

Interpolating Vancouver's Daily Ambient PM₁₀ Field

Li Sun

James V. Zidek

Nhu D. Le

Haluk Ozkaynak



NRCSE

Technical Report Series

NRCSE-TRS No. 033

The NRCSE was established in 1996 through a cooperative agreement with the United States Environmental Protection Agency which provides the Center's primary funding.



Interpolating Vancouver's Daily Ambient PM₁₀ Field

Li Sun¹, James V Zidek¹, Nhu D Le² and Halûk Özkaynak³

¹ University of BC, ² BC Cancer Agency and ³ US Environmental Protection Agency

January, 1999

Abstract

In this article we develop the first component of a spatial predictive distribution for the ambient space - time response field of daily ambient PM₁₀ in Vancouver, Canada. [That component deals with the prediction of daily levels, the second with hourly levels.] Observed responses have a remarkably consistent temporal pattern from one monitoring site to the next. We exploit this feature of the field by adopting a response model with two components, a common deterministic trend across all sites plus a stochastic residual. We are thereby able to whiten the temporal residuals without losing much of the spatial correlation in the original log-transformed series. This in turn enables us to develop an effective spatial predictive distribution for these residuals at unmonitored sites. By transforming the predicted residuals back to the original data scales we can impute Vancouver's daily PM₁₀ field for purposes such as human exposure and health impacts analysis.

KEY WORDS: PM₁₀; space-models; autoregressive processes; spatial interpolation; monitoring networks; spatial correlation.

1 Introduction.

This paper follows that of Li et al (1998) analyzing the hourly PM₁₀ field over the Greater Vancouver Regional District (GVRD). Our paper addresses problems arising from the spatial complexity of that field by turning from hourly to daily levels of this important air pollutant. We describe why our solution has been adopted and why it is satisfactory for certain purposes.

Interest in this pollutant derives from the recognition that elevated levels are associated with acute negative health impacts. A panel of experts appointed by the UK Department of the Environment, Transport and the Regions concludes: (<http://www.environment.detr.gov.uk/airq/aqs/particle/>, paragraph 25)

of PM₁₀ and health effects, ... that the higher the concentration of particles, the greater the effect on the health of the population and conversely, the lower the concentration, the smaller the effect.

Li et al (1998) analyze hourly ambient PM₁₀ concentrations collected in the Vancouver area from 1994 to 1996. Data come from 10 monitoring stations in the GVRD, different stations starting operation at different times. Tapered Element Oscillating Microbalance (TEOM) monitors generated the data, the data we use to construct the simple daily averages. (These "continuous" monitors use a tapered quartz element of conical shape. A detachable impervious filter is connected at the larger end and air is drawn onto that filter. The element oscillates at its resonant frequency when electrical current is passed through the element. However as the particle loading builds up that frequency changes unless the current is altered to maintain the resonant frequency. The change in current needed provides the surrogate measure of particulate concentration that gets converted and

Table 1: Cross-correlations ($\times 100$) for Hourly Average log PM10 de-AR'd Residuals Between Different Monitoring Stations.

site	1	2	3	4	5	6	7	8	9	10
1	–	16	24	22	11	23	19	16	11	9
2	16	–	17	15	24	22	17	12	7	5
3	24	17	–	17	9	23	22	14	7	8
4	22	15	17	–	12	23	25	27	15	10
5	11	24	9	12	–	14	17	12	9	7
6	23	22	23	23	14	–	29	14	10	9
7	19	17	22	25	17	29	–	18	16	10
8	16	12	14	27	12	14	18	–	17	13
9	11	7	7	15	9	10	16	17	–	15
10	9	5	8	10	7	9	10	13	15	–

Table 2: Cross-correlations ($\times 100$) For Hourly Average log PM10 De-trended Residuals Between Different Monitoring Stations.

site	1	2	3	4	5	6	7	8	9	10
1	–	44	61	55	37	59	53	48	41	41
2	44	–	43	45	58	53	46	40	35	29
3	61	43	–	54	31	63	57	45	33	37
4	55	45	54	–	40	62	67	65	48	39
5	37	58	31	40	–	40	41	40	35	30
6	59	53	63	62	40	–	67	50	40	38
7	53	46	57	67	41	67	–	55	45	36
8	48	40	45	65	40	50	55	–	48	45
9	41	35	33	48	35	40	45	48	–	45
10	41	29	37	39	30	38	36	45	45	–

averaged to yield the measurements. (Environmental Health Department, Warwick District Council, <http://www.warwickdceh.demon.co.uk/equip.htm#DataCapture>.)

The analysis of Li et al (1998) was to be a prelude to the development of a spatial prediction methodology for imputing unmeasured levels of PM10 at 299 additional locations. However, the intended interpolation methodology of Le and Zidek (1992) requires of the random field to be interpolated that its realizations: (1) have Gaussian distributions; and (2) be independent. Neither condition holds for our particulate fields. So the data were first subjected to logarithmic transformation to insure approximate attainment of (1). Denote the resulting response by $Y(x, t)$ at site x and time t .

Attaining (2) even approximately proves more challenging. First steps came from the discovery that the temporal pattern of the log-transformed measurements have a remarkable consistency across sites. So a trend model $T(t) = \mu + H_{hour} + D_{day} + W_{week}$ at time t was fitted across all ten monitoring sites. The “de-trended residuals” at site x and time t was then computed as $E(x, t) = Y(x, t) - T(t)$.

Next, autoregressive and other analyses of the $\{E(x, t)\}$ for each fixed x and varying t led to the adoption of a single AR(3) model for all sites:

$$E(x, t) = \alpha_1 E(x, t - 1) + \alpha_2 E(x, t - 2) + \alpha_3 E(x, t - 3) + e(x, t) \quad (1)$$

where α_1 , α_2 and α_3 are the model coefficients. For expository simplicity we refer to the $\{e(x, t)\}$ as the “de-AR’d” residuals to reflect the fact that the AR structure has been taken out. For fixed x , the de-AR’d residuals proved to be quite “white” for the 10 sites, having small auto-correlation.

Spatial prediction would then entail imputing the $\{e(x, t)\}$ for 299 sites x not among the original 10 to obtain say $\{\hat{e}(x, t)\}$ for such sites. Next the de-trended residuals for those sites would be constructed by taking them to be 0 at times $t = -2, -1$ and using (1) recursively to obtain the \hat{E} ’s for those sites. Finally the trend would be added to these imputed residuals to get series on the original log-PM10 scale.

However, the proposed procedure seemed likely to fail since for fixed t and the 10 stations, the de-AR’d residuals had not only the expected small auto-correlations but as well they had small between-site cross-correlations (see Table 1). These small correlations contrasted with their substantially larger counterparts for the $\{E(x, t)\}$ (see Table 2) Where had that correlation gone?

We found it had “leaked” into the lag-one hour cross-correlations between sites. Substantial correlation remained between $\{e(x, t)\}$ and $\{e(x', t - 1)\}$ at any two sites x and x' for varying t . This finding shows the 10 parallel time-series to have a complex spatial-temporal structure that cannot be modeled through univariate time-series methods applied one site at a time. Moreover, the seemingly obvious solution of using a multivariate auto-regressive approach cannot be used as we now explain.

The multivariate approach would involve use the de-trended residual series for all of the 10 sites in the model:

$$\mathbf{E}(t) = \mathbf{A}_1\mathbf{E}(t - 1) + \mathbf{A}_2\mathbf{E}(t - 2) + \mathbf{A}_3\mathbf{E}(t - 3) + \mathbf{e}(t) \quad (2)$$

where the $\{\mathbf{A}_i\}$ are the 10×10 matrices of model coefficients, $\mathbf{E}(t) = (E(1, t), \dots, E(10, t))'$ for all t and the $\{\mathbf{e}(t)\}$ are the multivariate de-AR’d residual vectors.

This approach would block the spatial correlation in the E-series from leaking into lag-1 spatial cross-correlation in the e-series. With the resulting increased cross-correlation in the lag-0 e-series success in constructing the \hat{e} series for the remaining 299 stations would be assured. To that extent the multivariate method would succeed.

However this approach makes difficult the construction of the \hat{E} series. Our 10×10 \mathbf{A} coefficient matrices would need to be extended to 309×309 matrices for all the sites. This seems to be an even more challenging problem than the interpolation of the series themselves. We see no way of solving it and alternative approaches are required.

The analysis of Li et al (1998) suggests three possible alternatives that we will now describe for completeness although a detailed description of the results obtained for the latter two must be deferred to a subsequent paper. The first is the subject of this paper and we turn from hourly to daily log average concentrations of PM10. Generally daily levels of airborne pollutants are of importance in environmental epidemiology where their association with morbidity and mortality are considered (*c.f.* Burnett et al 1994, Zidek et al 1998).

Since the de-trended hourly series was found to have an AR(3) structure it would be expected that successive de-trended log daily averages would have markedly smaller auto-correlations than their hourly counterparts. Theoretical considerations (see the Appendix) indicate that the possible lag 1 cross-correlations should be small as well. Thus the spatial-temporal complexities described above for the hourly series should be circumvented.

However for some purposes such as setting regulatory standards, hourly values may be needed. In this case one could treat each of the 24 hours one at a time. The AR(3) structure of Li et al (1998) suggests that these hourly values will be approximately independent since they are separated by 23 hours. Moreover we can fit different models for each hour thereby enabling us to capture such effects as the shifting wind field. That analysis will be further enhanced by fitting different models for the different seasons to allow for seasonal variations in the daily wind field. We present the results of that analysis elsewhere.

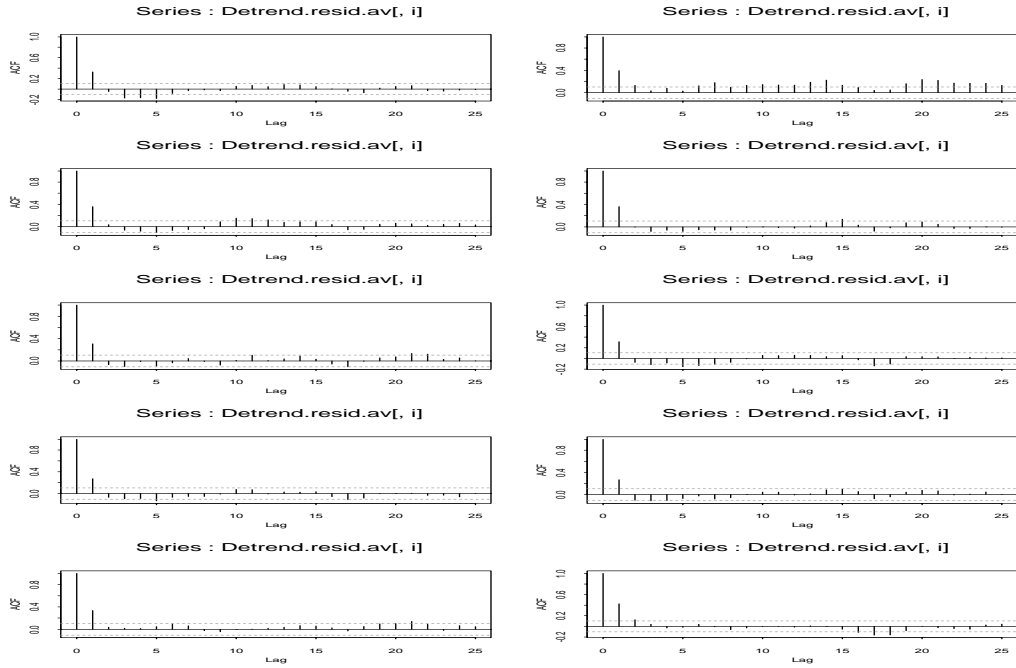


Figure 1: Auto-correlation Functions For De-trended Daily Log Averages At All Monitoring Sites.

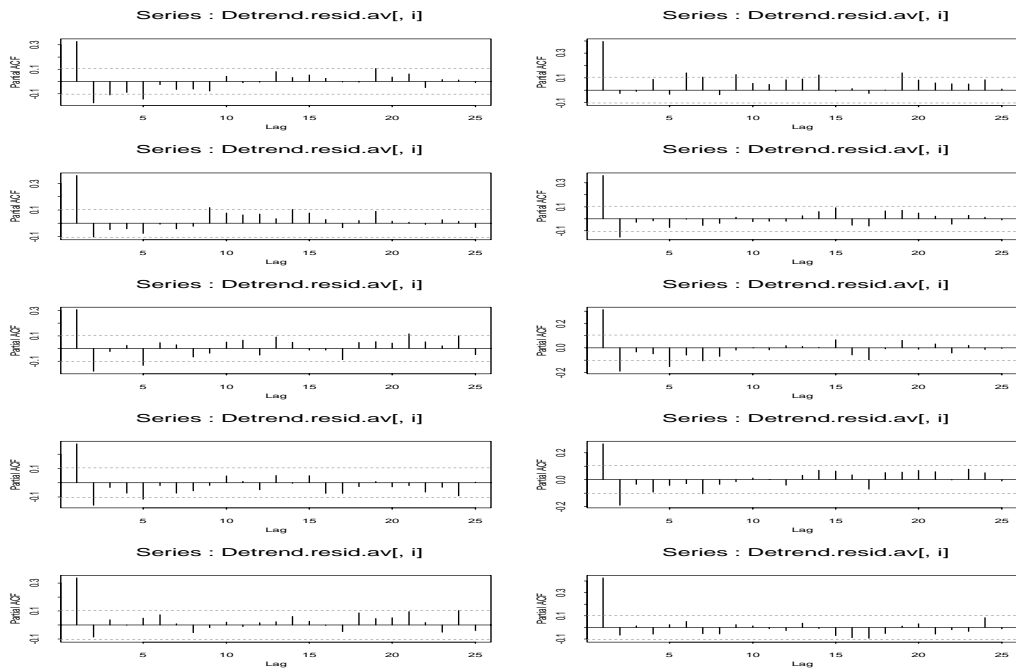


Figure 2: Partial Auto-correlation Functions For De-trended Daily Log Averages At All Monitored Sites.

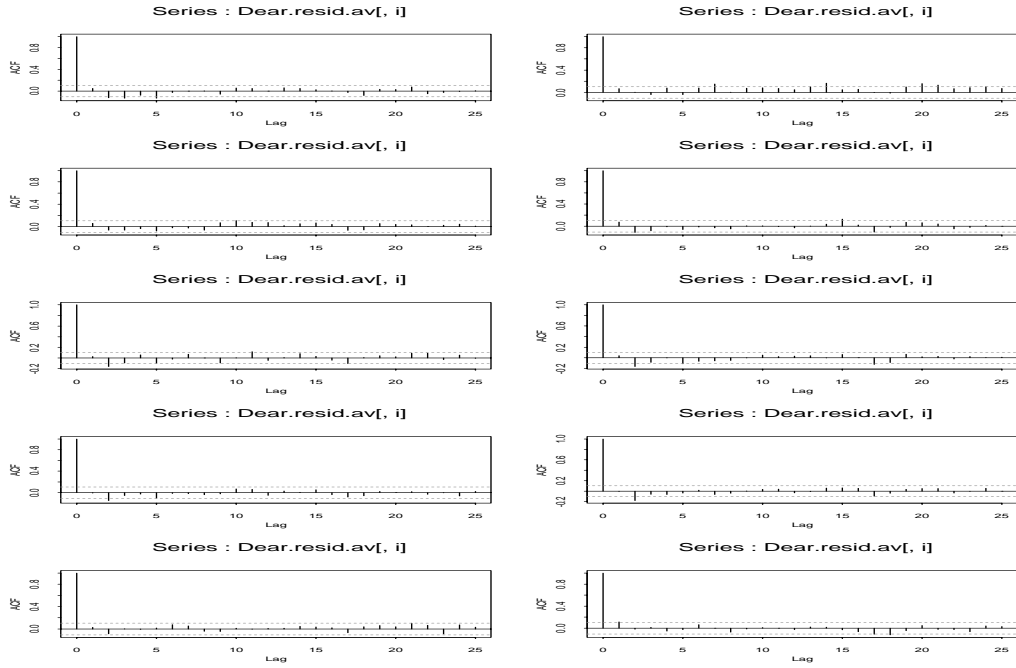


Figure 3: Auto-correlation Functions for AR(1) Residuals At All Monitoring Sites.

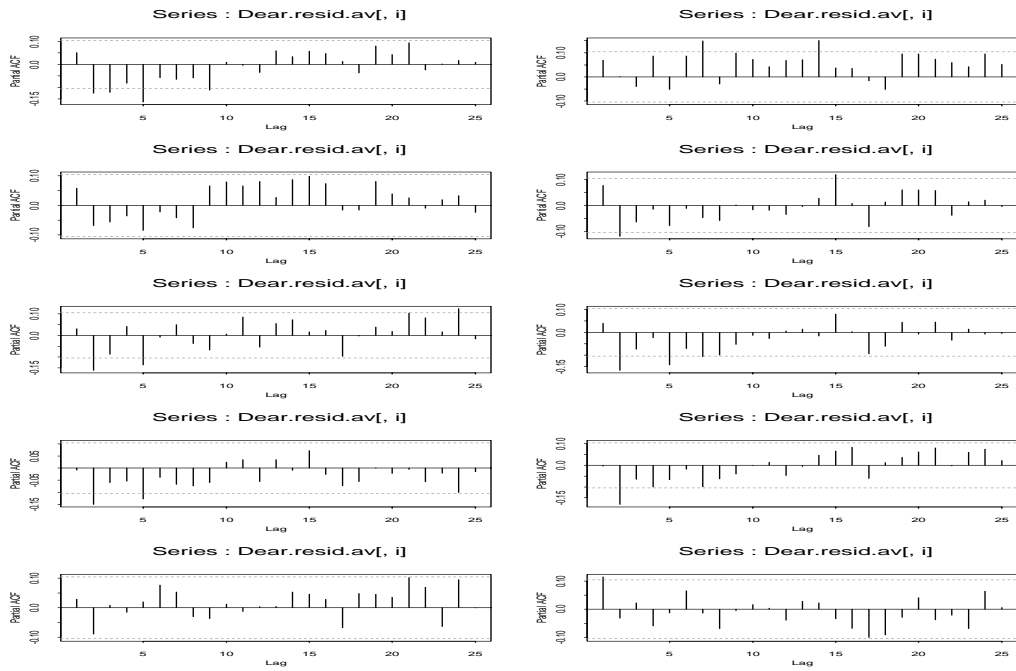


Figure 4: Partial Auto-correlation Functions for AR(1) Residuals At All Monitoring Sites.

Table 3: **Cross-correlation ($\times 100$) For De-trended Daily Average log PM₁₀ Residuals Between Different Monitoring Stations.**

site	1	2	3	4	5	6	7	8	9	10
1	–	61	81	75	60	78	75	74	64	63
2	61	–	56	58	79	67	59	60	49	44
3	81	56	–	76	51	84	77	72	54	59
4	75	58	76	–	62	84	87	86	67	53
5	60	79	51	62	–	62	59	62	54	48
6	78	67	84	84	62	–	85	79	62	59
7	75	59	77	87	59	85	–	81	63	53
8	74	60	72	86	62	79	81	–	68	63
9	64	49	54	67	54	62	63	68	–	61
10	63	44	59	53	48	59	53	63	61	–

The one-hour-at-a-time analysis suffers from the disadvantage that we cannot “borrow strength” from the hours adjoining that of interest. Therefore the last of our approaches uses the multivariate method of Brown et al (1994a). As a preliminary step we use the univariate approach above to cluster the hours according to their degree of similarity. We expect to see one, two or more hour-clusters that might be grouped into multivariate response vectors for subsequent analysis. That analysis is now underway.

2 Log PM₁₀ Concentrations In Vancouver.

In this section we show that an AR(1) model describes quite well the daily averages of de-trended log PM₁₀ concentrations in Vancouver. Hourly PM₁₀ measurements collected by a network of TEOM monitors across the GVRD in 1996 were used to calculate daily average values used in this analysis. For any location x and day d , let $X(x, d)$ represent the daily log PM₁₀ average concentration ($\mu g m^{-3}$). Furthermore let $S(d)$ represent the overall trend in these spatial averages for day d , *i.e.*

$$S(d) = \mu' + D_{day} + W_{week}$$

where $\mu' = \mu + (H_1 + \dots + H_{24})/24$ is the overall mean effect in the daily model with μ representing the average over all sites while H_j denotes the corresponding average for hour $j = 1, \dots, 24$ once μ has been subtracted from all responses.

To explore the nature of the temporal variation in the daily de-trended residuals $D(x, d) = X(x, d) - S(d)$ we estimated at each monitoring site x the auto-correlation in the D-series. In Figure 1 we see the resulting auto-correlation function plots. They indicate a strong first order auto-correlation at each site. The corresponding partial auto-correlation function plots (in Figure 2) confirm this observation. The latter in particular suggests we can rule out a moving average component in the series. The consistency across monitoring sites seen in the figures suggests the adoption of a single time-series model applicable for all sites.

Our analysis thus led us to fit to the de-trended residuals, a single first order autoregressive model

$$D(x, d) = \beta D(x, d - 1) + d(x, d),$$

with estimated coefficient $\hat{\beta} = 0.34$.

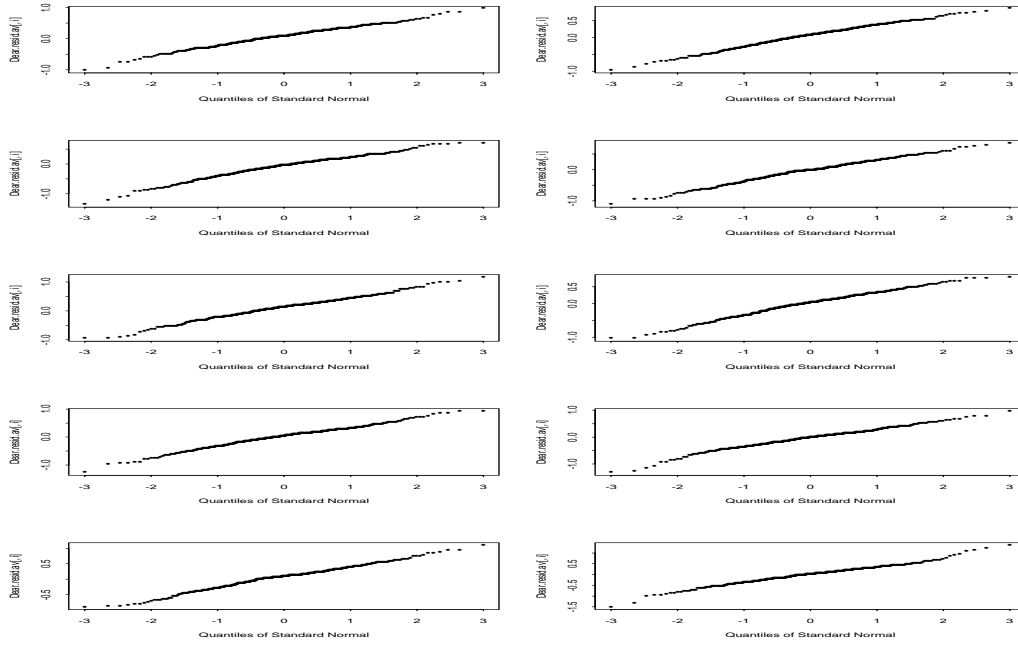


Figure 5: QQ-norm Plots For the AR(1) Residuals At All Monitoring Sites.

Table 4: Cross-correlation ($\times 100$) of AR(1) Residuals for De-trended Daily Average log PM₁₀ Residuals Between Different Monitoring Stations.

site	1	2	3	4	5	6	7	8	9	10
1	–	66	82	76	60	78	74	74	62	61
2	66	–	64	65	79	71	64	63	56	51
3	82	64	–	78	54	84	77	74	55	59
4	76	65	78	–	64	86	88	87	69	59
5	60	79	54	64	–	64	62	62	57	51
6	78	71	84	86	64	–	86	79	63	62
7	74	64	77	88	62	86	–	81	64	57
8	74	63	74	87	62	79	81	–	69	68
9	62	56	55	69	57	63	64	69	–	63
10	61	51	59	59	51	62	57	68	63	–

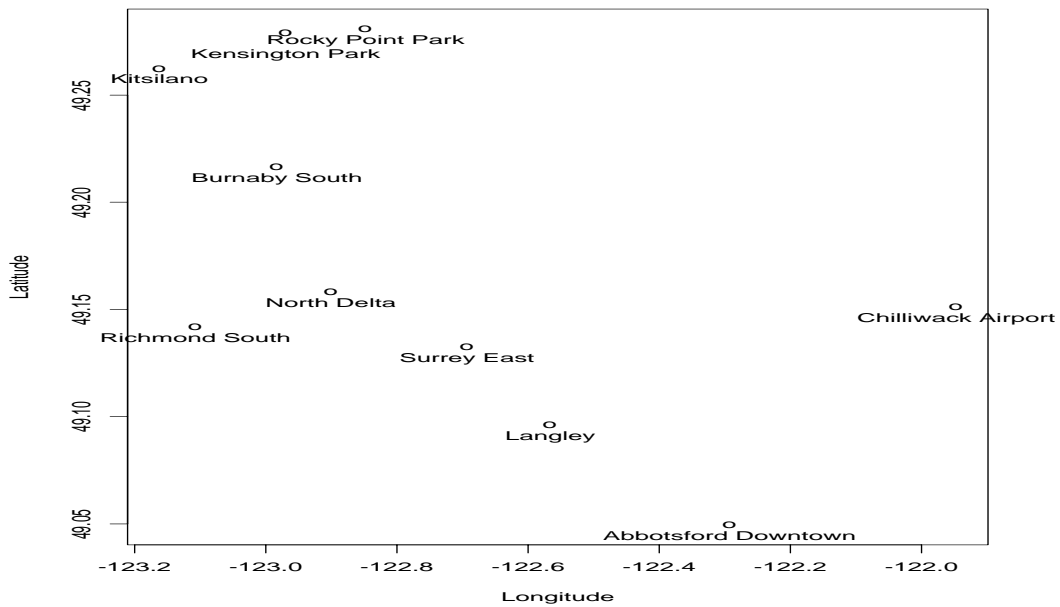


Figure 6: Locations of the 10 PM₁₀ Monitoring Stations in Vancouver.

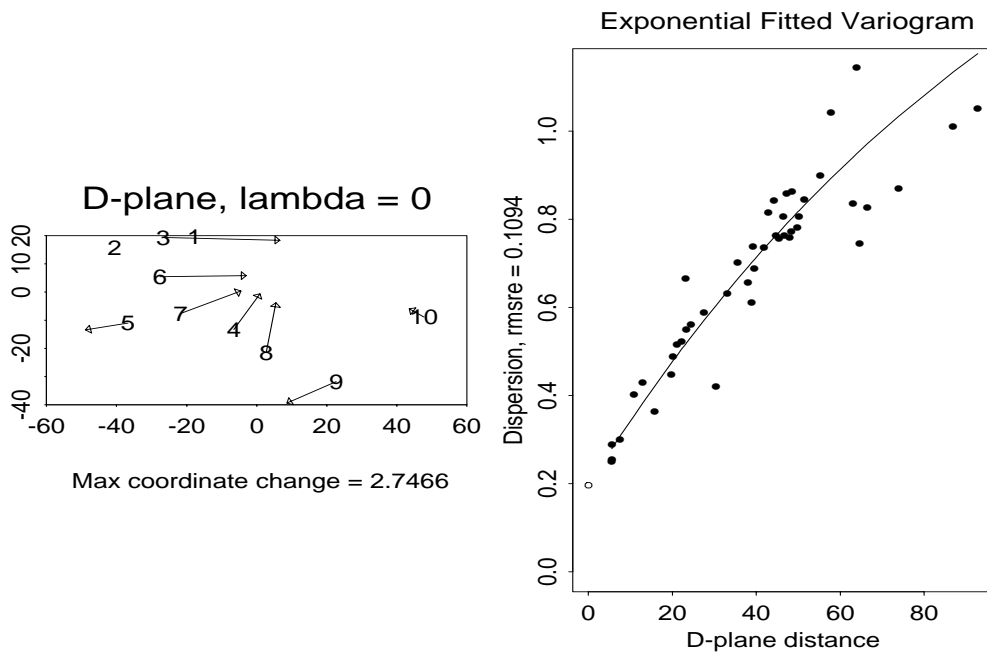


Figure 7: The Variogram Fit With No Smoothing.

To check the quality of this fit we estimated the auto-correlations in the resulting (“de-AR’d”) residual series. We present in Figures 3-4 the counterparts for these residuals, of Figures 1-2. They show little or no remaining auto-correlation.

The spatial prediction methodology used in the next section assumes a Gaussian joint distribution for the response field being predicted. The veracity of that assumption was examined through the use of QQ-norm plots, a standard way of diagnosing departures from a normal data distribution. In such plots sample quantiles are compared against those for a normal distribution, the assumption of normality being supported when these plots turn out to be linear. Our assumption is well supported by the strikingly linear QQ-norm plots in Figure 5 for these residuals at all monitoring sites.

We now face the issue of what we called spatial correlation “leakage” in the Introduction. Does spatial correlation in the AR(1) residual series for daily values leak in the same way as that noted in the Introduction for hourly series? To answer that question we first exhibit Table 3, the spatial cross-correlations between de-trended residuals at different sites for log-transformed daily averages.

We see in Table 4 the corresponding correlations for AR(1) residuals. It shows the answer to the question proves to be negative. Little or no decrease in spatial cross- correlation is seen on comparing the entries in Tables 3 and 4. This allows us to develop in the next section a spatial prediction methodology for use with the log PM10 data collected by the ten GVRD monitoring stations.

3 Interpolating the Daily Field.

In this section we develop a predictive distribution for unmeasured daily log PM10 concentration averages given data from 10 monitoring sites in Vancouver (more properly the GVRD). The locations of these sites may be seen in Figure 6. We will show how the mean of that distribution may be used to generate imputed values for the unmeasured pollution levels at each of about 300 census tracts. Its covariance in turn provides an indication of the (joint) reliability of those imputed values. Indeed that covariance can be used to generate 95% (or other level) credible sets for the unmeasured values.

We now describe the basic elements of the spatial predictor used in this analysis. Technical details of the method will be presented in Le, Sun and Zidek (1999).

The method was proposed by Le and Zidek (1992) for univariate random fields like the subject of this paper. It was extended to multivariate random fields by Brown, Le and Zidek (1994a) and further by Le, Sun and Zidek (1997) to enable the theory to contend with situations where data are systematically missing-by-design.

The method uses a hierarchical Bayesian approach. At level one the random field, for example the de-AR’d residual field in this paper is assumed to have a joint Gaussian distribution conditional on a model for the mean surface (and certain covariates in that model). As well at this level the spatial covariance matrix Σ is assumed to be known and fixed. Σ would be of dimension 299 in the application of this paper since in addition to the 10 monitoring sites an additional 299 unmonitored sites, centroids of census tracts representing population receptor sites for exposure assessment, are specified. Daily values for these other sites, are to be predicted.

A conjugate prior distribution is postulated at level 2 to account for uncertainty about Σ and other first level parameters. That prior for Σ is the Generalized Wishart distribution proposed by Brown, Le and Zidek (1994b). It requires the specification of its hyper-covariance matrix Ψ of the same dimension as Σ . With it specified the resulting predictive distribution becomes not Gaussian but a Multivariate -t (or matic-t in the multivariate case).

Specifying Ψ the hyper-covariance matrix corresponding to to the spatial cross-correlation of the de-AR’s residuals in this paper is a challenging problem. It sub-diagonal matrix corresponding to the monitoring stations can be estimated by type II maximum likelihood estimation using the EM

algorithm. However that sub-diagonal estimate Ψ_{gg} must be extrapolated to an estimate for the whole of Ψ .

That is done by using the method of Sampson and Guttorp (1992, hereafter SG). As described by Le, Sun and Zidek (1999) that method adopts a hypothetical Euclidean “D-plane” with respect to which the co-ordinates of the $\{\Psi_{ij}\}$ are isotropic. That is, the covariation between sites is a monotone decreasing function say ζ of their D-plane distances. The method estimates that monotone function and the D-plane location co-ordinates associated with each of the monitored sites $\mathbf{d}_i = (d_{i1}, d_{i2})$ for site i with geographical co-ordinates $\mathbf{g}_i = (g_{i1}, g_{i2})$. The method relates the $\{\mathbf{d}_i\}$ to the $\{\mathbf{g}_i\}$ through thin plate smoothing splines f by means of $d_j = f_j(\mathbf{g})$, $j = 1, 2$. These splines are fitted to the D- and G-plane co-ordinate pairs for the gauged or measured sites, the degree of fit depending on the so-called smoothing parameter λ or “lambda” in the figures below. The $\{\mathbf{d}_i\}$ are then replaced by the fits $\{\mathbf{f}(\mathbf{g}_i)\}$, where $\mathbf{f} = (f_1, f_2)$.

Large values of that parameter will entail poor G- to D-plane co-ordinate fits. However those splines will more faithfully maintain the character of the G-plane and lead to simplicity of interpretation of the results of the analysis. At the other extreme, small values can lead to splines that twist the G-plane into unrecognizable form while ensuring a good fit to the estimated D-plane co-ordinates.

The choice of this parameter is subjective. “Small” tends to be better because the co-ordinates of the estimated Ψ_{gg} will tend to be more closely isotropic in the \mathbf{f} image of the G-plane. On the other hand some smoothing is desirable to achieve a degree of interpretability in the relationship between the resulting plane and its G counterpart.

Once \mathbf{f} has been specified, the required extension of Ψ_{gg} to Ψ can easily be made. Represent the G-plane co-ordinates of sites i and j corresponding to Ψ_{ij} , \mathbf{g}_i and \mathbf{g}_j , by their \mathbf{f} images in the D-plane $\mathbf{d}_i = \mathbf{f}(\mathbf{g}_i)$ and $\mathbf{d}_j = \mathbf{f}(\mathbf{g}_j)$. Finally estimate Ψ_{ij} by $\zeta(\|\mathbf{d}_i - \mathbf{d}_j\|)$.

In practice the SG method is implemented through the so-called “variogram” in exactly the same way as described above for the covariance. In general for a random field $Z(x)$ the latter is defined for locations x and x' by $Var[Z(x) - Z(x')]$. It is closely related to the covariance and like the latter is easily estimated when independent replicates of Z over time are available at the two sites. The estimate is simply the sample average of squared differences in Z between the two sites. Software has been developed for implementing the SG method and we are indebted to Professors Guttorp and Sampson for supplying the version used here. Estimates of Ψ can easily be found from estimates of the variogram.

For any pair of the 10 monitored sites in our application we estimate the variogram in the manner described above. The resulting 45 estimates (for all possible site pairs) are assigned D-plane co-ordinates in the manner described above. Each of the 45 estimates can be plotted against the D-plane distance separating them. Of course the D-plane distance between them will depend on the selected size of the spline smoothing parameter. However, if no smoothing is used one can see that plot in the right hand panel of Figure 7. The scatter plot shows 45 plotted variogram estimates and the best fitting variogram plotted against them.

The left-hand panel of that figure shows two monitoring sites # 5 (Richmond) and # 9 (Abbottsford) that must move away from the remaining 8 sites to achieve an isotropic correlation field. In other words, these two stations tend to be un-correlated with the rest so must be moved away in the D-plane to achieve inter-station distances commensurate with the low spatial cross-correlations they have with the rest.

Figure 8: Transformation to the D-plane With No Smoothing.

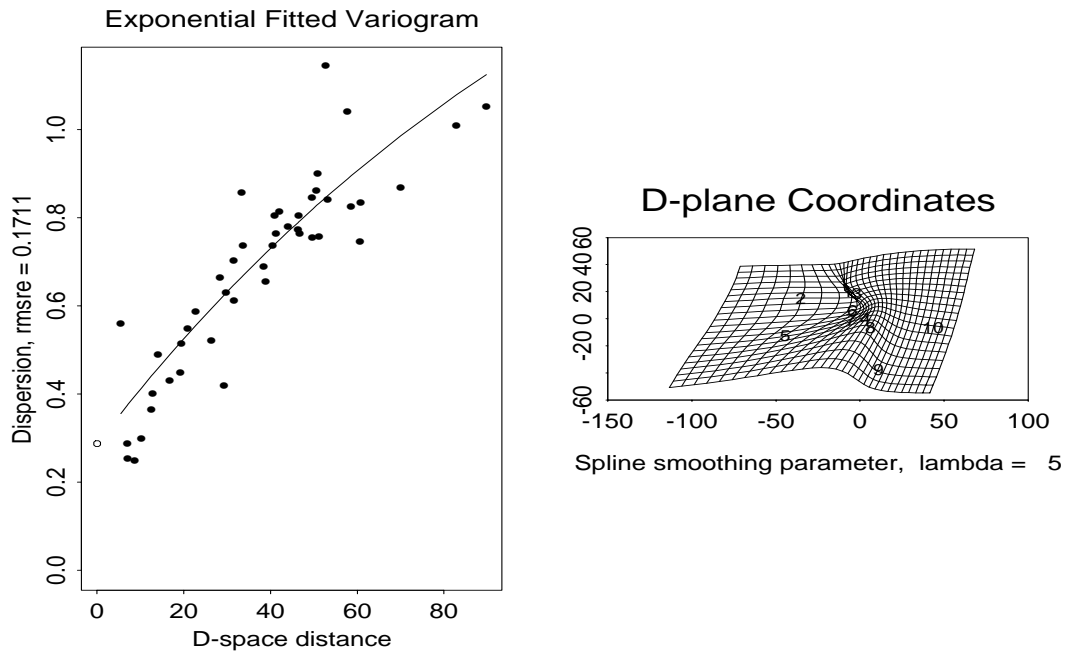
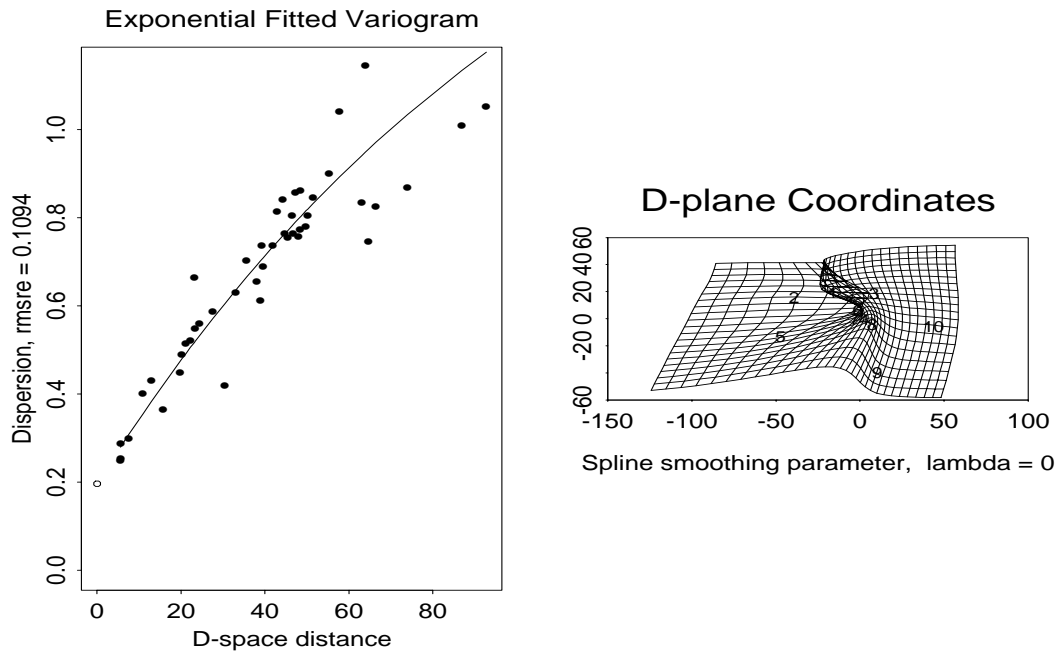


Figure 9: Transformation to D-plane with Moderate Smoothing.

Table 5: **Location of PM₁₀ Monitoring Sites**

Site	Location	Latitude	Longitude
1	Rocky Point Park	49.28083	122.8481
2	Kitsilano	49.26250	123.1625
3	Kensington Park	49.27917	122.9697
4	Surrey East	49.13278	122.6933
5	Richmond South	49.14194	123.1078
6	Burnaby South	49.21667	122.9833
7	North Delta	49.15833	122.9008
8	Langley	49.09611	122.5664
9	Abbotsford Downtown	49.04944	122.2925
10	Chilliwack Airport	49.15111	121.9469

Figure 8 offers a different view of the same situation. The scatterplot in that figure is identical to that in Figure 7. However it shows in a more pictorial way how the (geographic) G-plane must be folded to re-organize the G-surface so as to make the variogram separation between sites correspond to their D-plane distances.

That picture shows the surface must essentially be folded over on itself to achieve the desired state of isotropy. Interpreting the result is hard because of that folding. So an alternative is offered in Figure 9. There with a small amount of smoothing a flatter surface results with slight loss in the quality of the “fit”.

Figure 10 offers quite a different diagnostic. The bi-orthogonal grid depicted there shows again how the G-surface must be change to make the inter-station correlations come into line with the inter-station distances. This diagnostic will prove particularly valuable in the sequel to this report where the grids will be used to assess the hour-to-hour changes in the wind-fields that help to determine the spatial distribution of pollution fields. The solid curves in Figure 10 show the directions in which the surface must be contracted to achieve isotropic spatial correlation fields. In contrast expansion along the dotted lines is called for.

In any event, with the latter choice of the smoothing parameter we can apply the spatial prediction methodology described above. Figures 11-14 show the results of doing so for selected days. In particular, Figure 11 shows the variation in the imputed PM₁₀ field over days in two successive summer weeks (#31 and #32), Day # 7 from the former followed by Days # 1,2 and 3 in the latter. Figure 12 continues this sequence for days # 4, 5, 6 and 7.

The remaining Figures 13-14 show the counterparts for winter days of Figures 11-12.

Notice the substantial temporal and spatial variation in the interpolated fields in the Figures 11-14. In particular, the interpolated daily PM₁₀ surface is not flat on any given day.

4 Validation Study.

Since interpolated spatial fields can serve important societal purposes, the accuracy of an interpolation method must be assessable and high. Moreover, its degree of inaccuracy must itself be accurately assessable.

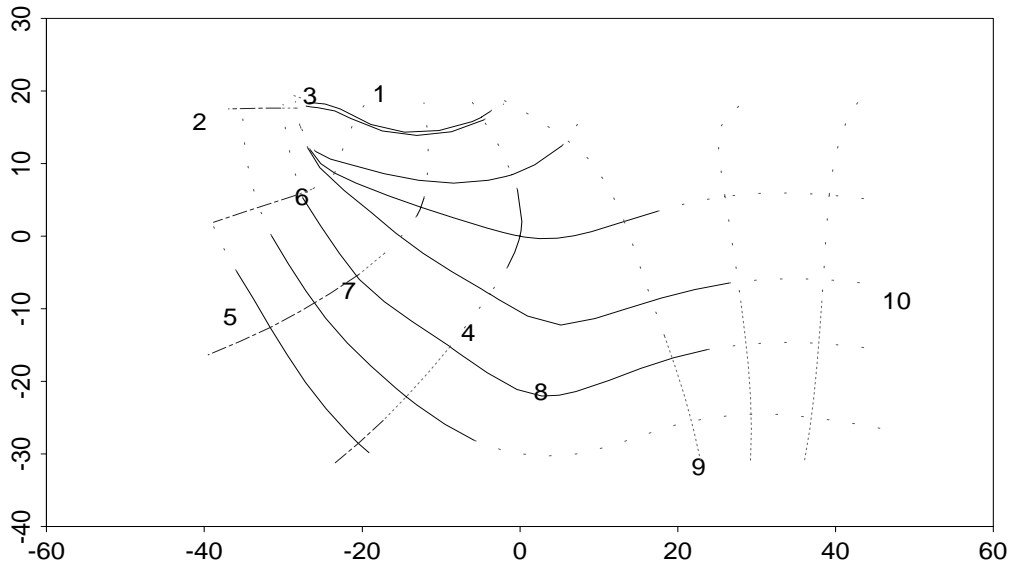


Figure 10: The Biorthogonal Plot After Moderate Smoothing.

Table 6: Coverage Probabilities and RMSPE

site	1	2	3	4	5	6	7	8	9	10
Coverage	94.8%	96.2%	94.0%	95.1%	95.4%	94.5%	94.5%	95.1%	94.0%	93.2%
RMSPE	0.161	0.174	0.161	0.123	0.201	0.133	0.145	0.150	0.224	0.258

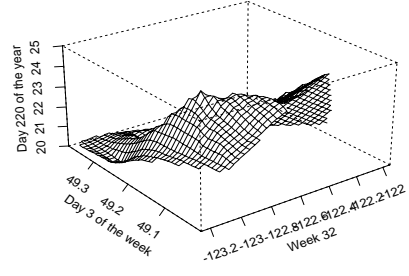
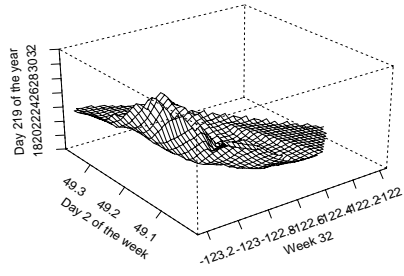
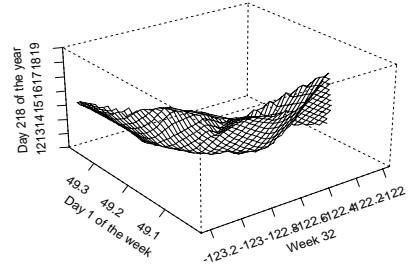
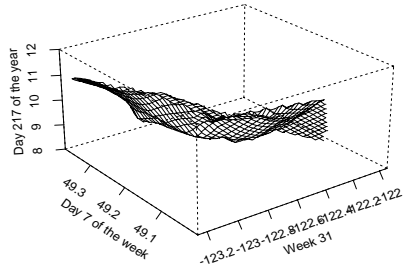


Figure 11: Interpolated PM10 Field For Selected Summer Days in 1996.

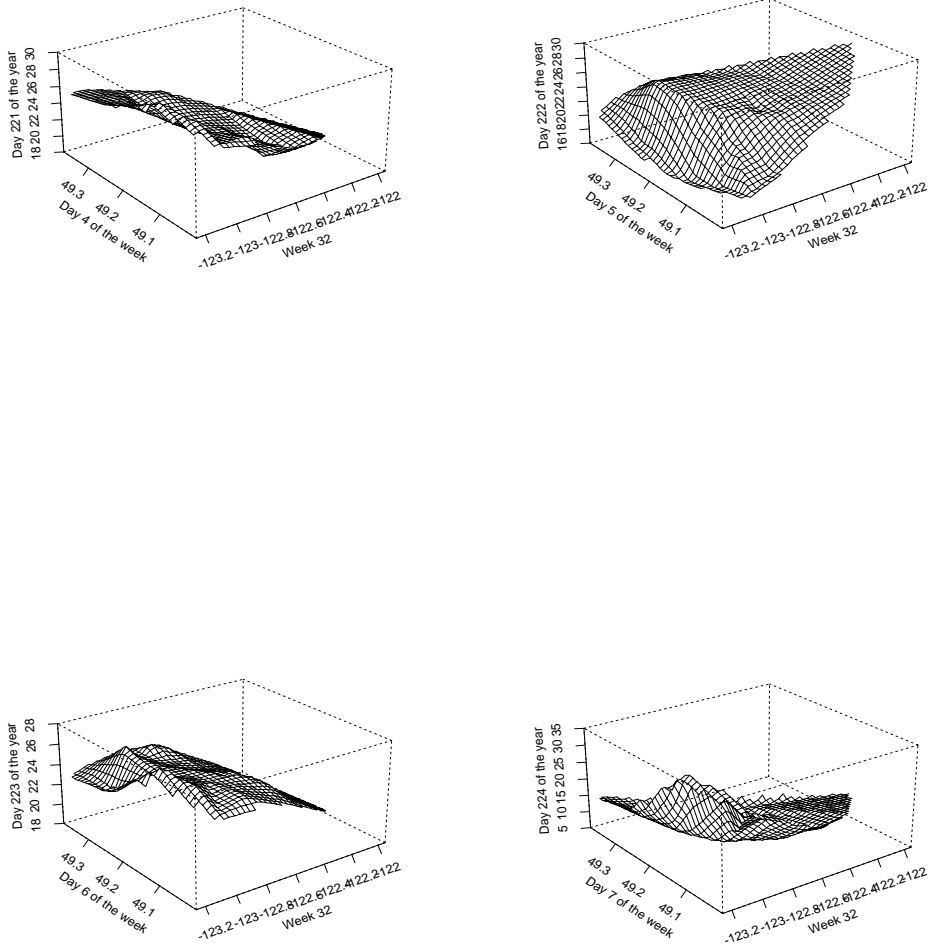


Figure 12: Interpolated PM10 Field For Selected Summer Days: Continuation of Figure 11.

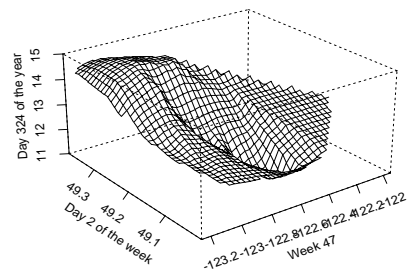
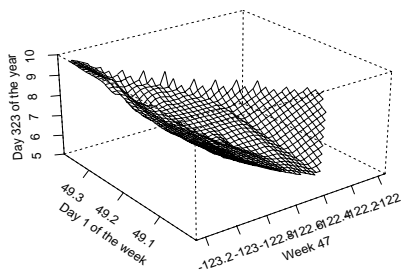
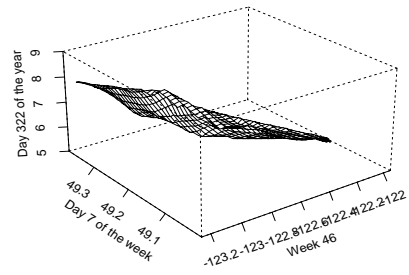
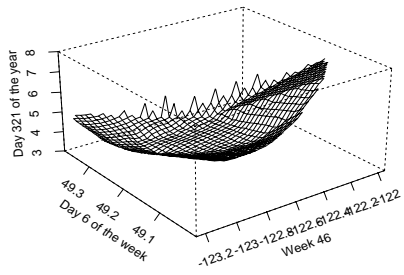


Figure 13: Interpolated PM10 Field For Selected Winter Days in 1996.

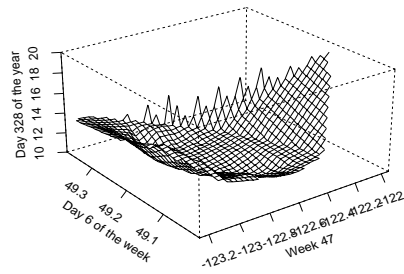
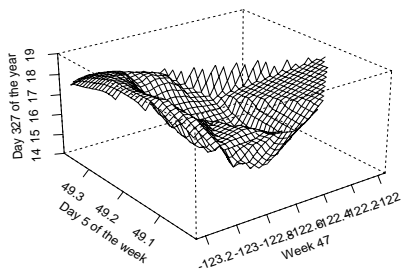
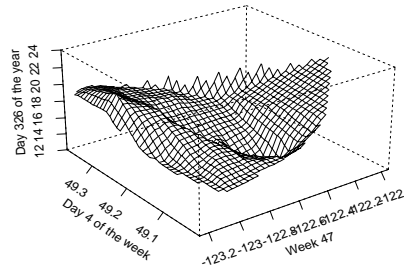
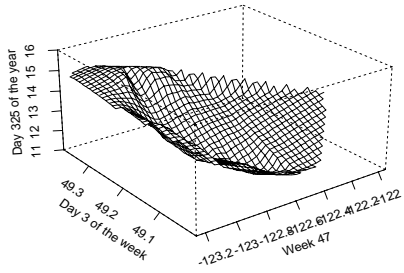


Figure 14: Interpolated PM10 Field For Selected Winter Days: Continuation of Figure 13.

Table 7: **Cross-correlation ($\times 100$) of De-trended Residuals Before Removal of Their AR(1) components: Daytime Series**

Site	1	2	3	4	5	6	7	8	9	10
1	–	50	80	71	47	78	71	65	65	58
2	50	–	51	47	71	55	46	46	52	40
3	80	51	–	75	44	85	73	68	53	56
4	71	47	75	–	49	80	84	86	64	54
5	47	71	44	49	–	48	48	49	48	38
6	78	55	85	80	48	–	82	71	62	59
7	71	46	73	84	48	82	–	75	60	51
8	65	46	68	86	49	71	75	–	63	62
9	65	52	53	64	48	62	60	63	–	62
10	58	40	56	54	38	59	51	62	62	–

We now present the results of a preliminary empirical validation study of our interpolation methodology. We assess accuracy and the accuracy of our assessments of inaccuracy. A more extensive study is currently underway.

Table 5 lists the 10 stations in our study. Note for future reference that the “Chilliwack Airport” station is well separated from the other 9 (see also Figure 6). In contrast, “Burnaby South” lies well within their geographical envelope.

In our validation study we fitted and fixed once and for all the 10×10 spatial covariance matrix of the de-AR’s residuals. [We comment on this fact in the next section.] We then commenced to remove the stations from the network one at a time. After a single station was removed we used the remaining 9 to develop the spatial predictive distribution for the missing site. The predicted values at the missing station could then be compared with the actual values, day-by-day throughout the year.

Table 6 shows the root mean square prediction error (RMSPE in $\log \mu m^{-3}$) of the predictive distribution’s mean over all the days of the year. Not surprising the de-AR’d residuals prove most difficult to predict at the two most remote stations, “Chilliwack Airport” and “Abbotsford Downtown.” While not geographically remote “Richmond South” values also prove difficult to predict. Though close geographically to the other stations it remains a outlier as our analysis of spatial correlation patterns in the last section has shown. Its values are not well-correlated with the remaining stations, making it hard to predict.

Figures 15-16 show the interpolated values themselves compared with the actual raw data values (μm^{-3}). Notice the difficulty encountered in predicting extremely large and small values in both cases. Although stations prove difficult to predict, their 95% prediction credibility intervals seem to correctly reflect that difficulty. Table 6 shows the actual coverage frequency of the 95% prediction credibility intervals are close to their nominal 95% level. “Chilliwack Airport’s actual level of 93.3% deviates most among the 10 stations from 95%.

We look more closely at results for “Burnaby South” (a proximate station) and “Chilliwack Station” (a remote station that challenges the prediction methodology). Their 95% error bands and actual de-AR’d log data series are plotted in Figures 17 and 18, respectively for just the particularly important summer months (so that details in the plots are more visible). By and large, the error bands do “bracket” the actual results in both cases.

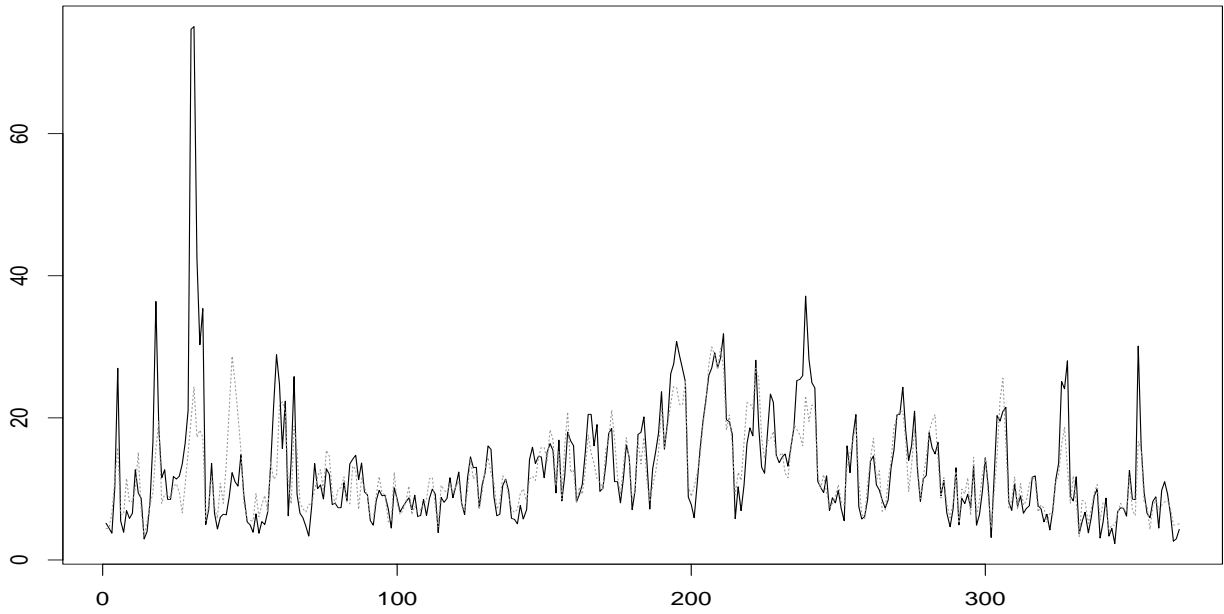


Figure 15: Estimated means and the true values (Chilliwack Airport).

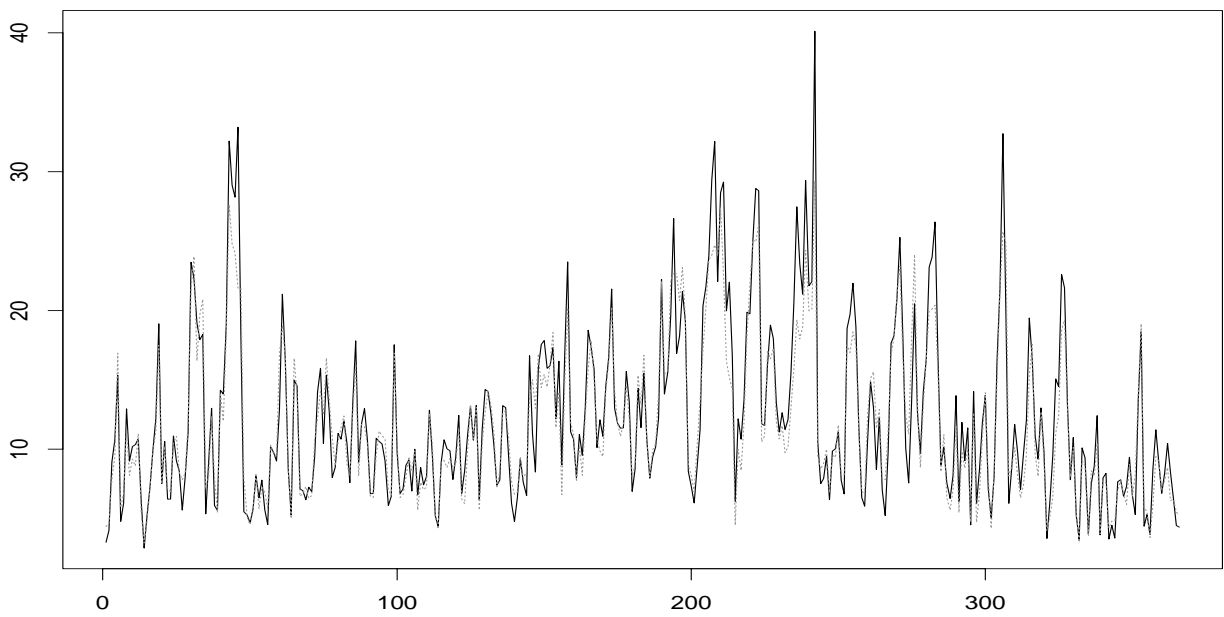


Figure 16: Estimated means and the true values (Burnaby South).

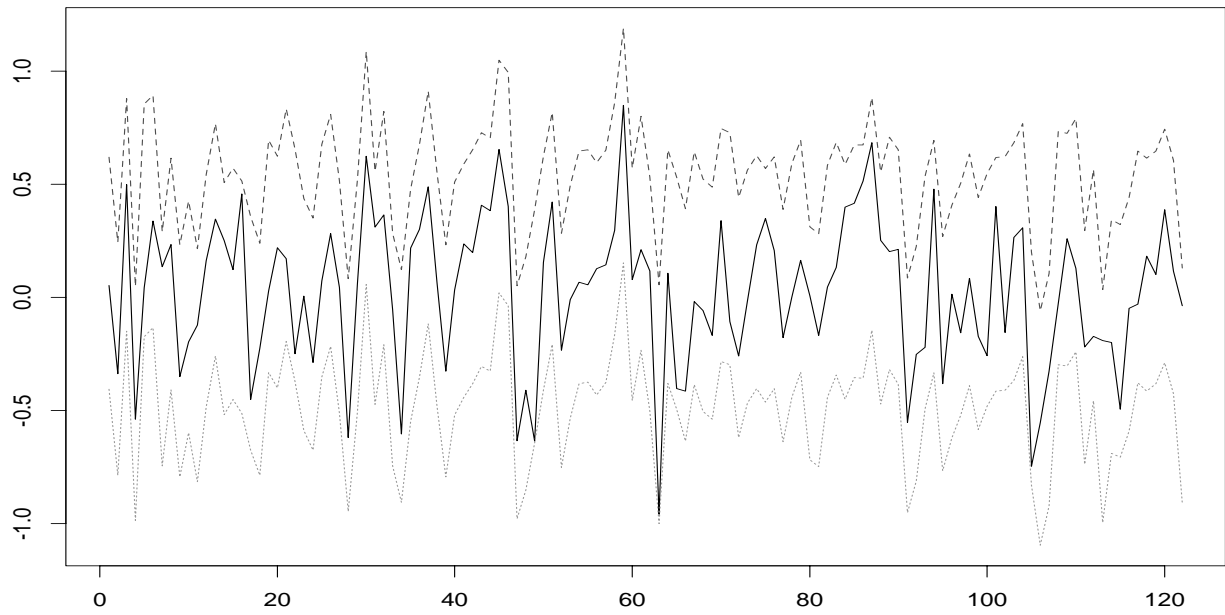


Figure 17: Validation of Chilliwack Airport - 95% confidence band and the true values of the de-AR'd residuals, June - September 1996.

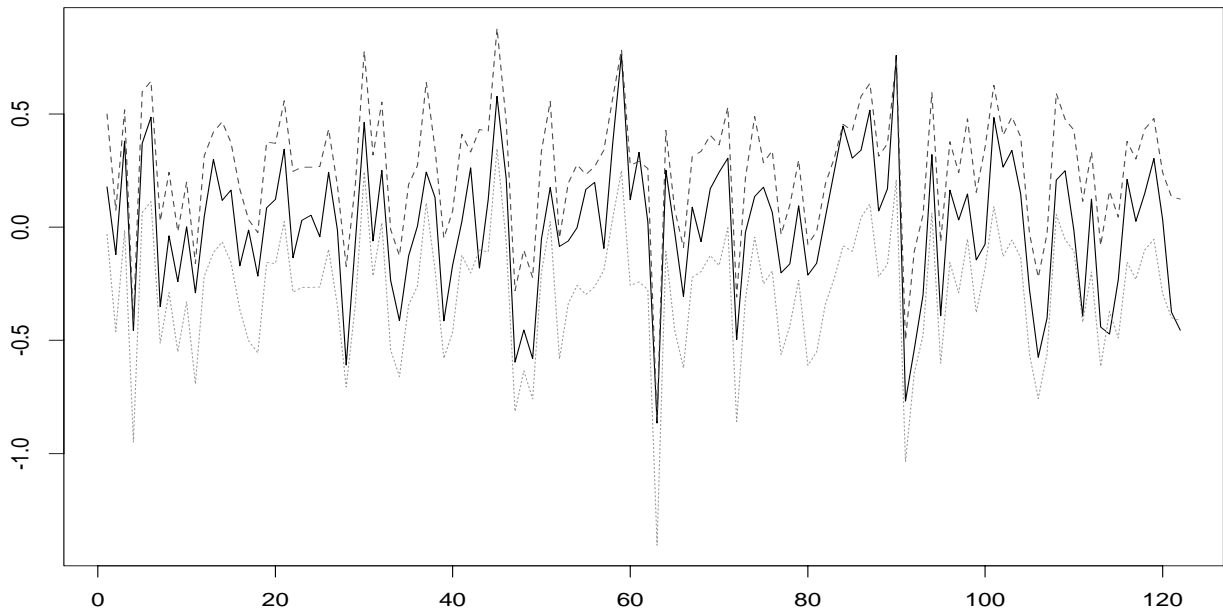


Figure 18: Validation of Burnaby South - 95% confidence band and the true values of the de-AR'd residuals, June - September 1996.

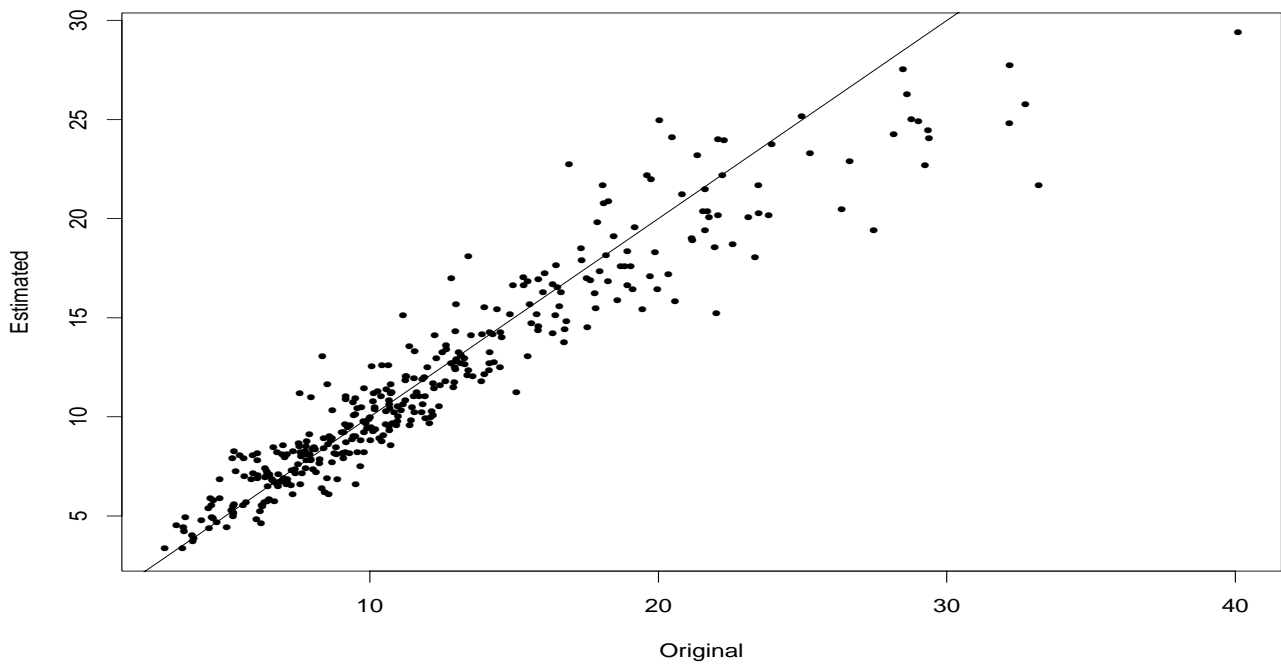


Figure 19: Scatter plot of the original against the estimates (Burnaby South).

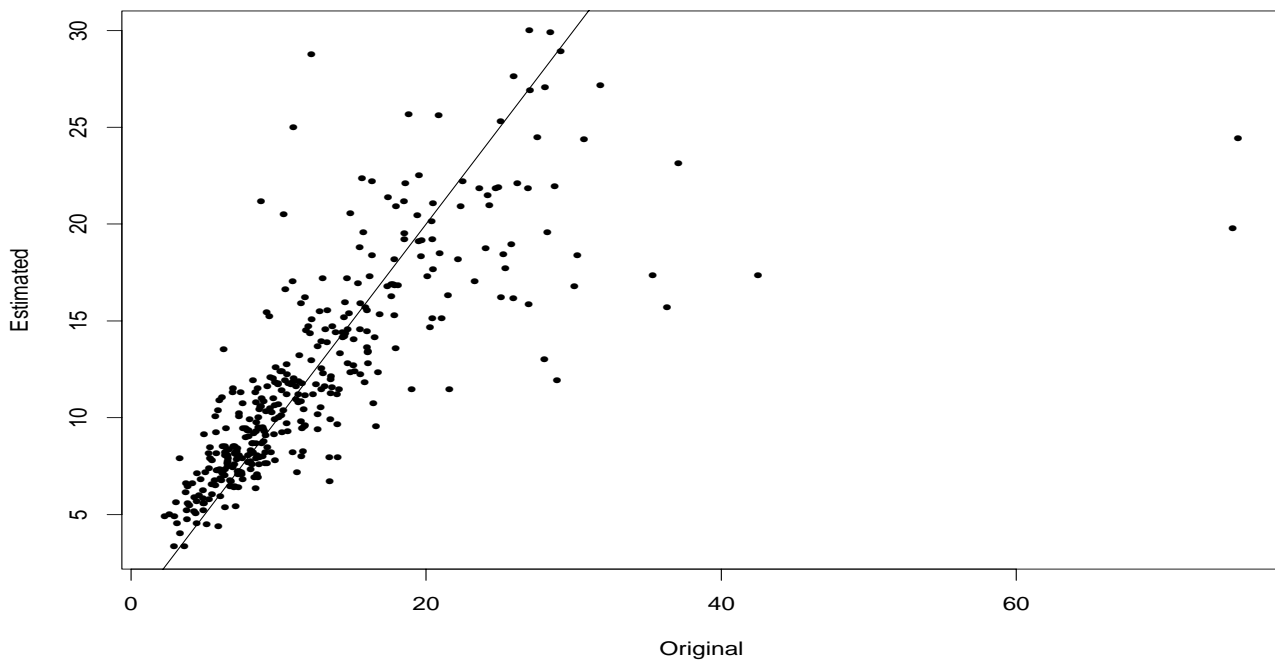


Figure 20: Scatter plot of the original against the estimates (Chilliwack Airport).

Table 8: **Cross-correlation ($\times 100$) of Detrended Residuals After Removal of Their AR(1) Components: Daytime Series**

Site	1	2	3	4	5	6	7	8	9	10
1	–	53	79	70	48	75	69	65	63	54
2	53	–	56	52	71	59	50	49	55	44
3	79	56	–	75	48	85	74	68	53	54
4	70	52	75	–	52	80	85	85	66	57
5	48	71	48	52	–	52	51	51	51	41
6	75	59	85	80	52	–	82	71	62	58
7	69	50	74	85	51	82	–	75	62	53
8	65	49	68	85	51	71	75	–	65	63
9	63	55	53	66	51	62	62	65	–	63
10	54	44	54	57	41	58	53	63	63	–

Table 9: **Cross-correlation ($\times 100$) of Detrended Residuals Before Removal of Their AR(1) Components: Nighttime Series**

Site	1	2	3	4	5	6	7	8	9	10
1	–	71	78	71	60	78	73	72	61	60
2	71	–	63	67	74	80	74	66	52	48
3	78	63	–	67	45	82	74	62	49	55
4	71	67	67	–	59	81	83	83	66	46
5	60	74	45	59	–	62	63	63	60	46
6	78	80	82	81	62	–	87	75	60	51
7	73	74	74	83	63	87	–	80	65	49
8	72	66	62	83	63	75	80	–	70	56
9	61	52	49	66	60	60	65	70	–	55
10	60	48	55	46	46	51	49	56	55	–

Figures 19-20 show this difficulty in another way. While the predictor is able to predict well actual values when they are small or moderate, it under-predicts large value demonstrating the regression effect anew. [We comment on this issue in the next section.]

5 Discussion.

The results of Section 3 demonstrates the use of a spatial predictive distribution for airborne particulate pollution fields. The validity of the assumptions leading to that distribution is demonstrated in Section 2. Section 4 suggests the method works reasonably well. It predicts the actual measurement field well unless unless the unmonitored stations are well separated from the primary domain of the data. Moreover it seems to quantify correctly its own level of prediction error.

However, we emphasize that we did not extrapolate the spatial covariance matrix from that of 9 stations to the 10th each time we constructed that distribution in Section 4. A proper cross validation study (the subject of current investigation) would require that. So the predictive distribution may not perform as well as our preliminary findings suggest. Nevertheless, we are fairly confident that out

Table 10: **Cross-correlation ($\times 100$) of De-trended Residuals After Removal of Their AR(1) components: Nighttime Series**

Site	1	2	3	4	5	6	7	8	9	10
1	-	74	79	72	60	79	74	72	61	60
2	74	-	68	70	74	82	76	69	56	52
3	79	68	-	70	47	84	76	65	52	56
4	72	70	70	-	60	82	84	83	68	52
5	60	74	47	60	-	62	62	63	62	45
6	79	82	84	82	62	-	88	75	62	55
7	74	76	76	84	62	88	-	80	66	53
8	72	69	65	83	63	75	80	-	72	60
9	61	56	52	68	62	62	66	72	-	57
10	60	52	56	52	45	55	53	60	57	-

Table 11: **Lag One Cross-correlation Leakage of Hourly Data - Spatial Cross-correlation ($\times 100$) Between deAR'd and Detrended Residuals**

site	1	2	3	4	5	6	7	8	9	10
1	-3	10	15	13	10	15	14	10	9	8
2	10	1	9	12	16	11	11	11	10	8
3	11	11	1	13	10	17	13	8	7	7
4	10	12	15	2	12	16	19	14	11	9
5	9	12	7	10	4	7	10	12	9	9
6	9	12	13	13	13	-1	14	11	9	8
7	10	12	13	15	9	14	3	11	8	8
8	10	9	14	13	8	13	15	-5	8	9
9	11	10	11	13	9	12	12	12	-5	11
10	16	9	15	14	9	14	13	15	13	5

findings do not exaggerate performance quality unduly. The cross-validators studies of Sun (1998) and Sun et al (1998) in other contexts generally agree with those above.

Figures 7-9 (Section 3) demonstrate how unrealistic the assumption of spatial non-stationarity can be. Clearly Vancouver’s 10 PM10 stations would have to move around a lot to get their inter-station geographic distances to correspond to their inter-station covariances. In our experience with environmental fields such non-stationary would be the rule rather than the exception since the monitors are designed mostly for compliance monitoring or ambient exposure assessment purposes.

Note that the interpolated surfaces in Figures 11-14 are not flat. Their irregularity comes in the first instance from variation in the daily levels of PM10 at the 10 monitored stations; the interpolated values must approximate the actual values at monitored stations, the “nugget effect” being quite small. However between stations the interpolated surface must regress towards the mean. The inevitable “regression-toward-the-mean” effect thus contributes to the impression of irregularity of the ambient particulate pollution field.

This finding shows that this interpolator under-predicts the extreme values in the pollution field. This could be quite significant in the analysis of population exposures and human health effects, for example. Here the contrast in pollution levels between geographical sub-regions should be preserved

Table 12: **Lag One Daily Data Cross-correlation ($\times 100$) Between de-AR'd and De-trended Residuals**

site	1	2	3	4	5	6	7	8	9	10
1	-1	-14	-2	-5	-11	-2	-1	-6	-5	-1
2	3	6	-2	-1	3	6	2	2	-7	-3
3	1	-16	2	-2	-13	1	1	-6	-7	1
4	4	-13	2	3	-10	2	2	-3	-8	-8
5	9	0	4	2	-3	6	4	3	-4	3
6	0	-14	-3	-6	-16	-3	-4	-7	-12	-7
7	-4	-17	-6	-8	-18	-7	-7	-10	-14	-11
8	1	-10	-2	-4	-9	-2	-2	-7	-11	-8
9	12	-9	6	2	-6	7	7	2	0	2
10	12	-9	6	-4	-7	4	2	-3	-1	10

to maximize the power of the method to detect association between air pollution exposures and health outcomes, such as admission to hospitals for respiratory morbidity. However because we have based our interpolation methodology on a spatial predictive distribution, the methodology recognizes these extremes (implicitly) and allows for our uncertainty about their size in health impact analysis.

To conclude we consider one other issue concerning the strategy we have developed for space-time analysis. That issue revolves around the level of temporal aggregation needed to avoid the spatial correlation leakage effect described in the Introduction. To that end we experimentally reran our analysis for 12-hour aggregates rather than 24-hour aggregates as in this paper. The result is the same: no leakage through lagged cross-correlation. We can see this by comparing Tables 7-10 for the the residual series before and after the AR(1) effect has been removed. We see in particular only a small resulting drop in the spatial correlation between stations. The result is the same whether we look at the daytime or nighttime series.

We should add a final point that much to our surprise the Sampson-Guttorp spatial covariance model for day- and nighttime series were quite similar when a moderate amount of smoothing is done. Our surprise stems from our prior belief that the big differences between day and night in the atmospheric processes would induce different levels of spatial correlation for the two periods.

In contrast we have found in current work on hourly levels of PM10 that the SG spatial covariance estimates change quite dramatically from one hour to the next particularly during the period of 12 hours following 3am.

6 Acknowledgements.

We are indebted to Drs Jianping Xue and Jack Spengler for valuable observations about results obtained with an early version of our interpolation procedure.

References

- [1] Brown, PJ, Le, ND and Zidek, JV (1994a). Multivariate spatial interpolation and exposure to air pollutants. *Canadian Journal of Statistics*, 22, 489-509.

- [2] Brown, PJ, Le, ND and Zidek, JV (1994b). "Inference for A Covariance Matrix", *Aspect of Uncertainty: A Tribute to D.V. Lindley* AFM Smith and PR Freeman (Eds). John Wiley & Sons.
- [3] Burnett, RT, Dales, RE, Raizenne, MR, Krewski, D, Summers, PW, Roberts, GR, Raad-Young, M, Dann, T, Brook, J (1994). Effects of low ambient levels of ozone and sulfates on the frequency of respiratory admissions to Ontario hospitals. *Environmental Research*, **65**, 172-94.
- [4] Li, K, Le, ND, Sun, L and Zidek, JV (1998). Spatial-temporal Models for ambient hourly PM10 in Vancouver. *Environmetrics*. To appear.
- [5] Le, ND and Zidek, JV (1992). Interpolation with Uncertain Spatial Covariance: a Bayesian Alternative to Kriging. *Journal Multivariate Analysis.*, 43, 351-374.
- [6] Le, ND, Sun, W and Zidek, JV (1997). Bayesian multivariate spatial interpolation with data missing by design, *J.R.Statist.Soc.B* 59,5-1-510.
- [7] Le, ND, Sun, W and Zidek, JV (1999). Bayesian spatial interpolation and backcasting using the Gaussian inverted Wishart model. Submitted.
- [8] Sampson, P and Guttorp, P (1992). "Nonparametric estimation of nonstationary spatial covariance structure." *J. Amer. Statist. Assoc.* Vol.87 No. 417, 108-119.
- [9] Sun, W.(1998). "Comparison of a CoKriging Method With a Bayesian Alternative". *Environmetrics*, 9, 445-457.
- [10] Sun, W, Le, ND, Zidek, JV and Burnett, R (1998) Assessment of a Bayesian multivariate spatial interpolation approach for health impact studies. *Environmetrics*, 9, 565-586.
- [11] Zidek, JV, White, R, Le, ND, Sun, W and Burnett,(1998). Imputing Unmeasured Explanatory Variables in Environmental Epidemiology With Application To Health Impact Analysis of Air Pollution. *Environmental and Ecological Statistics*, 5, 99-115.

Cross-correlation leakage

Let

$$E(x, t) = \alpha E(x, t - 1) + e(x, t)$$

be an AR(1) model, with

$$\sigma_E^2 = \text{var}[E(x, t)], \quad \sigma_e^2 = \text{var}[e(x, t)].$$

It follows that

$$(1 - \alpha^2)\sigma_E^2 = \sigma_e^2.$$

We now look at the spatial cross-correlations between x' 's:

$$\begin{aligned}
& cov[E(x, t), E(x', t)] \\
&= cov[\alpha E(x, t-1) + e(x, t), \alpha E(x', t-1) + e(x', t)] \\
&= \alpha^2 cov[E(x, t-1), E(x', t-1)] \\
&\quad + \alpha cov[E(x, t-1), e(x', t)] \\
&\quad + \alpha cov[e(x, t), E(x', t-1)] \\
&\quad + cov[e(x, t), e(x', t)].
\end{aligned}$$

Thus

$$\begin{aligned}
\sigma_E^2 cor[E(x, t), E(x', t)] &= \alpha^2 \sigma_E^2 cor[E(x, t-1), E(x', t-1)] \\
&\quad + \alpha \sigma_E \sigma_e cor[E(x, t-1), e(x', t)] \\
&\quad + \alpha \sigma_E \sigma_e cor[e(x, t), E(x', t-1)] \\
&\quad + \sigma_e^2 cor[e(x, t), e(x', t)].
\end{aligned}$$

It turns out that

$$\begin{aligned}
cor[E(x, t), E(x', t)] &= cor[e(x, t), e(x', t)] \\
&\quad + \frac{\alpha}{\sqrt{1-\alpha^2}} (cor[E(x, t-1), e(x', t)] + cor[e(x, t), E(x', t-1)]).
\end{aligned}$$

To implement these results we need the spatial cross-correlations in Table 12.

In search of spatial cross- correlation leakage we now present the compute correlations in the following table.

A Descriptive plots.

In this section we present some of the exploratory analysis that underlies the work reported in this paper. In Figure 21 we see a plot of ambient levels of untransformed hourly concentrations of PM10 at 10 Vancouver stations.

The same display for just the summer, a period of particular interest, appears in Figure 22.

The boxplots for the data in Figure 21 appear in Figure 23. These show the data distribution to be skewed at all sites, justifying the log-transformation used in our analysis. The distinctive nature of the Richmond (# 5) and Abbotsford (# 9) is revealed anew in this display.

Figures 24-26 are the counterparts of 15-17 for log-transformed daily PM10 concentrations. In particular, Figure 26 reveals the symmetric distributions achieved by the logarithmic transformation of the (geometric) daily averages.

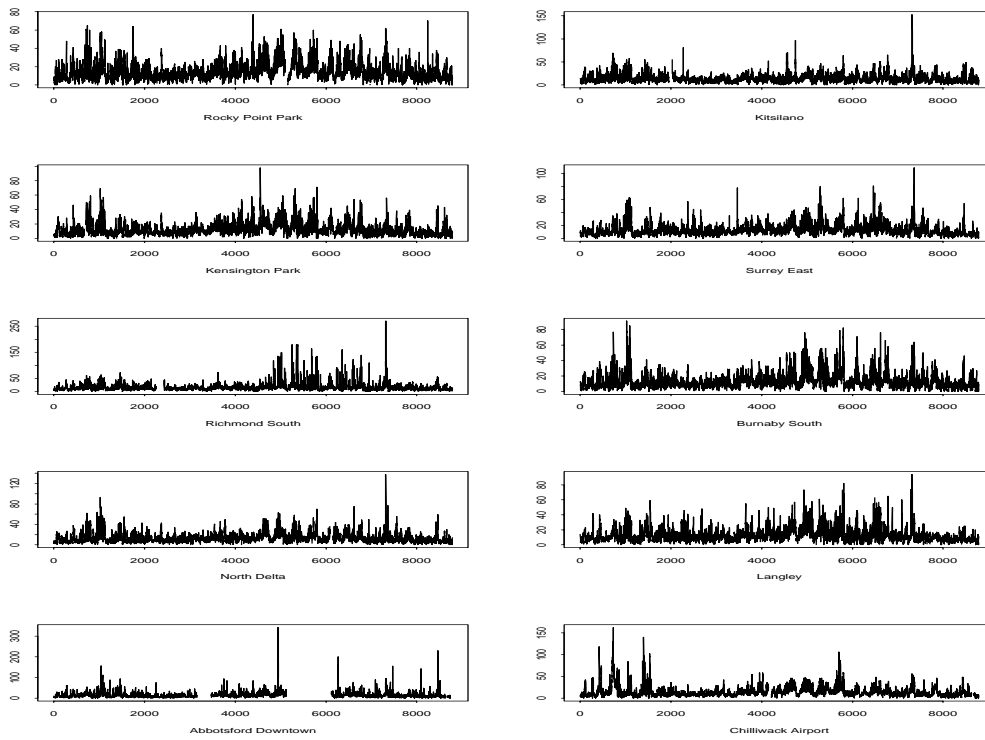


Figure 21: Ambient Hourly PM10 Levels at 10 Sites in 1996.

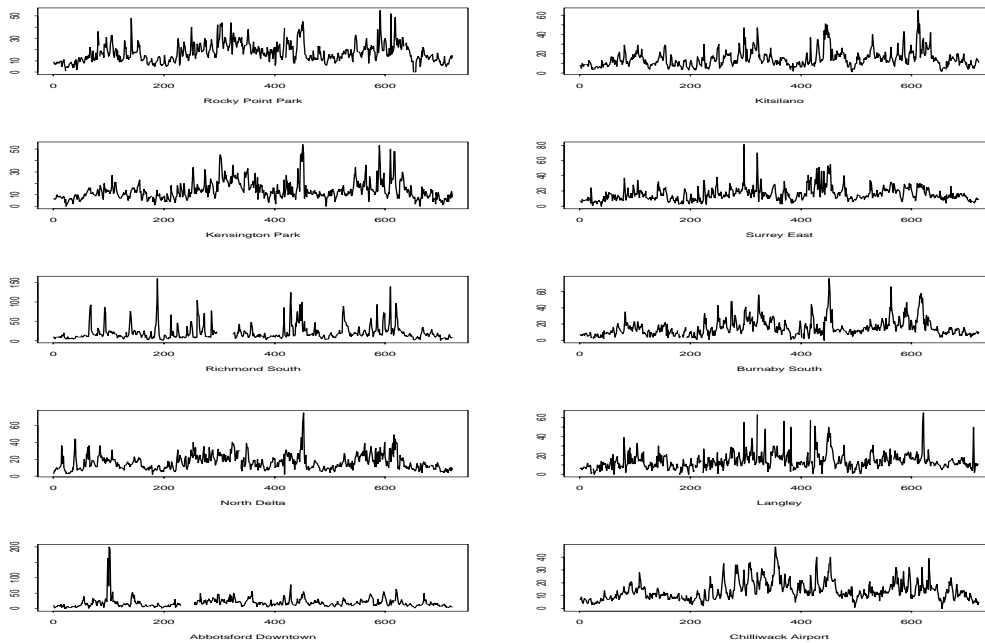


Figure 22: Ambient Hourly PM10 Levels During the Mid-August to Mid-September Period.

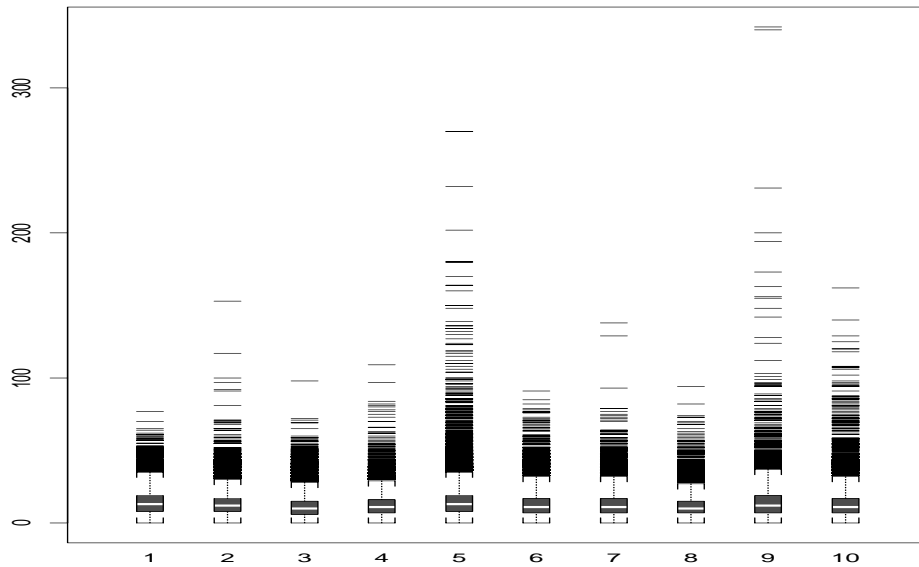


Figure 23: Boxplots of PM10 Levels in 10 Sites, 1996.

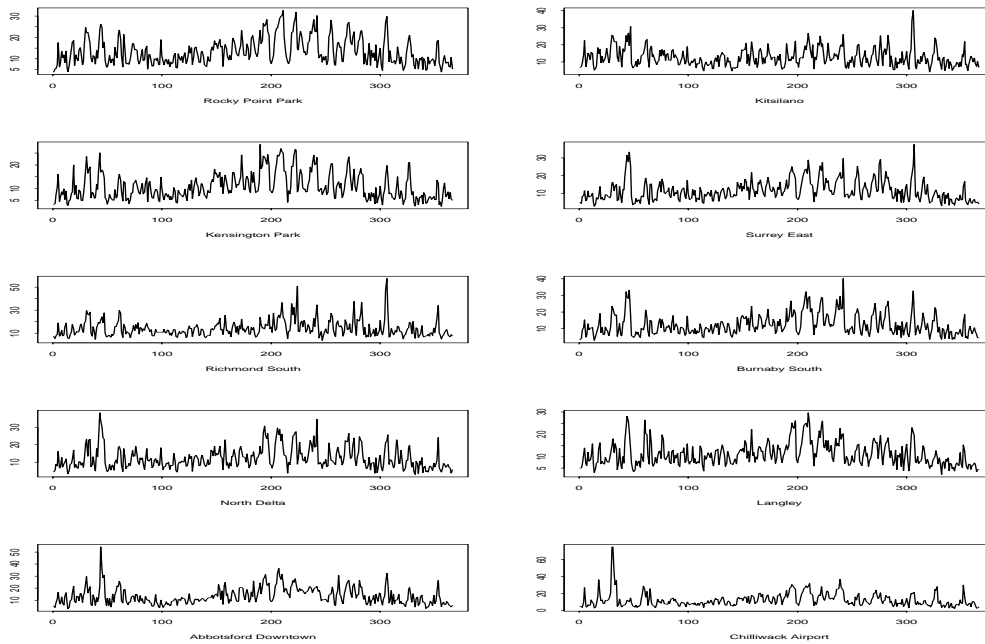


Figure 24: Daily Averages in 10 Sites, 1996.

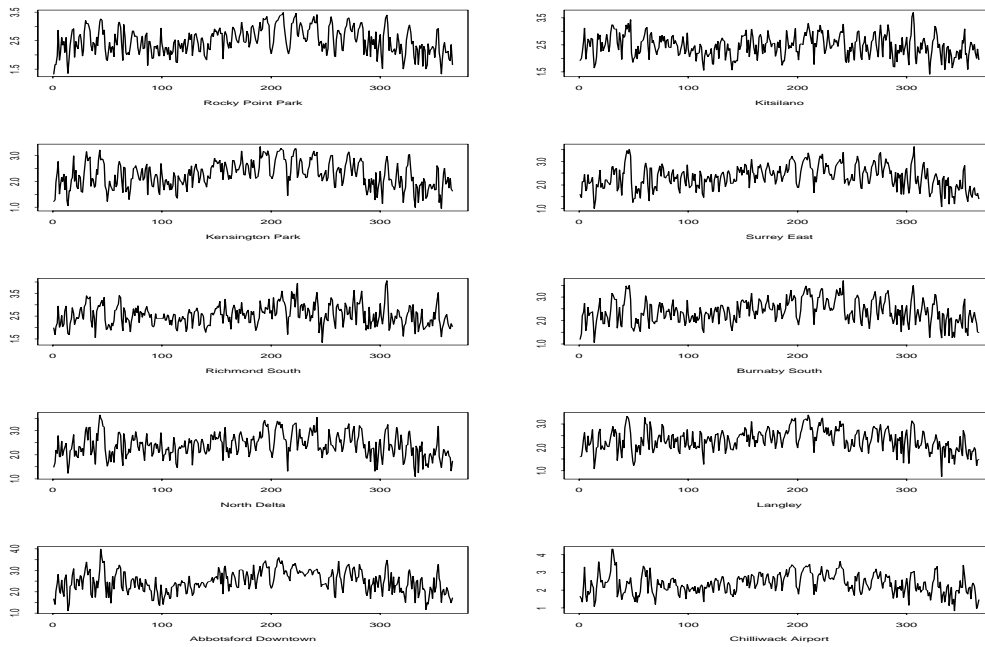


Figure 25: Log-transformation of PM10 Daily Averages in 10 Sites, 1996

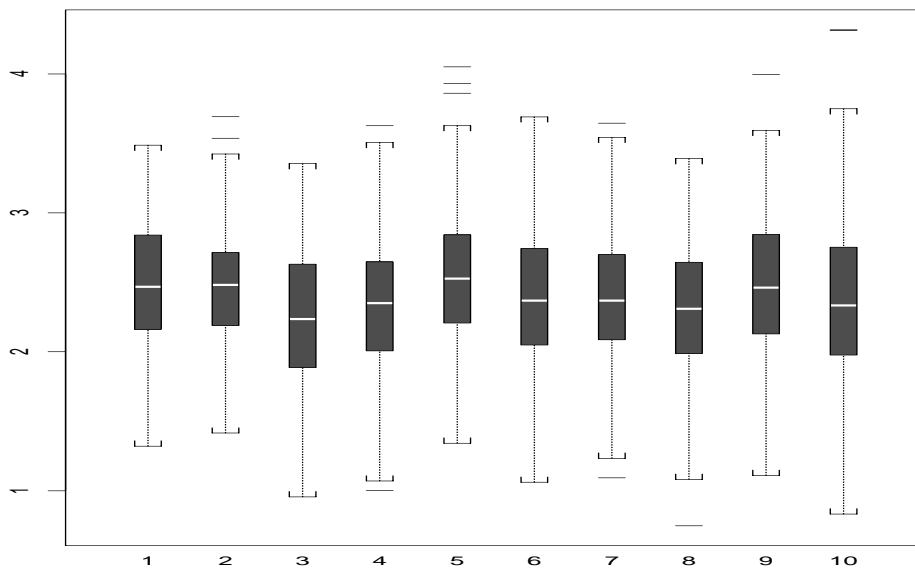


Figure 26: Boxplots of the Log-transformed Daily PM10 in 10 Sites, 1996.