drawn approximately from their full joint distribution.

The results of this paper are based on the output of 10,000 iterations, the first 1,000 discarded for burn-in. The code for the Gibbs sampler was written in FORTRAN 77, and when compiled with a standard optimization option, runs in 25 to 45 seconds on a Sun Ultra 60 2360 otherwise idle. Convergence and independence from the starting values were checked by standard tools (i.e. CODA, Best et al., 1995). The starting values were randomly generated from the prior distributions adopted.

As for sensitivity to the prior assumptions, we point out that we adopt non-informative priors, as can be assessed for example in Figure 2. This way we are truly letting the data guide the posterior analysis. More generally, this kind of hierarchical analysis for linear models is well consolidated and has been thoroughly analysed with respect to this issue, and we refer again to Lindley and Smith (1972), Smith (1973), Gelfand et al (1990), and chapters in Gilks et al (1996) for details.

6.2. FITTING DYNAMIC HIERARCHICAL REGRESSION MODELS

The details of dynamic model analysis are standard, following general theory in West and Harrison (1997, chapter 4). Assume the models of equations (4) and (5), and note that the first equation can be expressed as $y_{ijt} = x'_{ijt}\theta_{ijt} + \epsilon_{ijt}$. Assume initial priors $p(\theta_{ij0}, \sigma_i^2 | D_{j0})$ of the conjugate form $N(\theta_{ij0} | m_{ij0}, \sigma_i^2 C_{ij0})$ multiplied by $G(\sigma_i^{-2} | n_{ij0}/2, n_{ij0} S_{ij0}/2)$. For every link i this structure holds independently of other links j , for $j \neq i$, so analysis proceeds in parallel across links, though the flow data is “transferred” across links via the regression terms. Beginning at time $t = 1$, the data is sequentially analysed and at any time t the following component distributions are easily calculated (for full details of the updating equations and other components of analysis, see West and Harrison 1997, chapter 4).

- $p(\theta_{ijt}, \sigma_i^2 | D_{jt})$ has the conjugate form

$$N(\theta_{ijt} | m_{ijt}, \sigma_i^2 C_{ijt}) \times G(\sigma_i^{-2} | n_{ijt}/2, n_{ijt} S_{ijt}/2)$$

where

$$\begin{aligned} m_{ijt} &= m_{ij,t-1} + R_{ijt} x_{ijt} e_{ijt} / q_{ijt}, \\ C_{ijt} &= R_{ijt} - R_{ijt} x_{ijt} x'_{ijt} R'_{ijt} / q_{ijt}, \\ R_{ijt} &= C_{ij,t-1} + W_{ijt}, \\ q_{ijt} &= x'_{ijt} R_{ijt} x_{ijt} + 1, \\ e_{ijt} &= y_{ijt} - x'_{ijt} m_{ij,t-1}, \\ n_{ijt} &= n_{ij,t-1} + 1, \text{ and} \\ S_{ijt} &= S_{ij,t-1} (1 + e^2_{ijt} / (q_{ijt} n_{ijt})). \end{aligned}$$

- $p(\theta_{ijt} | D_{jt})$ is the density of a multivariate Student T distribution with n_{ijt} degrees of freedom, namely

$$T_{n_{ijt}}(\theta_{ijt} | m_{ijt}, S_{ijt} C_{ijt}).$$

- Forecasting one-step ahead at time t , the one-step ahead forecast or predictive distribution for $(y_{ij,t+1} | D_{jt})$ is a univariate Student T distribution, namely

$$T_{n_{ijt}}(y_{ij,t+1} | x'_{ij,t+1} m_{ijt}, S_{ijt} q_{ij,t+1}).$$

- $p(\theta_{ijt-k}|D_{jt})$, the k -step smoothed distribution, is a univariate Student T distribution, namely

$$T_{n_{ijt}}(\theta_{ijt-k}|a_{ijt}(-k), \frac{S_{ijt}}{S_{ijt-k}}R_{ijt}(-k))$$

where

$$\begin{aligned} a_{ijt}(-k) &= m_{ij,t-k} + B_{ijt-k}(a_{ijt}(-k+1) - a_{ijt-k+1}), \\ a_{ijt}(0) &= m_{ijt}, \\ a_{ijt-k}(1) &= a_{ijt-k+1}, \\ a_{ijt} &= m_{ijt-1}, \\ R_{ijt}(-k) &= C_{ijt-k} + B_{ijt-k}(R_{ijt}(-k+1) - R_{ijt-k+1})B'_{ijt-k}, \\ R_{ijt}(0) &= C_{ijt}, \\ R_{ijt-k}(1) &= R_{ijt-k+1}, \text{ and} \\ B_{ijt} &= C_{ijt}R_{ijt+1}^{-1}. \end{aligned}$$

The lines in Figure 15 and 16 connect the values $f_{ijt}(-k) = x_{ijt-k}a_{ijt}(-k)$ where t is the maximum time recorded (the 228th minute) and k runs from $t - p_i - 1$ to 1, (p_i being the number of predictors used in the regression model for link i).

7. ACKNOWLEDGEMENTS

Research supported in part by National Science Foundation under grant DMS-9313913 to the National Institute of Statistical Sciences, Research Triangle Park NC, (www.niss.org).

8. REFERENCES

- Best, N.G., Cowles, M.K. and Vines, S.K., "CODA: Convergence diagnosis and output analysis software for Gibbs sampler output" (Version 0.3). Medical Research Council Biostatistics Unit, Cambridge.
- Gelfand, A.E., Hills, S.E., Racine-Poon, A. and A.F.M. Smith, 'Illustration of Bayesian inference in Normal data models using Gibbs sampling', *Journal of the American Statistical Association*, **85** (1990), 972-985.
- Gilks, W.R., Richardson, S. and D.J. Spiegelhalter, *MCMC in Practice*, (1996), Chapman and Hall: London.
- Graves, T.L., Karr, A.F., Roupail, N.M. and P. Thakuriah, P., 'Predicting incipient congestion on freeways from detector data', (1997), Technical Report, National Institute of Statistical Sciences, Research Triangle Park, NC.
- Lindley, D.V. and A.F.M. Smith, 'Bayes estimates for the linear model' (with discussion), *Journal of the Royal Statistical Society (Ser. B)*, **34** (1972), 1-41.
- Smith, A.F.M., 'A general Bayesian linear model', *Journal of the Royal Statistical Society (Ser. B)*, **35** (1973), 1-41.
- West, M. and P.J. Harrison, *Bayesian Forecasting and Dynamic Models*, (1997), 2nd edn. Springer Verlag: New York.
- Whittaker, J., Garside S. and Lindveld K., 'Tracking and Predicting a Network Traffic Process', *International Journal of Forecasting*, **13** (1997), 51-61.

Authors' biographies:

Claudia Tebaldi is Data Modeler in the Advanced Technology Group at Athene Software inc, Boulder, CO. Her recent work has been in regression and classification using nonlinear techniques, Bayesian model selection, with application in the atmospheric sciences. Most recently, she is working on data mining and statistical modelling with large data sets.

Mike West is the Arts and Sciences Professor of Statistics and Decision Sciences, and Director of the Institute of Statistics and Decision Sciences, Duke University. Recent publications include regression models for large data sets and high dimensional covariates,

latent structures in multivariate time series, multiresolution models for time series and spatial processes, semiparametric and simulation-based Bayesian methods. Areas of application include bioinformatics, natural sciences, finance.

Alan F. Karr is Director, National Institute of Statistical Sciences, and Professor in the Statistics and Biostatistics departments at the University of North Carolina at Chapel Hill. His current research focuses on the interface between statistics and information technology. Among the projects he manages at NISS are transportation, materials properties, software engineering and data confidentiality.

Figure 1 : The section of I-5 under study and detector layout. The numbers at the bottom represent the distances, in miles, between subsequent detectors.

Figure 2 : Detector 2: prior and posterior marginal distributions for components of the regression hyperparameter vector.

Figure 3 : Detector 1: posterior means of the spline functions. The solid line corresponds to the mean value of the hyperparameter vector; the dashed lines to the mean value of the parameters for the single days. The line showing a “dip” in the interval (50, 100) corresponds to the peculiar Wednesday mentioned in the analysis.

Figure 4 : Detector 2: posterior summaries for regression parameters, in the top frame and the posterior means of the spline functions, in the bottom frame. For the latter, the solid line corresponds to the mean value of the hyperparameter vector; the dashed lines to the the mean value of the parameters for the single days, with the two Wednesdays corresponding to the two lines at the top of the panel.

Figure 5 : Detector 3: posterior summaries for regression parameters, in the top frame and the posterior means of the spline functions, in the bottom frame. For the latter, the solid line corresponds to the mean value of the hyperparameter vector; the dashed lines to the mean value of the parameters for the single days, with the peculiar Wednesday corresponding to the line at the very top.

Figure 6 : Detector 4: posterior summaries for regression parameters, in the top frame and the posterior means of the spline functions, in the bottom frame. For the latter, the solid line corresponds to the mean value of the hyperparameter vector; the dashed lines to the mean value of the parameters for the single days, with the peculiar Wednesday corresponding to the line at the very top.

Figure 7 : Detector 5: posterior summaries for regression parameters, in the top frame and the posterior means of the spline functions, in the bottom frame. For the latter, the solid line corresponds to the mean value of the hyperparameter vector; the dashed lines to the mean value of the parameters for the single days, with the peculiar Wednesday corresponding to the line at the very top.

Figure 8 : Detector 6: posterior summaries for regression parameters, in the top frame and the posterior means of the spline functions, in the bottom frame. For the latter, the solid line corresponds to the mean value of the hyperparameter vector; the dashed lines to the mean value of the parameters for the single days, with the peculiar Wednesday corresponding to the line at the very bottom.

Figure 9 : Tuesday 18th June: Data (dotted lines) and fitted values (solid lines) for the entire sequence of detectors. From top to bottom, detector 1 through detector 6.

Figure 10 : Thursday 27th June: Data (dotted lines) and fitted values (solid lines) for the entire sequence of detectors. From top to bottom, detector 1 through detector 6.

Figure 11 : Tuesday 2nd July: Data (dotted lines) and out-of-sample fitted values (solid lines) for the entire sequence of detectors. From top to bottom, detector 1 through detector 6.

Figure 12 : Thursday 8th August: Data (dotted lines) and out-of-sample fitted values (solid lines) for the entire sequence of detectors. From top to bottom, detector 1 through detector 6.

Figure 13 : Tuesday 2nd July: Data (dotted lines) and one-step-ahead forecasts (solid lines) for the 6 detectors. From top to bottom, detector 1 through detector 6.

Figure 14 : Thursday 8th August: Data (dotted lines), and one-step-ahead forecasts (solid lines) for the 6 detectors. From top to bottom, detector 1 through detector 6.

Figure 15 : Tuesday 2nd July: Data (dotted lines) and smoothed estimates of the counts posterior means (solid lines) for the 5 detectors that utilise upstream data as explanatory variables. From top to bottom, detector 2 through detector 6.

Figure 16 : Thursday 8th August: Data (dotted lines) and smoothed estimates of the counts posterior means (solid lines) for the 5 detectors that utilise upstream data as explanatory variables. From top to bottom, detector 2 through detector 6.

Figure 17 : Comparison of performances between standard model and dynamic model. Out-of-sample data from Wednesday, 31st July. The real data are the dotted line.

Figure 18 : Flows at detector 1 for the in-sample (top 2 frames) and out-of-sample (bottom 2) days. Fitted/predicted values by the original model (smooth curves) and by the model that includes recent downstream flows. Dotted lines are the real data.

Figure 19 : Flows at detector 3 for the in-sample (top 2 frames) and out-of-sample (bottom 2) days. Fitted/predicted values by the original model and by the model that includes recent downstream flows. The real data are the dotted lines, the two solid lines are indistinguishable.

Figure 20 : Improving the predictions by downstream data inclusion (bottom frame) in a case of major slowdown.

Table 1 : Comparison of static ($\delta = 1$) and dynamic ($\delta = 0.95$) models.

A.M. data, Southbound direction

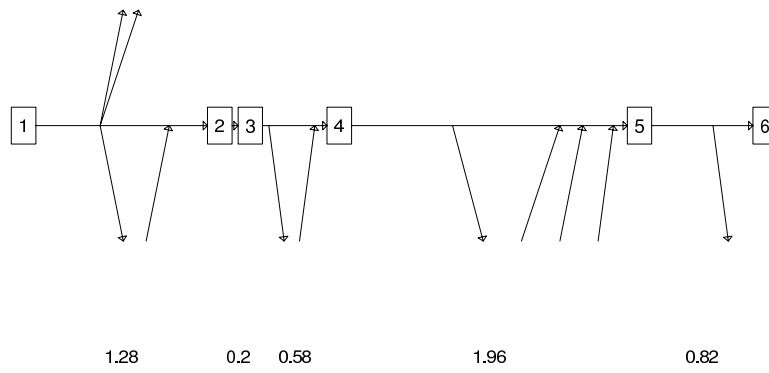


Figure 1:

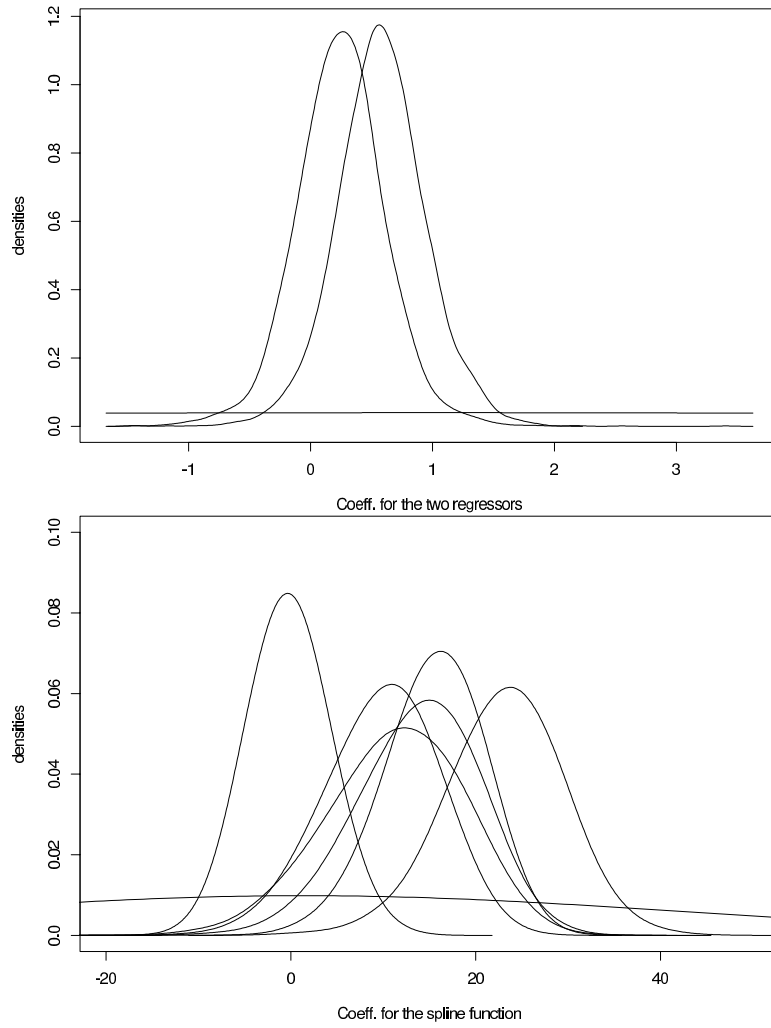


Figure 2:

the daily-mean splines and the overall-mean spline at detector 1

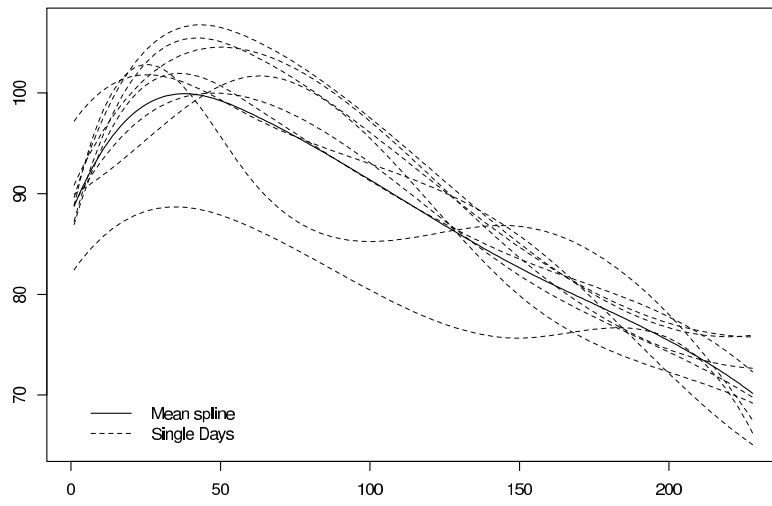


Figure 3:

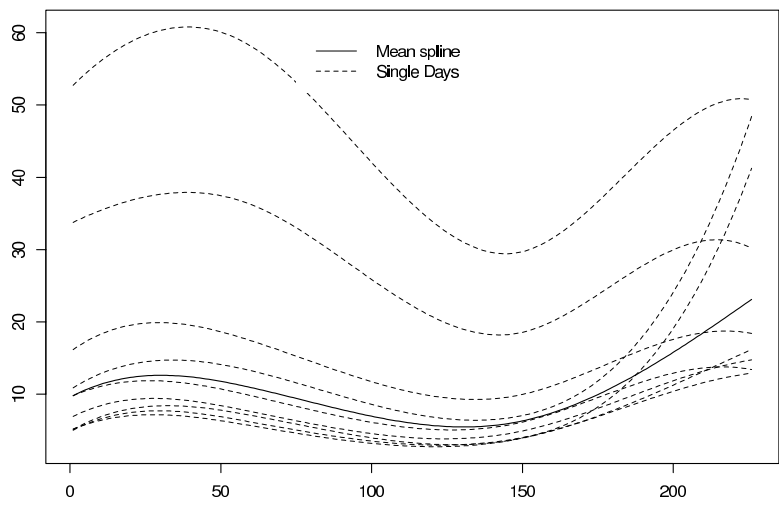
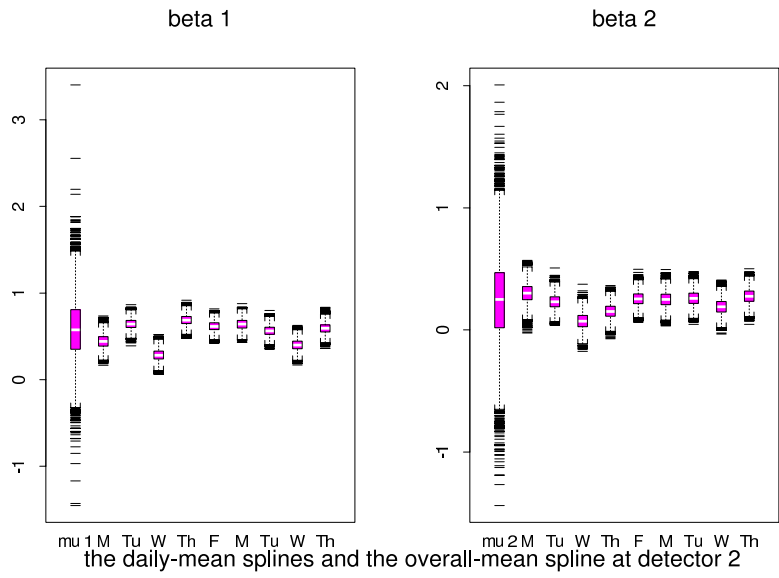
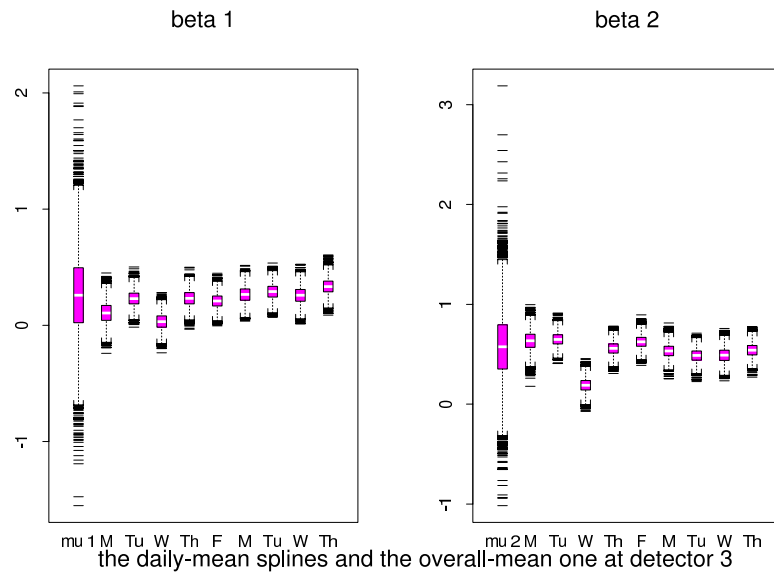


Figure 4:



the daily-mean splines and the overall-mean one at detector 3

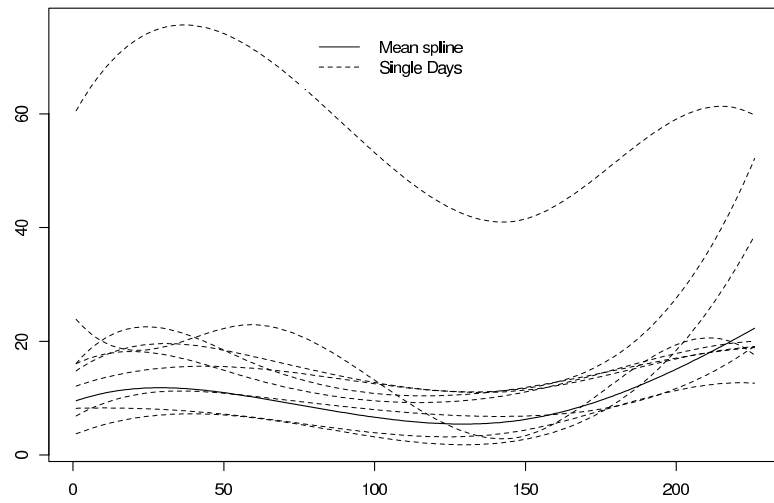
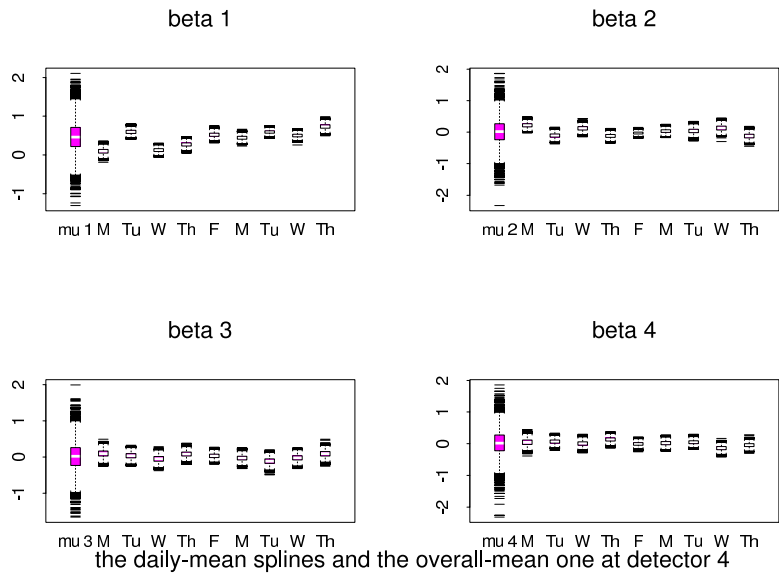


Figure 5:



the daily-mean splines and the overall-mean one at detector 4

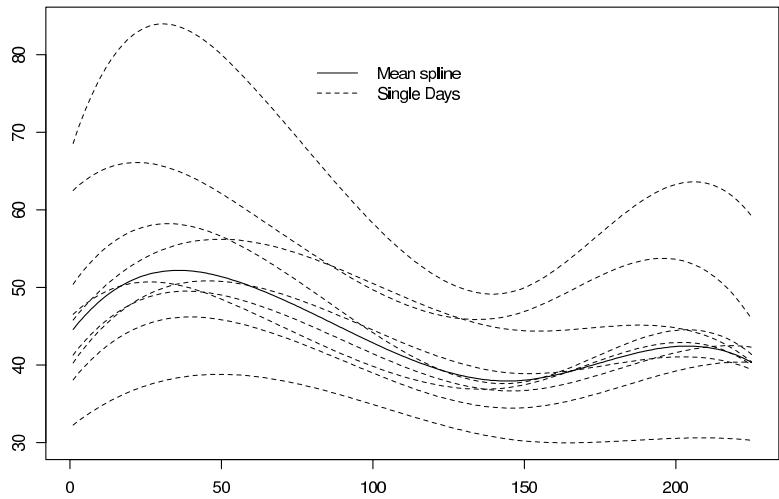


Figure 6:

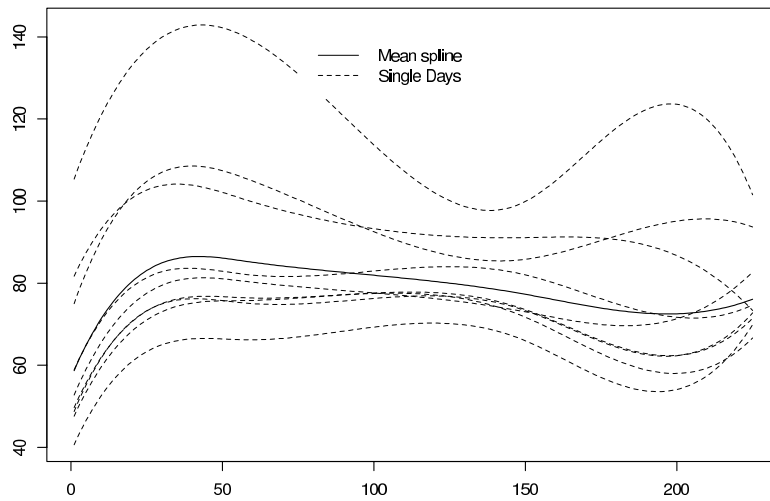
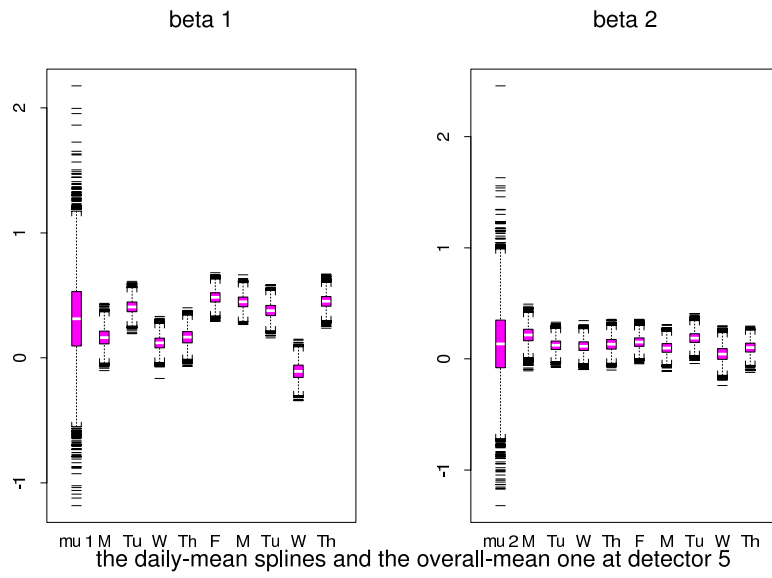


Figure 7:

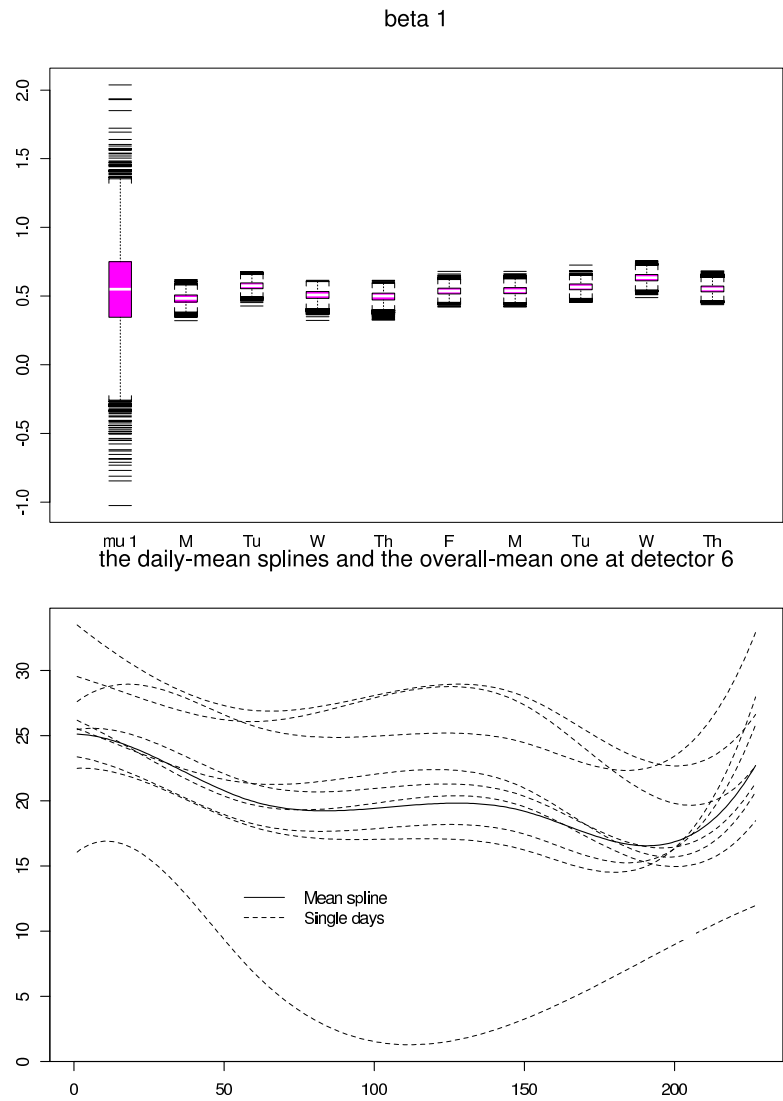


Figure 8:

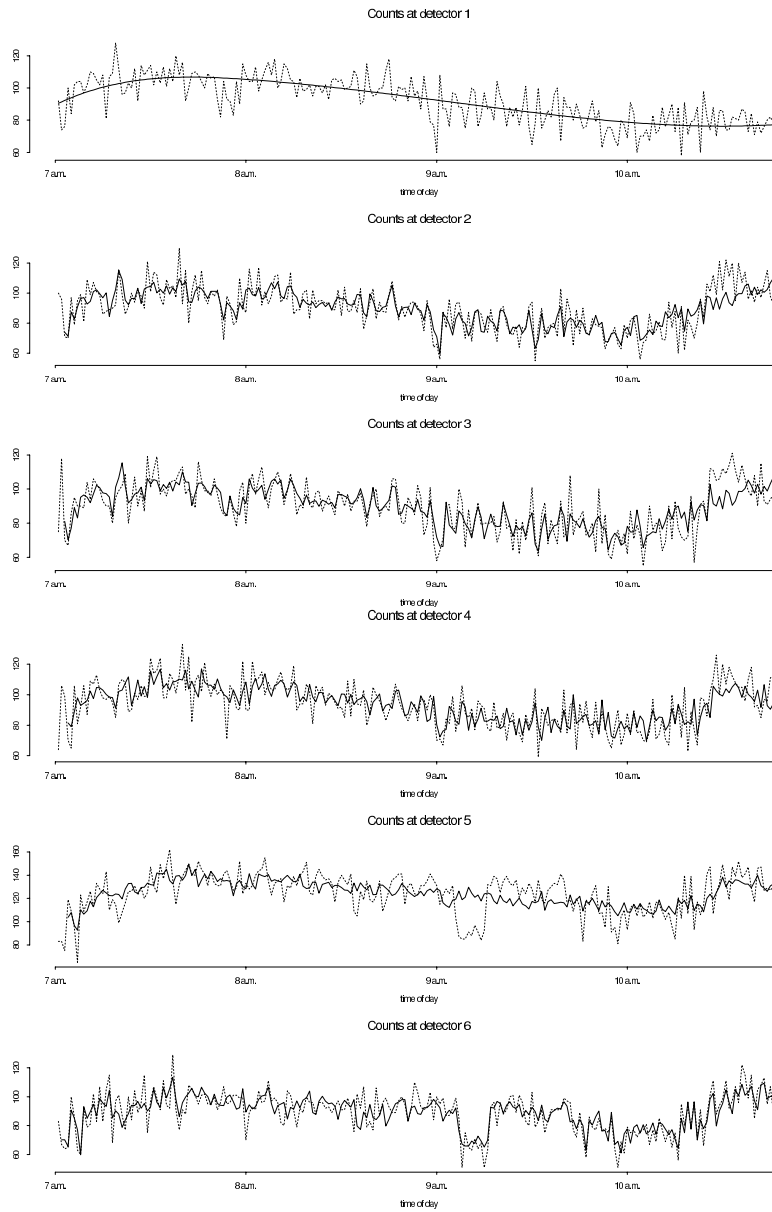


Figure 9:

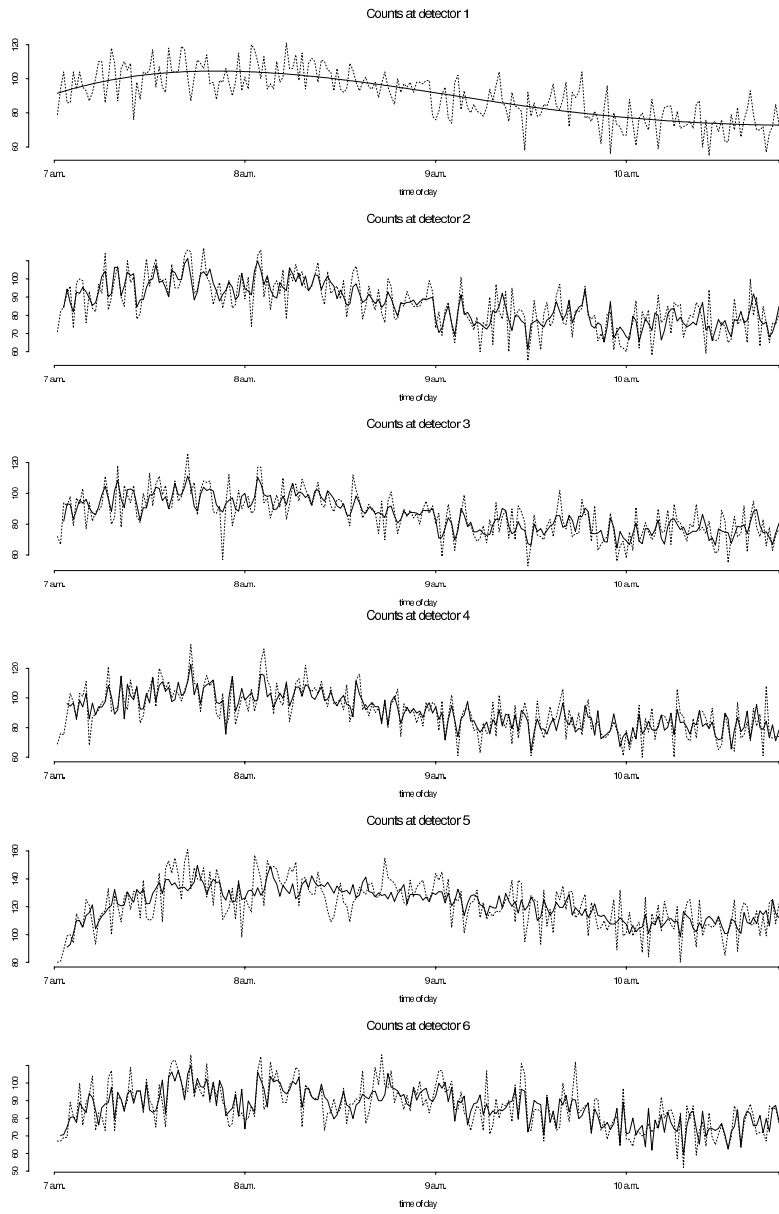


Figure 10:

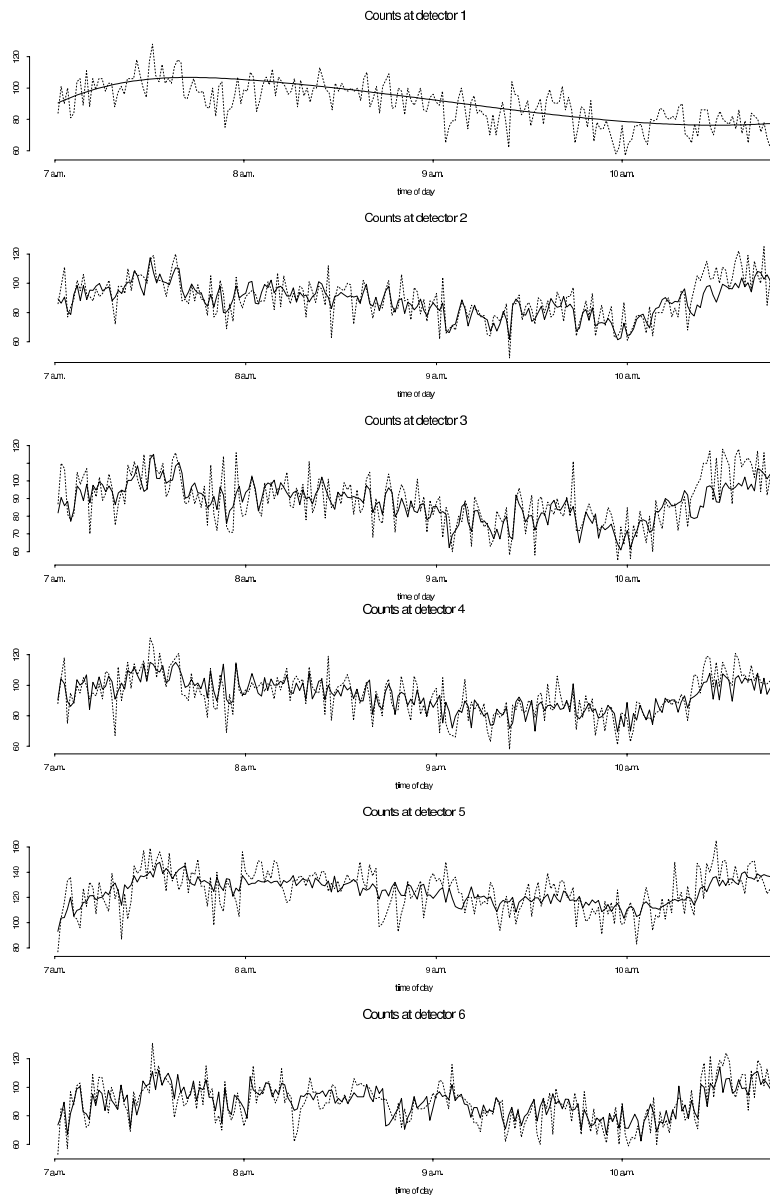


Figure 11:

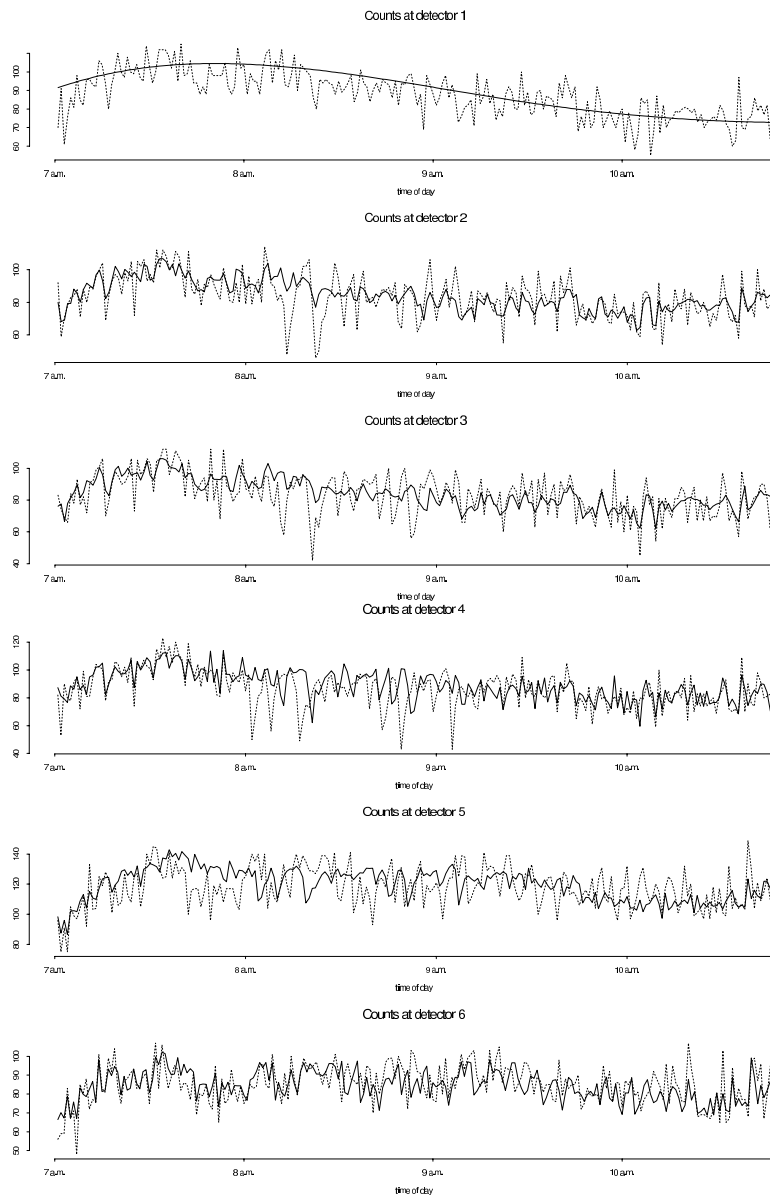


Figure 12:

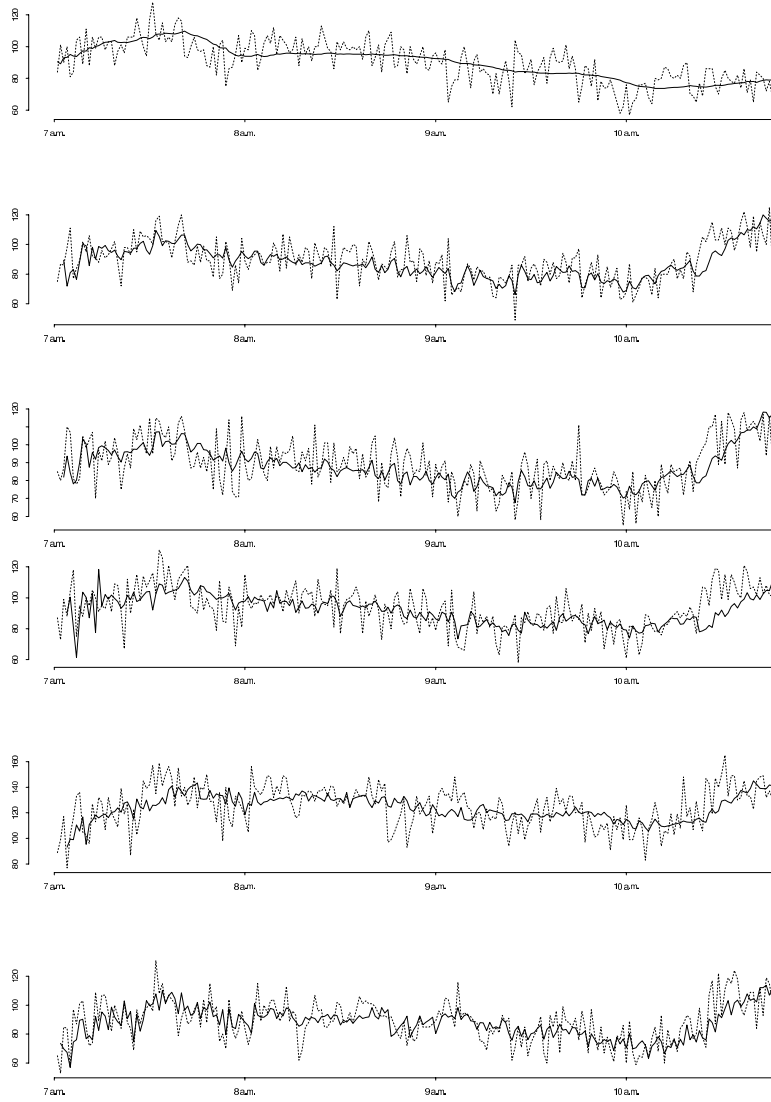


Figure 13:

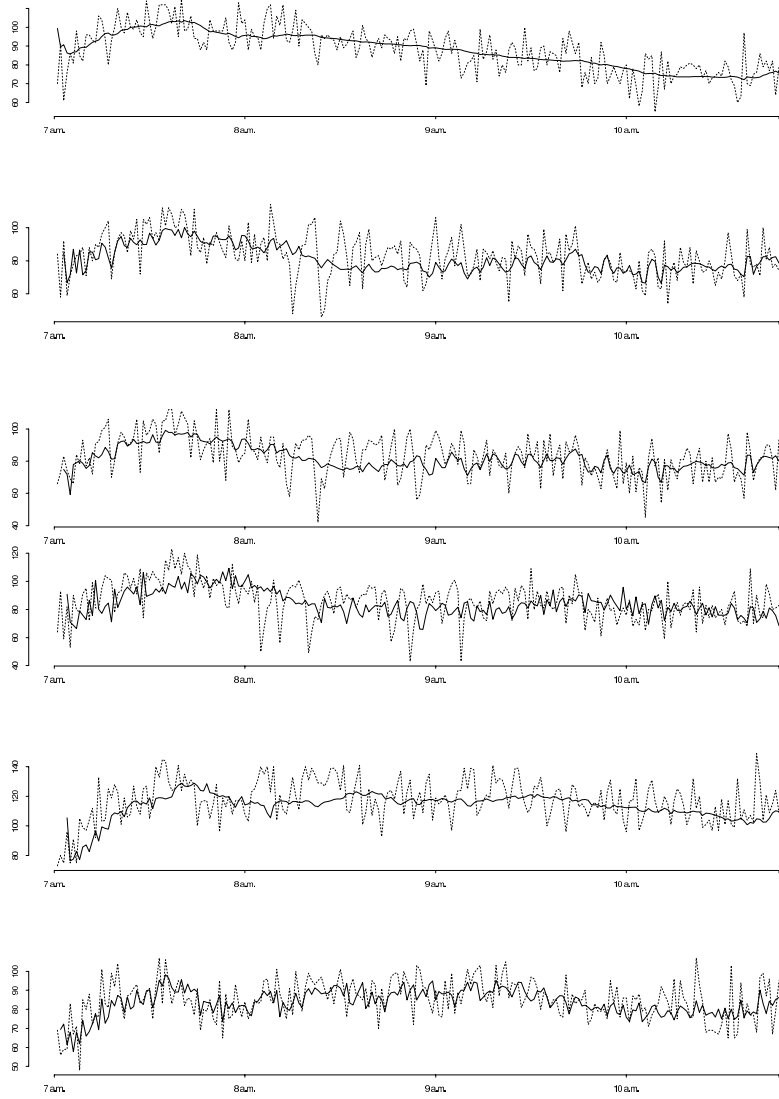


Figure 14:

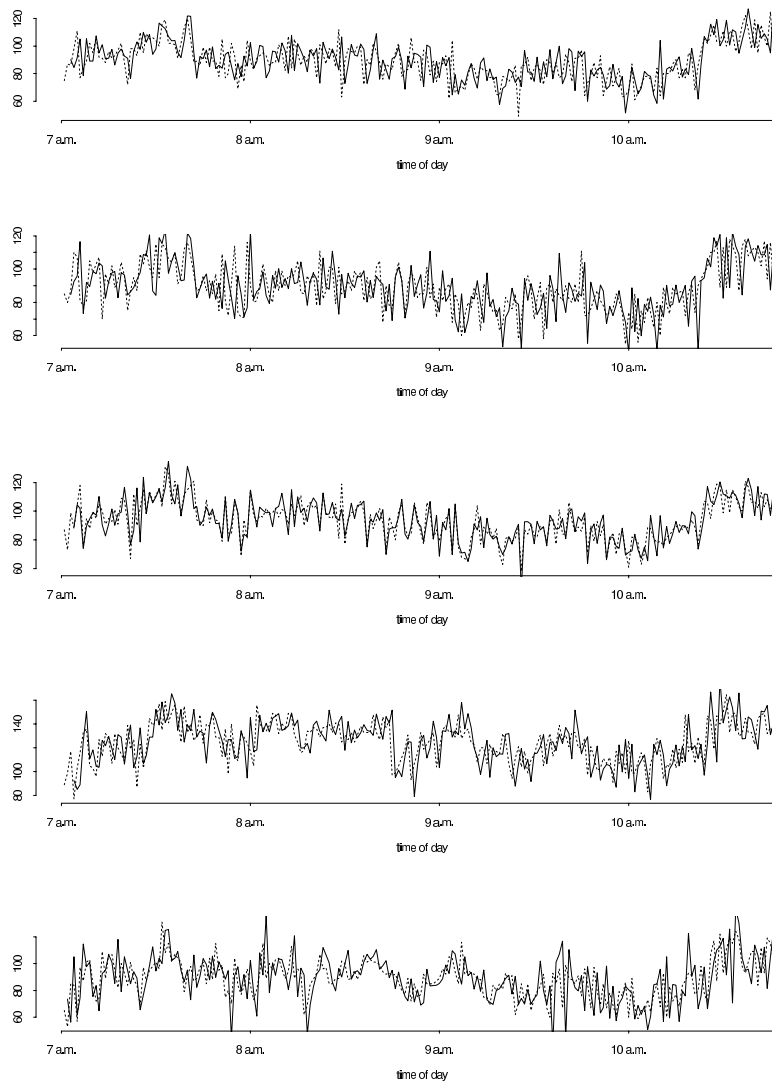


Figure 15:

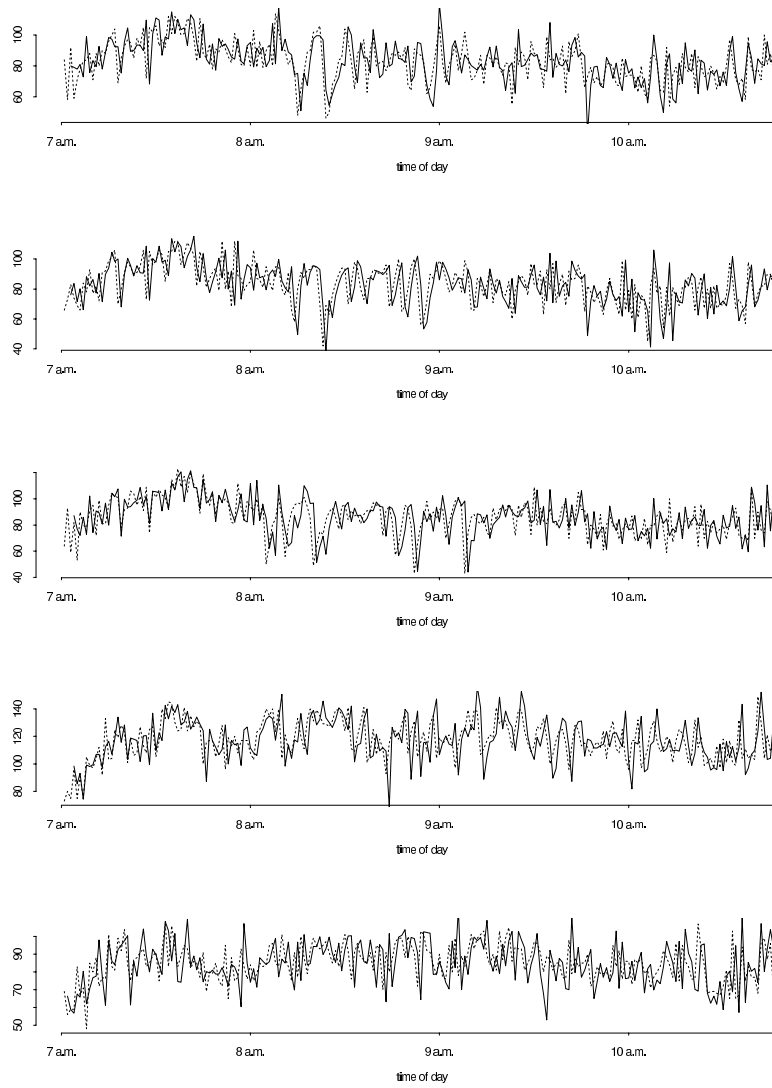


Figure 16:

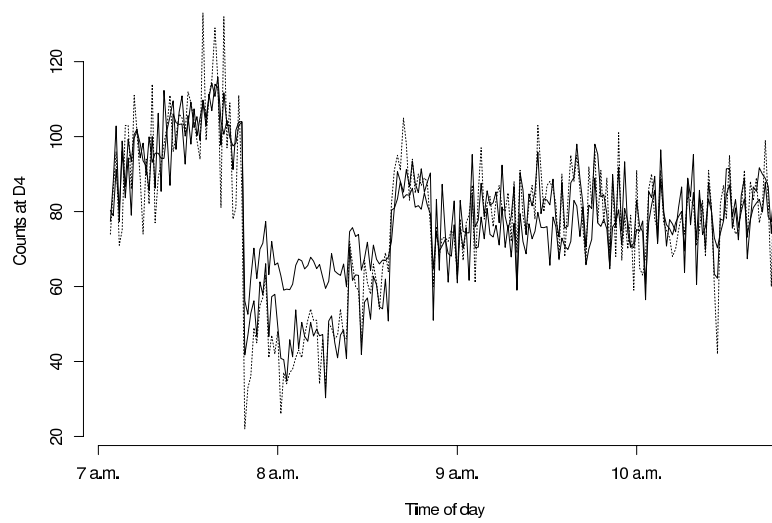


Figure 17:

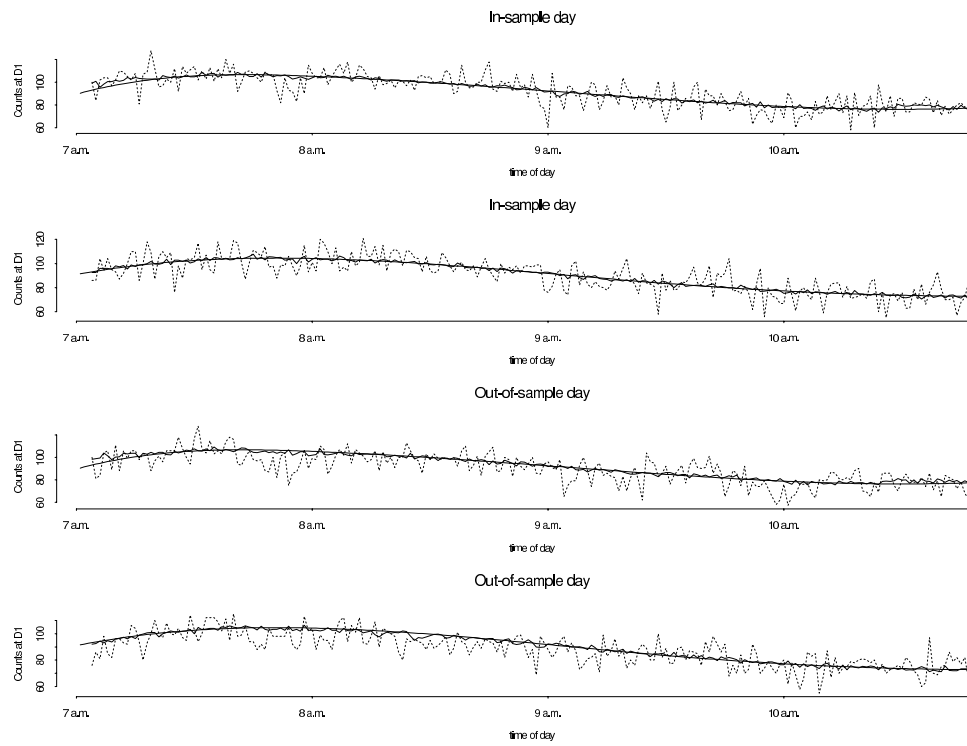


Figure 18:

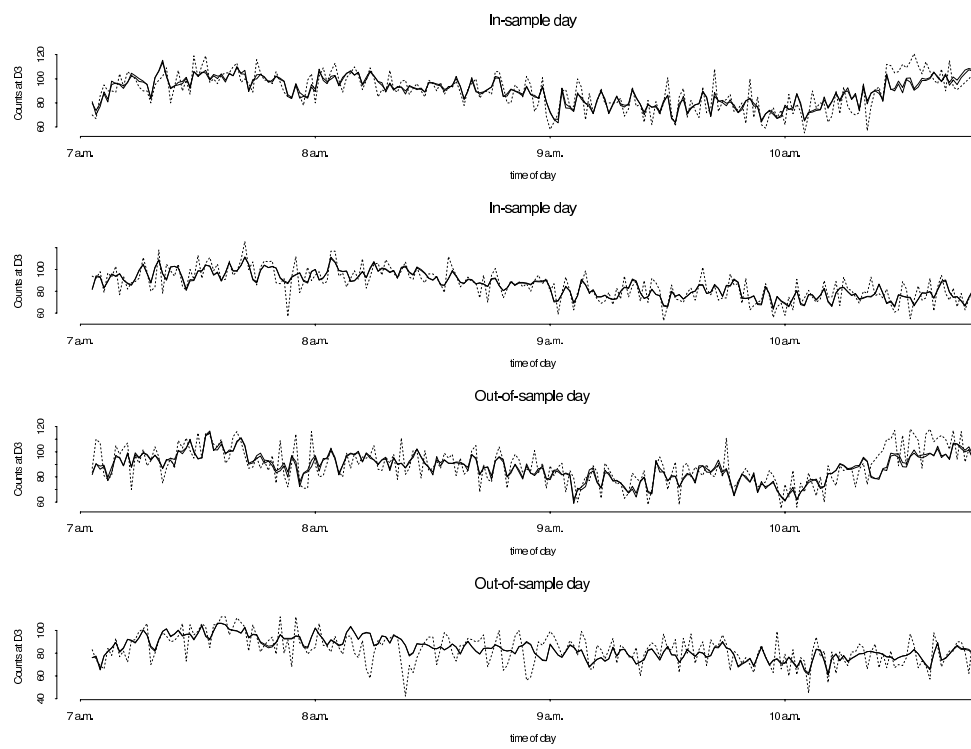


Figure 19:

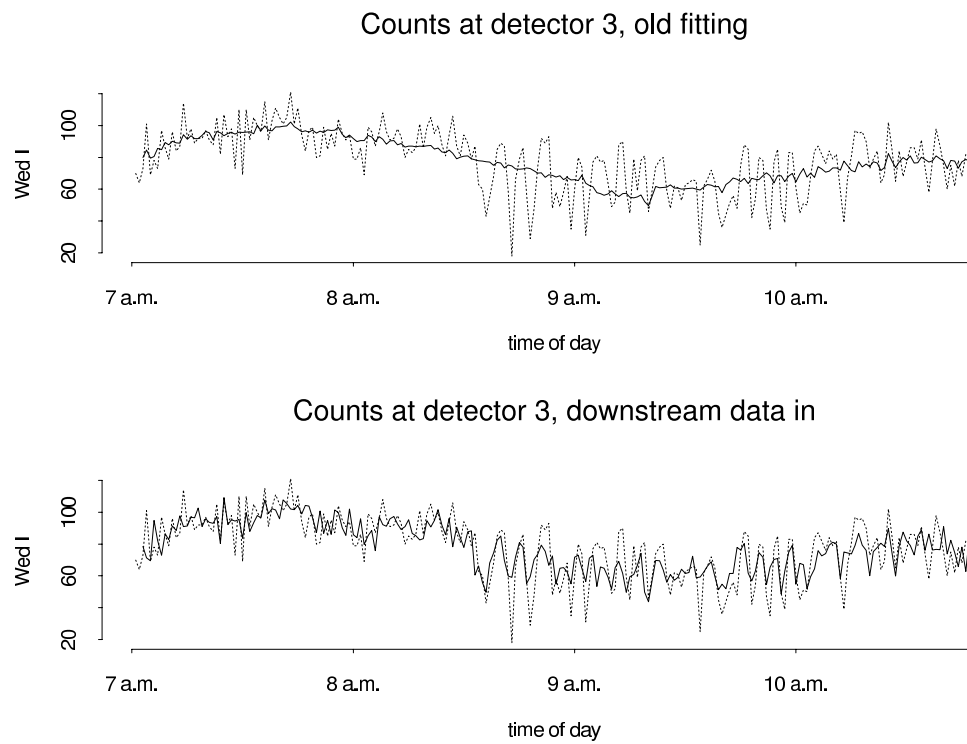


Figure 20:

Detector	log-likelihood	log-likelihood
07.02	$\delta = 1$	$\delta = 0.95$
1	-327.6	...
2	-331.5	-326.9
3	-329.0	-324.6
4	-327.4	-318.8
5	-326.5	-325.0
6	-338.4	-336.8

Detector	log-likelihood	log-likelihood
08.08	$\delta = 1$	$\delta = 0.95$
1	-327.1	...
2	-328.9	-327.6
3	-326.9	-325.9
4	-330.9	-329.3
5	-330.1	-327.6
6	-335.8	-335.5