Combining Spatial Statistical and Ensemble Information in Probabilistic Weather Forecasts

VERONICA J. BERROCAL, ADRIAN E. RAFTERY, AND TILMANN GNEITING

Department of Statistics, University of Washington, Seattle, Washington

(Manuscript received 27 February 2006, in final form 10 July 2006)

ABSTRACT

Forecast ensembles typically show a spread-skill relationship, but they are also often underdispersive, and therefore uncalibrated. Bayesian model averaging (BMA) is a statistical postprocessing method for forecast ensembles that generates calibrated probabilistic forecast products for weather quantities at individual sites. This paper introduces the spatial BMA technique, which combines BMA and the geostatistical output perturbation (GOP) method, and extends BMA to generate calibrated probabilistic forecasts of whole weather fields simultaneously, rather than just weather events at individual locations. At any site individually, spatial BMA reduces to the original BMA technique. The spatial BMA method provides statistical ensembles of weather field forecasts that take the spatial structure of observed fields into account and honor the flow-dependent information contained in the dynamical ensemble. The members of the spatial BMA ensemble are obtained by dressing the weather field forecasts from the dynamical ensemble with simulated spatially correlated error fields, in proportions that correspond to the BMA weights for the member models in the dynamical ensemble. Statistical ensembles of any size can be generated at minimal computational cost. The spatial BMA technique was applied to 48-h forecasts of surface temperature over the Pacific Northwest in 2004, using the University of Washington mesoscale ensemble. The spatial BMA ensemble generally outperformed the BMA and GOP ensembles and showed much better verification results than the raw ensemble, both at individual sites, for weather field forecasts, and for forecasts of composite quantities, such as average temperature in National Weather Service forecast zones and minimum temperature along the Interstate 90 Mountains to Sound Greenway.

1. Introduction

Ensemble prediction systems have been developed to generate probabilistic forecasts of weather quantities that address the two major sources of forecast uncertainty in numerical weather prediction: uncertainty in initial conditions, and uncertainty in model formulation. Originally suggested by Epstein (1969) and Leith (1974), ensemble forecasts have been operationally implemented on the synoptic scale (Toth and Kalnay 1993; Houtekamer et al. 1996; Molteni et al. 1996) and are under development on the mesoscale (Stensrud et al. 1999; Wandishin et al. 2001; Grimit and Mass 2002; Eckel and Mass 2005). In a wide range of applications, probabilistic forecasts based on ensembles provide

E-mail: veronica@stat.washington.edu

DOI: 10.1175/MWR3341.1

higher economic and societal value than a single deterministic forecast (Richardson 2000; Palmer 2002; Gneiting and Raftery 2005).

While showing significant spread-error correlations, ensemble forecasts are often biased and underdispersive (Buizza 1997; Hamill and Colucci 1997; Grimit and Mass 2002: Scherrer et al. 2004: Eckel and Mass 2005). Hence, to realize the full potential of an ensemble forecast it is necessary to apply some form of statistical postprocessing, with the goal of generating probabilistic forecasts that are calibrated and yet sharp. In the spirit of the pioneering work of Glahn and Lowry (1972), who introduced regression-type model output statistics approaches to a meteorological audience, various statistically based ensemble postprocessing techniques have been proposed. In this paper, we introduce a postprocessing technique that combines two of these methods, Bayesian model averaging (Raftery et al. 2005) and the geostatistical output perturbation technique (Gel et al. 2004a), to generate calibrated probabilistic forecasts

Corresponding author address: Veronica J. Berrocal, Department of Statistics, University of Washington, Box 354320, Seattle, WA 98195-4320.

of whole weather fields simultaneously, rather than just weather quantities at individual locations.

Bayesian model averaging (BMA) is a statistical technique originally developed for social and health science applications in situations with several competing statistical models (Hoeting et al. 1999). Raftery et al. (2005) proposed the use of BMA to calibrate forecast ensembles and generate predictive probability density functions (PDFs) for future weather quantities. The BMA predictive PDF is a weighted average of predictive PDFs associated with each individual ensemble member, with weights that reflect the member's relative skill. However, each location in the forecast domain is considered individually, and spatial correlations among errors are ignored.

The geostatistical output perturbation (GOP) method dresses a single deterministic weather field forecast with simulated error fields, to obtain statistical ensembles of weather fields that take spatial correlations into account (Gel et al. 2004a). This resembles the perturbation approach in Houtekamer and Mitchell (1998, 2001), but in the GOP technique spatially correlated perturbations are applied to the outputs of numerical weather prediction models, rather than the inputs.

In essence, the BMA technique honors ensemble information but ignores spatial correlation. The GOP method takes spatial dependencies into account, but applies to a single deterministic forecast, rather than to an ensemble of weather field forecasts, and fails to honor the flow-dependent spread that derives from the nonlinear evolution of the atmosphere and is characteristic for dynamical ensembles.

Spatial BMA addresses these shortcomings by combining the two techniques. As in the original BMA technique, the spatial BMA predictive PDF is a weighted average of forecast PDFs centered at biascorrected versions of the ensemble member models, with weights that relate to each member's performance. However, in spatial BMA the forecast PDFs are multivariate densities with covariance structures designed to honor the spatial structure of weather observations. The spatial BMA technique can be used to generate statistical ensembles of whole weather fields simultaneously, of any size, and at minimal computational cost. At any location individually, spatial BMA reduces to the original BMA technique.

The paper is organized as follows. In section 2, we review the BMA and GOP methods and describe the spatial BMA technique in detail. In section 3 we give an example of spatial BMA forecasts of surface temperature over the North American Pacific Northwest, using the University of Washington mesoscale ensemble (Grimit and Mass 2002; Eckel and Mass 2005). Section 4 presents verification results for spatial BMA forecasts in the calendar year 2004, focusing on spatial and composite quantities. The paper ends with a discussion in section 5, in which we compare the spatial BMA technique to the dressing approaches of Roulston and Smith (2003) and Wang and Bishop (2005).

2. Methods

We now describe the BMA, GOP, and spatial BMA techniques, and we explain our approach to parameter estimation.

a. Bayesian model averaging

We consider an ensemble of K weather field forecasts. In our examples, this is the eight-member University of Washington mesoscale ensemble (UWME; Eckel and Mass 2005), but BMA applies to all forecast ensembles with physically distinguishable member models, such as the poor person's or multimodel ensembles. With small modifications, BMA also applies to ensembles with exchangeable members, including bred and singular-vector ensembles (Raftery et al. 2005).

We write y for the weather quantity of interest, and f_1, \ldots, f_K for the respective ensemble member forecasts. With each ensemble member, we associate a conditional PDF, $g_k(y|f_k^0)$, which we interpret as the conditional PDF of y given that member k is the best among the ensemble member forecasts, as indicated by the superscript. The BMA predictive PDF for the weather quantity then is

$$p(y|f_1,\ldots,f_k) = \sum_{k=1}^{K} w_k g_k(y|f_k^0), \qquad (1)$$

where w_k is the probability of ensemble member k being the best. In the implementation of Raftery et al. (2005), which applies to forecasts of surface temperature and sea level pressure, the conditional PDFs are univariate normal densities centered at a linearly biascorrected forecast. Hence, $g_k(y|f_k^0)$ is a univariate normal PDF with mean $a_k + b_k f_k$ and standard deviation σ_0 , assumed to be constant across ensemble members. We denote this situation by

$$y|f_k^0 \sim \mathcal{N}(a_k + b_k f_k, \sigma_0^2). \tag{2}$$

The BMA weights in (1) and the bias and variance parameters in (2) are estimated from training data using a two-stage procedure. The bias parameters a_k and b_k are estimated for each ensemble members separately via linear least squares regression: they are the values of a_k and b_k that minimize the residual sum of squares over the entire domain. As such, they are domain specific and do not vary with location. The BMA weights w_1, \ldots, w_k and the BMA variance σ_0^2 are estimated simultaneously for all the K ensemble members using the Expectation Maximization (EM) algorithm (Dempster et al. 1977). The BMA weights reflect the relative performance of the ensemble member models during the training period; since they are probabilities, they are nonnegative and their sum is equal to 1. The best member interpretation is intuitively appealing, but it should be noted that the model in (1) is also a mixture model, where w_k represents the weight of the kth mixture component.

The BMA method as specified by (1) is implemented in the ensembleBMA package for the R language (Ihaka and Gentleman 1996), which is available online at http://cran.r-project.org.

In its original formulation, BMA yields predictive PDFs for one location at a time, and thus ignores correlations between the errors in forecasts of the same weather quantity at different locations. Typically, however, there are strong spatial correlations between these errors. If we seek the predictive PDF of a spatially aggregated quantity such as the average or minimum temperature across a region, then the spatial correlation is important. One way to proceed using BMA output would be to obtain the predictive PDF assuming that forecast errors at different locations are statistically independent, and thus uncorrelated. However, this would give an erroneous predictive PDF for an aggregated quantity if the spatial correlation was strong, as it often is.

b. The geostatistical output perturbation technique

The GOP technique dresses a single deterministic weather field forecast with Gaussian error fields that are generated using geostatistical methods (Gel et al. 2004a). Here, the deterministic weather field forecast is taken to be a member of the dynamical ensemble.

Specifically, let *S* denote a possibly large but finite set of distinct model grid points or scattered observation sites. If our intention is to produce postage stamp maps of weather field forecasts, this set is the model grid. For verification purposes, it is a collection of observation locations, and the forecasts are bilinearly interpolated from the model grid to the observation sites. We write

$$\mathbf{Y} = \{ y(s) \colon s \in S \}$$

for the weather field at the sites of interest, and $\mathbf{F}_k = \{f_k(s): s \in S\}$ for the corresponding deterministic weather field forecast. The GOP technique employs a statistical model, which assumes that

$$\mathbf{Y}|\mathbf{F}_{k} \sim \mathcal{MVN}(a_{k}\mathbf{1} + b_{k}\mathbf{F}_{k}, \mathbf{\Sigma}_{k}), \qquad (3)$$

where **1** is the vector with all components equal to 1. The right-hand side of (3) denotes a multivariate normal PDF centered at the bias-corrected member forecast, $a_k \mathbf{1} + b_k \mathbf{F}_k$, with covariance matrix $\boldsymbol{\Sigma}_k$, whose entries are specified in (4). Superficially, one might think of (3) as a spatial version of (2), but the relationships differ fundamentally: in (3), we consider \mathbf{F}_k as a single deterministic forecast without reference to any of the other ensemble members; in (2), we consider $f_k(s)$ conditionally on this member being the best among the ensemble member forecasts. This latter assumption of forecast k being the best generally implies a deflated variance in (2), when compared to (3), as will be seen below. For surface temperature and sea level pressure, the use of a multivariate normal PDF seems reasonable as an approximation, but this may not be true for other weather variables, such as precipitation or wind speed.

From here on, we refer to the difference between the observation and the bias-corrected forecast as the error. The covariance matrix in (3) describes the spatial structure of the error field and needs to be estimated from training data. Gel et al. (2004a) used a parametric, stationary, and isotropic geostatistical model, which assumes that the (i, j)th element of the covariance matrix Σ_k is

$$\rho_k^2 \delta_{ij} + \tau_k^2 \exp\left(-\frac{\|s_i - s_j\|}{r_k}\right),\tag{4}$$

where $||s_i - s_j||$ denotes the Euclidean distance between the respective locations, s_i and s_j , and δ_{ij} equals 1 if $s_i = s_j$ and is 0 otherwise. In geostatistical terminology, ρ_k^2 is called the nugget effect and represents the variance of the measurement error as well as small-scale variability, $\rho_k^2 + \tau_k^2$ is known as the sill, and r_k is called the range and indicates the rate at which the spatial correlations of the errors decay (Cressie 1993; Chilès and Delfiner 1999). In meteorological terminology, measurement error is often referred to as instrument error, and representativeness errors correspond to small-scale variability. Covariance structures that are more complex can be accommodated, and we discuss some of the options in section 5.

Note that (3) and (4) give a fully specified, multivariate normal predictive PDF for the weather field \mathbf{Y} . To generate statistical ensembles from this PDF, we express (3) and (4) in the form of the stochastic representation

$$\mathbf{Y}|\mathbf{F}_k \sim a_k \mathbf{1} + b_k \mathbf{F}_k + \mathbf{E}_{1k} + \mathbf{E}_{2k}, \tag{5}$$

where \mathbf{F}_k is the deterministic weather field forecast, a_k and b_k are scalar bias parameters, and $\mathbf{E}_{1k} = \{\epsilon_{1k}(s): s \in S\}$

April 2007

and $\mathbf{E}_{2k} = \{\epsilon_{2k}(s): s \in S\}$ are independent random vectors with mean zero, satisfying

$$\operatorname{cov}[\boldsymbol{\epsilon}_{1k}(s_i), \boldsymbol{\epsilon}_{1k}(s_j)] = \tau_k^2 \exp\left(-\frac{\|\boldsymbol{s}_i - \boldsymbol{s}_j\|}{r_k}\right),$$

and $\operatorname{cov}[\epsilon_{2k}(s_i), \epsilon_{2k}(s_j)] = \rho_k^2 \delta_{ij}$, respectively. In this representation, \mathbf{E}_{1k} is a spatially correlated error field that varies continuously with distance, and we refer to it as the continuous component of the error field. In contrast, \mathbf{E}_{2k} is a noise vector that stands for instrument and representativeness errors, and we refer to it as the discontinuous component of the error field. Statistical GOP ensembles of any size can be obtained by simulating \mathbf{E}_{1k} and \mathbf{E}_{2k} from their respective multivariate PDFs, and adding the simulated errors to the biascorrected forecast, as directed by (5). For the simulations, we use the circulant embedding technique (Wood and Chan 1994; Gneiting et al. 2006) as implemented in the RandomFields package for the R language (Schlather 2001). The GOP method is itself implemented in the ProbForecastGOP package for the R language. (All R packages are available online at http:// cran.r-project.org.)

c. Spatial BMA

We now show how to combine the BMA and GOP methods into the spatial BMA technique. Again, we consider a weather field $\mathbf{Y} = \{Y(s): s \in S\}$ at a possibly large but finite collection *S* of locations, but now conditionally on an ensemble,

$$\mathbf{F}_1 = \{ f_1(s) : s \in S \}, \dots, \mathbf{F}_K = \{ f_K(s) : s \in S \},\$$

of K weather field forecasts simultaneously, rather than just a single deterministic weather field forecast. The spatial BMA predictive PDF for the weather field is

$$p(\mathbf{Y}|\mathbf{F}_1,\ldots,\mathbf{F}_K) = \sum_{k=1}^K w_k g_k(\mathbf{Y}|\mathbf{F}_k^0), \qquad (6)$$

where w_k is the BMA weight, equal to the probability that member k is the best among the ensemble member forecasts, and $g_k(\mathbf{Y}|\mathbf{F}_k^0)$ is the conditional PDF of \mathbf{Y} given that member k is the best, as indicated by the superscript. In our implementation, the conditional PDFs are multivariate normal densities centered at the bias-corrected ensemble member forecast, $a_k \mathbf{1} + b_k \mathbf{F}_k$, and having a spatially structured covariance matrix, $\boldsymbol{\Sigma}_k^0$. By analogy to (2), we denote this situation by

$$\mathbf{Y}|\mathbf{F}_{k}^{0} \sim \mathcal{MVN}(a_{k}\mathbf{1} + b_{k}\mathbf{F}_{k}, \boldsymbol{\Sigma}_{k}^{0}).$$
(7)

In (7),

$$\boldsymbol{\Sigma}_{k}^{0} = \frac{\sigma_{0}^{2}}{\rho_{k}^{2} + \tau_{k}^{2}} \boldsymbol{\Sigma}_{k}, \qquad (8)$$

where σ_0^2 is the BMA variance in (2), Σ_k is the spatially structured GOP covariance matrix with entries specified in (4), and ρ_k^2 and τ_k^2 are the respective GOP covariance parameters. The quantity

$$\alpha_k = \frac{\sigma_0^2}{\rho_k^2 + \tau_k^2}$$

is the ratio of the BMA variance to the GOP variance for the errors, and we call it as the deflation factor for member model k, where k = 1, ..., K. Spatial BMA generalizes both the original BMA method and the GOP technique: it reduces to the former when the set S consists of a single location only, and it reduces to the latter for an ensemble of size K = 1, that is, a deterministic weather field forecast.

Similarly to GOP, the spatial BMA Eqs. (6)–(8) give a fully specified, multivariate predictive PDF for the weather field. However, it is more practical to generate a statistical ensemble of weather field forecasts, by sampling from the spatial BMA predictive PDF. Conditionally on ensemble member k being the best, we can write (7) as

$$\mathbf{Y}|\mathbf{F}_{k}^{0} \sim a_{k}\mathbf{1} + b_{k}\mathbf{F}_{k} + \mathbf{E}_{1k}^{0} + \mathbf{E}_{2k}^{0}, \qquad (9)$$

where \mathbf{E}_{1k}^{0} and \mathbf{E}_{2k}^{0} denote the continuous and the discontinuous parts of the conditional error field, respectively, with multivariate normal PDFs equal to those described in section 2 for the unconditional counterparts, \mathbf{E}_{1k} and \mathbf{E}_{2k} , except that the covariance matrix is rescaled by the deflation factor, α_k .

The following algorithm generates a member of the spatial BMA ensemble:

- Sample a number k ∈ {1,..., K}, with probabilities given by the BMA weights, w₁,..., w_K. This specifies the member of the dynamical ensemble to be dressed.
- 2) Simulate realizations of the continuous and discontinuous parts, \mathbf{E}_{1k}^{0} and \mathbf{E}_{2k}^{0} , of the conditional error field from the respective conditional PDFs.
- 3) Use the right-hand side of (9) to dress the biascorrected weather field forecast, $a_k \mathbf{1} + b_k \mathbf{F}_k$, with the simulated conditional error fields, \mathbf{E}_{1k}^0 and \mathbf{E}_{2k}^0 .

Proceeding in this manner, we obtain spatial BMA ensembles of weather field forecasts, of any desired ensemble size, and at minimal computational cost.

d. Parameter estimation

The estimation of a spatial BMA model for an underlying dynamical ensemble requires the fitting of a BMA model as well as GOP models for the individual ensemble members. This is done using prior observations and ensemble forecasts for the same prediction horizon and forecast cycle, with forecasts that are bilinearly interpolated from the model grid to the observation sites. We use a sliding training period consisting of the recent past. In deciding how long this training period should be, there is a trade-off: with a short training period the method adapts more quickly to changes in the ensemble and in its component members as well as seasonal changes. With a longer training period, on the other hand, estimation tends to be less variable. Raftery et al. (2005) showed that for 48-h BMA forecasts of surface temperature in the North American Pacific Northwest there are substantial gains in increasing the length of the training period to 25 days, but there is little gain beyond. In the examples below, we adopt this choice of a sliding 25-day training period. Other weather variables, domains, and forecast lead times may require different choices.

To fit the BMA models (1) and (2), we follow Raftery et al. (2005) in estimating the bias parameters, a_k and b_k , by linear least squares regression of the observations on the respective ensemble member forecast. The BMA weights, w_k , and the BMA variance, σ_0^2 , are estimated using the maximum likelihood technique in the form of the EM algorithm (Dempster et al. 1977). The estimate of the BMA variance σ_0^2 is then refined by searching numerically for the value of σ_0^2 that minimizes the continuous ranked probability score (CRPS; Hersbach 2000; Gneiting et al. 2005; Wilks 2006, his section 7.5.1) of BMA over the training data. This is done keeping all the other BMA parameters fixed, while searching over a range of values of σ_0^2 centered around the maximum likelihood estimate given by the EM algorithm.

It remains to fit the GOP models for the weather field forecasts using member model k, where k = 1, ..., K. Estimation of the spatial covariance parameters, ρ_k^2 , τ_k^2 and r_k in (4), is based on the fact that the GOP error field, $\epsilon_k(s) = \epsilon_{1k}(s) + \epsilon_{2k}(s)$, satisfies

$$\frac{1}{2}E[\boldsymbol{\epsilon}_k(s_i) - \boldsymbol{\epsilon}_k(s_j)]^2 = \rho_k^2 + \tau_k^2 \left[1 - \exp\left(-\frac{\|\boldsymbol{s}_i - \boldsymbol{s}_j\|}{\tau_k}\right)\right],$$

where E denotes expectation. In geostatistical language, the error field has variogram

$$\gamma_k(d) = \rho_k^2 + \tau_k^2 (1 - e^{-d/r_k}),$$

where $d = ||s_i - s_j||$ denotes the Euclidean distance between two distinct observation sites, and $\gamma_k(d)$ is onehalf the expected squared difference between errors at stations that are distance *d* apart.

We now compute the sample version of the vari-

ogram, $\hat{\gamma}_k(d)$, using data from the sliding training period, as follows:

- 1) Use the estimates of the bias-correction terms a_k and b_k previously obtained by fitting a linear least squares regression to the data from the 25-day training period.
- 2) For each day in the training period, find the empirical error field, by subtracting the bias-corrected forecast field from the corresponding field of verifying weather observations.
- 3) For each day in the training period, and for all pairs of observation locations on that day, find the distance between the sites, and compute one-half the squared difference between the errors.
- 4) Group the distances into bins B_l with midpoints d_l .
- 5) Compute the empirical variogram value $\hat{\gamma}_k(d_l)$ at distance d_l , by averaging the respective one-half squared differences over the distance bin B_l .

With this, we apply the weighted least squares technique to estimate the GOP parameters. Specifically, if n_l denotes the total number of pairs of observation sites whose distance falls into bin B_l , the weighted least squares estimates of the covariance parameters ρ_k^2 , τ_k^2 and r_k are the values that minimize

$$S(\rho_k^2, \tau_k^2, r_k) = \sum_l n_l \left\{ \frac{\hat{\gamma}_k(x_l) - [\rho_k^2 + \tau_k^2(1 - e^{-d_l/r_k})]}{\rho_k^2 + \tau_k^2(1 - e^{-d_l/r_k})} \right\}^2.$$

To solve this optimization problem, we use the quasi-Newton and conjugate-gradient techniques described by Byrd et al. (1995) and implemented in the R language (Ihaka and Gentleman 1996).

Following the estimation of the spatial covariance parameters for ensemble members k = 1, ..., K, we combine the GOP and BMA models into the spatial BMA model, using (6)–(8). We do the estimation using the previously mentioned R packages, ensembleBMA and ProbForecastGOP.

3. Example

We now give an example of 48-h spatial BMA forecasts of surface temperature over the North American Pacific Northwest, which includes Oregon, Washington, southern British Columbia, and part of the northeastern Pacific Ocean, using the UWME (Grimit and Mass 2002; Eckel and Mass 2005). In the 2004 version used here, the UWME is an eight-member multianalysis ensemble. The members use the fifth-generation Pennsylvania State University–National Center for Atmospheric Research (PSU–NCAR) Mesoscale Model (MM5; Grell et al. 2004) driven by initial and lateral

Member		Land				Ocean			
	W _k	a_k	b_k	$\sigma_0^2 (^{\circ}\mathrm{C})^2$	Wk	a_k	b_k	σ_0^2 (°C) ²	
AVN	0.11	0.93	0.90	7.78	0.03	1.13	0.87	5.20	
CMCG	0.12	0.97	0.88	7.78	0.49	1.21	0.86	5.20	
Eta	0.19	1.05	0.91	7.78	0.08	1.23	0.86	5.20	
GASP	0.00	0.88	0.87	7.78	0.00	1.05	0.87	5.20	
JMA	0.15	0.98	0.92	7.78	0.05	1.17	0.89	5.20	
NGPS	0.27	1.04	0.90	7.78	0.15	1.18	0.87	5.20	
TCWB	0.00	0.85	0.83	7.78	0.00	1.08	0.83	5.20	
UKMO	0.16	0.97	0.88	7.78	0.20	1.14	0.86	5.20	

TABLE 1. Estimates of BMA parameters for 48-h forecasts of surface temperature verifying at 0000 UTC 16 Feb 2004, using UWME.

boundary conditions supplied by eight distinct global models. Specifically; the AVN member uses initial and lateral boundary conditions from the Global Forecast System run by the National Centers for Environmental Prediction (NCEP); the CMCG member is based on the Global Environmental Multiscale model run by the Canadian Meteorological Center; the Eta Model member uses the limited-area mesoscale model run by NCEP; the GASP member is based on the Global Analysis and Prediction model run by the Australian Bureau of Meteorology; the JMA member is based on the Global Spectral Model run by the Japan Meteorological Agency; the NGPS member uses the Navy Operational Global Atmospheric Prediction System run by the Fleet Numerical Meteorology and Oceanography Center; the TCWB member is based on the Global Forecast System run by the Taiwan Central Weather Bureau; and the UKMO member derives from the Unified Model run by the Met Office. Eckel and Mass (2005) give a detailed description of UWME.

Our example is for 48-h forecasts of the surface (2 m) temperature field over the North American Pacific Northwest, initialized at 0000 UTC 14 February 2004. To deal with nonstationarities in the error fields, we divided the 12-km UWME forecast grid into two sub-

domains, land and ocean, and estimated separate spatial BMA models for the two domains, using the aforementioned 25-day sliding training period.

Table 1 shows estimates of the BMA variance, σ_0^2 , the BMA weights, w_k , and the additive and multiplicative bias, a_k and b_k , respectively, for the eight UWME members. The BMA weights differed substantially between land and ocean. The member with the highest BMA weight on land was the NGPS model, and the CMCG model had the highest weight over the ocean. The GASP and TCWB models performed poorly relative to the other members during the training period and received negligible weights in both domains. The two domains also differed in terms of the BMA variance, which was smaller over the Pacific Ocean, likely because of a decrease in the representativeness error.

Table 2 shows estimates of the GOP covariance parameters, ρ_k^2 , τ_k^2 and r_k , for the error fields, along with estimates of the deflation factor, α_k . The estimates of the nugget effect, ρ_k^2 , which subsumes instrument and representativeness errors, were much larger on land than over the Pacific Ocean. The estimates of τ_k^2 were generally somewhat larger on land than over ocean. The range, r_k , corresponds to the correlation length of the continuous component of the error field, with spa-

 TABLE 2. Estimates of spatial BMA covariance parameters and deflation factors for 48-h forecasts of surface temperature verifying at 0000 UTC 16 Feb 2004, using UWME.

	Land				Ocean			
Member	$ ho_k^2$ (°C) ²	$ au_k^2$ (°C) ²	r_k (km)	α_k	ρ_k^2 (°C) ²	$ au_k^2 (^{\circ}\mathrm{C})^2$	r_k (km)	α_k
AVN	2.26	6.30	129	0.91	1.08	5.87	258	0.75
CMCG	2.32	6.06	134	0.93	1.07	5.10	246	0.84
Eta	2.24	6.08	124	0.94	1.06	5.58	245	0.78
GASP	2.31	7.25	163	0.81	1.02	6.11	265	0.73
JMA	2.29	6.24	134	0.91	1.12	5.96	277	0.73
NGPS	2.20	5.37	105	1.03	1.05	5.16	245	0.84
TCWB	2.35	6.67	149	0.86	1.03	6.98	312	0.65
UKMO	2.29	6.39	141	0.90	0.98	5.29	211	0.83



FIG. 1. Empirical variograms of 48-h errors for surface temperature over a 25-day training period ending 14 Feb 2004, using UWME: (a) NGPS member on land and (b) CMCG member over ocean.

tial correlations decaying to about 0.05 at distance $3r_k$. The estimates of the range were larger over the ocean than on land, suggesting stronger correlations over water.

The deflation factor α_k reflects the skill of each ensemble member, with the more accurate members receiving the higher estimates. Indeed, if a member model generally performs well, then its unconditional error variance will not be very different from its conditional error variance given that it is the best among the ensemble member forecasts, and its deflation factor will be close to 1. Still, caution is needed in interpreting estimates of deflation factors. For instance, the estimated deflation factors in Table 2 were generally higher on land than they were over the ocean, and the land deflation factor for the NGPS model was larger than 1, counter to intuition. These patterns can be explained by Fig. 1, which illustrates the estimation of the GOP covariance parameters for the NGPS member on land and the CMCG member over the ocean. Each panel shows both the empirical variogram of the error field, composited over the training period, and the fitted exponential variogram. The intercept of the fitted exponential variogram equals the estimate of the nugget effect, ρ_k^2 , and corresponds to the variance of instrument and representativeness errors. The horizontal asymptote is at the estimated sill, $\rho_k^2 + \tau_k^2$, and equals the estimated marginal variance of the GOP error field. The weighted least squares technique seems to underestimate the sill for the NGPS member on land, resulting in a deflation factor that exceeds 1.

The exponential variograms fit quite well over the first 400 km, and the fit deteriorates thereafter. This is quite typical of geostatistical applications, and is not a matter of concern. Generally, when fitting a parametric

variogram model, attention is focused on the smaller distances, which are particularly relevant in characterizing the spatial statistical properties of the error fields.

Figures 2 and 3 illustrate the generation of a member of the spatial BMA ensemble on land and over the ocean, respectively. Figures 2a and 3a show the biascorrected member of the dynamical ensemble that is to be dressed. On land, this is the NGPS member, and over the ocean it is the CMCG member. Figures 2b,c and 3b,c show simulated realizations of the continuous and discontinuous components of the error field, respectively. Figures 2d and 3d show the member of the spatial BMA ensemble as the sum of the three components. Repeating this process, statistical ensembles of any size can be generated.

A characteristic feature of Figs. 2 and 3, and in general of the spatial BMA (and GOP) ensemble member fields, is an increase in roughness compared to weather fields generated by numerical weather prediction models. This stems from spatial BMA aiming to reproduce the spatial structure of weather observations, including instrument and representativeness errors, represented by the discontinuous component of the error field (Figs. 2c and 3c). The discontinuous component can be ignored, if desired, and the spatial BMA technique can be implemented by adding the bias-corrected weather field forecast (Figs. 2a and 3a) and the continuous component of the simulated error field (Figs. 2b and 3b) only. This is an implementation decision that needs to be made depending on the prediction problem at hand. In our implementation, we added both components of the error to the bias-corrected forecast field.

The continuous component of the error field generally contributes more than the discontinuous component, since the estimates of the covariance parameter



FIG. 2. A member of the spatial BMA ensemble for 48-h forecasts of surface temperature over the land portion of the Pacific Northwest, initialized at 0000 UTC 14 Feb 2004: adding (a) the bias-corrected UWME NGPS weather field forecast, (b) the continuous, and (c) the discontinuous component of the simulated error field, we obtain (d) a member of the spatial BMA ensemble. Note that different color scales are used in (a)–(d) to make it easier to see the patterns in each one.

 τ_k^2 , which represents the marginal variance of the continuous component, are substantially larger than the estimates of the nugget effect ρ_k^2 , the marginal variance of the discontinuous component.

For each region, we generated a spatial BMA ensemble of 19 weather fields. These could be displayed in the form of a postage stamp plot, but this would be likely to overwhelm users, and plots that summarize the spatial BMA ensemble are likely to be more useful. Ensemble forecasts for all types of composite quantities can be derived from the statistical ensemble. For instance, we might be interested in predicting the empirical variogram of the temperature field verifying at 0000 UTC 16 February 2004. We computed the empirical variogram for each of the 19 members of the spatial BMA ensemble, using 300 distance bins. At each bin, the minimum and the maximum of the respective 19 values envelop a 95% prediction interval for the verifying variogram value, which we computed from the observed temperature field. Figure 4 shows the results of this experiment. The prediction intervals generally cover the verifying empirical variogram values.

4. Verification results

In calendar year 2004, the 0000 UTC cycle for the 12-km domain of the eight-member UWME was run on 245 days. For each day, we fitted BMA, GOP, and spatial BMA models for 48-h forecasts of surface (2 m) temperature over the North American Pacific Northwest, separately on land and over the ocean, using a sliding 25-day training period. We then generated original BMA, GOP, and spatial BMA forecast ensembles for each day. The original BMA ensembles were created by sampling from the univariate original BMA predictive PDFs at each location separately, incorrectly assuming spatial independence of the error fields. The GOP ensembles were based on the UWME UKMO



FIG. 3. Same as in Fig. 2, but for the UWME CMCG member and over the Pacific Ocean.

model, which had the best aggregate performance among the ensemble member models, both on land and over the ocean. In the interest of a fair comparison to the eight-member UWME, our GOP, original BMA, and spatial BMA ensembles also had eight members only. However, the statistical ensembles allow for ensembles of any size, and larger ensembles frequently show better verification results.

We now assess and rank the performance of the UWME, GOP, original BMA, and spatial BMA ensembles, emphasizing spatial and composite quantities. On average, observations of surface temperature were available at 761 stations on land and 196 stations over the Pacific Ocean. We verified bilinearly interpolated ensemble forecasts against the temperature observations.

In contrast to the statistical ensembles, UWME is not designed to take instrument and representativeness errors into account. Hence, we consider a fifth ensemble, which we call UWME + noise. To create the UWME + noise ensemble, we added Gaussian noise to each of the eight UWME members, at each site independently, with mean zero and a variance that equals the estimated nugget effect, ρ_k^2 , for the corresponding member model.

a. Temperature forecasts at individual sites

We begin by assessing surface temperature forecasts at individual sites. For forecasts at single sites, spatial BMA, and original BMA are equivalent; hence the results for the two ensembles are essentially identical, with any differences due to chance variability in the generation of the ensemble members. All verification statistics were spatially and temporally composited over the Pacific Northwest and the 2004 calendar year.

Table 3 shows the mean absolute error (MAE) and the average CRPS (Hersbach 2000; Gneiting et al. 2005; Wilks 2006, his section 7.5.1) for the various ensemble methods. The MAE assesses the accuracy of the deterministic forecasts. The UWME and UWME + noise



FIG. 4. Empirical variogram values (dots) for the verifying surface temperature field at 0000 UTC 16 Feb 2004, and pointwise minimum and maximum of the empirical variogram values (lines) from the 19-member spatial BMA weather field ensemble (a) on land and (b) over ocean.

deterministic forecast is the raw ensemble mean; for the GOP method it is the bias-corrected UWME UKMO forecast; and for the original BMA and spatial BMA techniques it is a weighted average of the bias-corrected ensemble member forecasts. The CRPS is a scoring rule for predictive PDFs that addresses calibration as well as sharpness, and is proper, that is, discourages hedging. The CRPS generalizes the absolute error, to which it reduces for deterministic forecasts; it is also reported in degrees Celsius, and average CRPS values can be directly compared to the MAE (Gneiting et al. 2005). A clear rank order can be observed, in that the BMA ensembles showed substantially lower CRPS values than the GOP ensemble, followed by the UWME + noise and UWME ensembles.

To assess the calibration of the ensemble forecasts, we use the verification rank histogram (Anderson 1996; Talagrand et al. 1997; Hamill and Colucci 1997; Hamill 2001). Figure 5 shows the histograms for the various ensembles. We also computed the respective discrepancy from uniformity,

$$D = \sum_{j=1}^{K+1} \left| p_j - \frac{1}{K+1} \right|, \tag{10}$$

where K = 8 is the number of ensemble members and p_j is the observed relative frequency of rank *j*. The smaller the discrepancy, the smaller the deviation from a uniform rank histogram, and the better the calibration. In both domains, the UWME and UWME + noise ensembles were underdispersive, while the GOP, original BMA, and spatial BMA ensembles had rank histograms that were nearly uniform. The slight overdispersion of the rank histograms of the original and spatial

BMA ensembles over ocean can be attributed to the smaller number of cases available over ocean, compared to land.

b. Temperature field forecasts

To assess the calibration of the ensembles as weather field forecasts, rather than as forecasts of weather quantities at individual sites, we use a variant of the verification rank histogram that is tailored to this task, namely the minimum spanning tree (MST) rank histogram (Smith and Hansen 2004; Wilks 2004). An MST rank $k \in \{1, \ldots, K + 1\}$ is computed based on each day's ensemble of weather field forecasts and the verifying weather field. This yields 245 MST ranks for each of the five ensemble techniques, and the corresponding histogram is uniform if the ensemble is calibrated. Figure 6 shows the MST rank histograms for the UWME, UWME + noise, GOP, original BMA, and spatial BMA weather field ensembles, separately on land and over the ocean, along with the discrepancy (10) that measures the departure from uniformity. The UWME, UWME + noise, and GOP weather field ensembles were severely underdispersive. The original BMA en-

TABLE 3. MAE and average CRPS for 48-h forecasts of surface temperature over the Pacific Northwest in 2004 (°C).

	La	and	Ocean		
Ensemble	MAE	CRPS	MAE	CRPS	
UWME	2.94	2.58	2.44	2.12	
UWME + noise	2.94	2.23	2.44	1.89	
GOP	2.71	2.13	2.35	1.82	
Original BMA	2.70	1.95	2.35	1.72	
Spatial BMA	2.70	1.95	2.35	1.72	



FIG. 5. Verification rank histograms for 48-h forecasts of surface temperature over the Pacific Northwest in 2004 for individual stations (top) on land and (bottom) over the ocean.

semble also was underdispersive, but to a lesser extent. The MST rank histograms for the spatial BMA weather field ensemble departed the least from uniformity. The difference between the GOP and spatial BMA ensembles corroborates the widely held perception that it is advantageous to take account of the flow-dependent information contained in the dynamical ensemble.

As an alternative approach to the spatial verification of ensembles of weather field forecasts, we repeated the variogram computations in Fig. 4 for the 245 available days in 2004 and the five types of weather field ensembles. Each eight-member ensemble supplies nominal $7/9 \times 100\% = 77.8\%$ prediction intervals for variogram values computed from the verifying temperature field. If an ensemble is faithfully reproducing the spatial structure of the observed weather field, we would expect that the prediction intervals for variogram values constructed using the variogram of the ensemble forecasts would contain the variogram of the verifying temperature field 78 times out of 100. If this is not the case, the ensemble may not be reproducing the spatial structure of the observed weather field.

Table 4 shows the empirical coverage of the prediction intervals when composited over the 245 days and 300 distance bins. For all five types of ensembles, the empirical coverage was lower than desired, but the coverage for the GOP and spatial BMA ensembles was closest to the nominal 77.8%.



FIG. 6. MST rank histograms for 48-h weather field forecasts of surface temperature over the Pacific Northwest in 2004 (top) on land and (bottom) over the ocean.

 TABLE 4. Coverage of nominal 77.8% prediction intervals for variogram values.

Ensemble	Land	Ocean
UWME	20.5	28.8
UWME + noise	36.3	42.8
GOP	56.6	58.7
Original BMA	30.9	46.3
Spatial BMA	60.1	57.1

c. Average temperature in National Weather Service forecast zones

Spatial correlations play crucial roles in the prediction of a number of composite quantities. Here, we present verification results for ensemble forecasts of spatial averages of temperature. Figure 7 shows the 44 National Weather Service (NWS) forecast zones in the state of Washington. For each zone and each day, we considered ensemble forecasts of average surface temperature, understood as the mean of the temperature observations at the stations within the zone.

Figure 8 summarizes verification statistics for the various types of eight-member ensembles in the 44 zones. The performance of the GOP ensemble was almost identical to that of the spatial BMA ensemble, and we omit the corresponding results. Figure 8a shows the discrepancy (10) that measures the departure of the verification rank histogram from uniformity. In almost all zones, the spatial BMA ensemble showed the lowest discrepancy. Figure 9 illustrates this for forecast zone 7,

which has one of the highest numbers of stations and contains the city of Seattle, Washington. Both UWME, UWME + noise, and the original BMA ensemble were underdispersive, while the GOP and spatial BMA ensembles had verification rank histograms that were similar to each other and close to being uniform. The underdispersion of the original BMA ensemble is not surprising, in that the assumption of spatial independence of errors implies an underestimation of the variance of temperature averages. The slight overdispersion of spatial BMA and GOP may reflect the small number of cases used to construct the rank histogram (only 245).

Figure 8b shows the average range of the forecast ensemble for the various types of ensembles. The range quantifies the sharpness of the predictive distributions and is simply the difference between the maximum and the minimum of the eight ensemble values. The UWME had the sharpest predictive distributions, but it was underdispersive, and therefore uncalibrated. A similar comment applies to the original BMA ensemble. The spatial BMA ensemble was the least sharp, but it was better calibrated than the other types of ensembles. Finally, to assess calibration and sharpness simultaneously, Fig. 8c shows the aggregate CRPS values. Despite being sharpest, the UWME generally had the highest, least desirable CRPS values. The original BMA and the spatial BMA ensembles had CRPS values that were lower, and quite similar to each other, even though the ensembles behaved quite differently in



FIG. 7. NWS forecast zones in the state of WA, bordered by the Pacific Ocean to the west, and British Columbia, ID, and OR to the north, east, and south, respectively. (See www.atmos. washington.edu/data/images/zone.gif.)



FIG. 8. Verification statistics for 48-h forecasts of average surface temperature in NWS forecast zones in 2004. (a) Verification rank histogram discrepancy, (b) mean ensemble width (°C), and (c) mean CRPS value (°C).

terms of calibration and sharpness. There is a trade-off between calibration and sharpness, in that the goal of probabilistic forecasting is to produce a predictive distribution that is as concentrated as possible, yet calibrated, meaning that it is statistically consistent with the distribution of the observations (Gneiting et al. 2005). From the perspective of maximizing sharpness subject to calibration, the performance of the spatial BMA ensemble is superior.

d. Minimum temperature along Interstate 90

We now present verification results for another composite quantity: minimum temperature along the Interstate 90 Mountains to Sound Greenway, Washington's primary east–west-bound highway. Accommodating 20 million travelers annually, Interstate 90 crosses the Cascade Mountains in a dramatic mountain landscape with substantial altitude differentials. Accurate and reliable forecasts of minimum temperature are critical to highway maintenance operations.

Figure 10 shows the locations of 13 meteorological stations along the Cascades section of Interstate 90, some of which are very near each other. We consider ensemble forecasts of the minimum temperature among these 13 stations. The UWME forecasts were available on the 12-km model grid and were bilinearly interpolated to the observation locations. However, In-

terstate 90 and the meteorological stations are generally located at lower altitudes, while the surrounding grid points are at higher altitudes. On average, there is a difference of 264 m between the station height and the respective heights of the surrounding model grid points. Hence, altitude is a critical consideration, and we applied a standard lapse rate correction of 0.65° C $(100 \text{ m})^{-1}$ to all five types of forecast ensembles at all 13 stations.

Figure 11 shows verification rank histograms for the eight-member UWME, UWME + noise, GOP, original BMA, and spatial BMA forecast ensembles. The UWME and UWME + noise ensembles were underdispersive. The original BMA ensemble was strongly biased, tending to underestimate the minimum temperature along Interstate 90. Indeed, the minimum of a collection of independent forecasts tends to be smaller than the minimum of a collection of forecasts that are spatially correlated. The GOP and spatial BMA ensembles had rank histograms that were close to being uniform, and their slight overdispersion can be explained in terms of the small number of events used to construct the histogram. Table 5 shows the verification rank histogram discrepancy, the mean ensemble range, and the mean CRPS value for the forecast ensembles. The UWME, GOP, and spatial BMA ensembles showed similar CRPS values, thereby illustrating a



FIG. 9. Verification rank histograms for 48-h forecasts of average surface temperature in NWS forecast zone 7 in 2004.



FIG. 10. Meteorological stations along the Cascades corridor of Interstate 90.

trade-off between calibration and sharpness. In view of our goal of maximizing sharpness under the constraint of calibration, we contend that the GOP and spatial BMA ensembles are preferable for most users.

5. Discussion

We have introduced the spatial BMA method, a statistical postprocessing technique for calibrating forecast ensembles of whole weather fields simultaneously. Spatial BMA generalizes and combines Bayesian model averaging (BMA) and the geostatistical output perturbation (GOP) technique, and it honors ensemble as well as spatial statistical information. The spatial BMA predictive PDF for the weather field is a weighted average of multivariate normal PDFs centered at biascorrected members of the dynamical forecast ensemble. At any single location, spatial BMA reduces to the original BMA technique. It is computationally inexpensive and can be used to generate statistical ensembles of any size.

In experiments with the University of Washington mesoscale ensemble, the Spatial BMA ensemble compared favorably to the raw dynamical ensemble, the raw ensemble with added observational noise, the GOP ensemble, and the original BMA ensemble. In particular, the minimum spanning tree histogram, a key tool in assessing the calibration of ensembles of weather field forecasts (Smith and Hansen 2004; Wilks 2004), was closest to being uniform for the spatial BMA ensemble. For forecasts of composite quantities, such as temperature averages over NWS forecast zones and minimum temperature along the Cascades corridor of Interstate 90, the GOP ensemble and the spatial BMA ensemble showed similar performances, and outperformed the other types of ensembles. While our experiments were with surface temperature fields, spatial BMA in its present form applies to all weather variables with forecast error distributions that are approximately Gaussian, including sea level pressure. Further research is needed to extend spatial BMA to other weather variables, such as precipitation or wind speed. Sloughter et al. (2007) presented a non-Gaussian version of BMA that yields calibrated quantitative probabilistic precipitation forecasts at individual sites, but not for weather fields.

There are several directions in which the spatial BMA technique could be developed. One is bias correction. In the current implementation, we use a simple linear bias correction that does not take altitude, land use, latitude, longitude, or distance from the ocean into account. More sophisticated regression based bias removal techniques might include some or all of these quantities as predictor variables. Another possibility is to use a nearest-neighbor approach based on distance, altitude, and land use categories.

Another possibility would be to reduce the bias correction to a simple additive bias correction, including only the term a_k and fixing b_k to 1. This would reduce the number of parameters to estimate. In their implementation of the original BMA technique for the Canadian Ensemble System, Wilson et al. (2007) found that for training periods of up to 50 days and for forecast lead times up to 7 days, using only a_k in the bias removal step performed as well as using both a_k and b_k .

Another way to reduce the number of parameters to estimate on each day would be to estimate the covari-



FIG. 11. Verification rank histograms for 48-h forecasts of minimum temperature along Interstate 90.

TABLE 5. Verification rank histogram discrepancy, mean ensemble width, and mean CRPS value for ensemble forecasts of minimum temperature along Interstate 90.

Ensemble	Discrepancy	Range (°C)	CRPS (°C)
UWME	0.68	2.76	1.55
UWME + noise	0.47	4.37	1.67
GOP	0.18	7.75	1.53
Original BMA	0.95	6.48	2.76
Spatial BMA	0.22	8.04	1.54

ance parameters for each ensemble member only once, using data from a previous year. The regression parameters a_k and b_k , the BMA weights, w_1, \ldots, w_8 , and the BMA standard deviation σ_0 would still have to be estimated using the training period, and the new estimate of the BMA variance would be employed to compute the deflaction factor α_k for each ensemble member.

In modeling the covariance structure of the error fields, we used a stationary and isotropic exponential correlation function. There are several ways in which more complex, and potentially more realistic, covariance structures could be used. Stationary and isotropic correlation functions that are more versatile than an exponential function are available (Mitchell et al. 1990; Gneiting 1999). Anisotropic covariance structures could also be used (Purser et al. 2003); however, in the case of surface temperature over the Pacific Northwest, Gel et al. (2004b) did not find any significant differences between longitudinal and latitudinal empirical variograms of the forecast error fields. Finally, nonstationary covariance models, that is, models that are not translation invariant, could be used. In our experiments, we dealt with nonstationarities between the land and the ocean by fitting and generating two distinct spatial BMA ensembles, each of which used a stationary and isotropic covariance structure. This was a fairly simple way to resolve nonstationarities, and yet produced good results. The methods of Paciorek and Schervish (2006) could be used to fit valid covariance structures that are stationary on homogeneous domains, yet nonstationary globally, thereby allowing for the generation of a single spatial BMA ensemble over all domains simultaneously, without incurring discontinuities along the boundaries.

An issue not explicitly considered in the spatial BMA approach is that of phase or displacement errors. These could perhaps be addressed by partitioning the errors of the ensemble member weather field forecasts into displacement, distortion, amplitude, and residual fields, as in Du et al. (2000), and applying the spatial BMA technique to the residual component only, while developing parametric statistical models for displacement, distortion, and amplitude errors. This would be an interesting avenue for future research, with potential rewards in the form of sharper yet calibrated forecast PDFs, but may require impracticably large sets of training data.

Another issue that calls for discussion is the choice of the training period. In the current implementation, we use forecast and observation data from a sliding window consisting of the 25 most recent days available to estimate the spatial BMA parameters. This allows the method to adapt rapidly to seasonal changes in the atmosphere as well as changes in the design of the ensemble, but limits the availability of training data. However, even with limited data, we did not have the problem of overfitting when estimating all the spatial BMA parameters, as our results on the quality of the spatial BMA out-of-sample predictions indicate. In addition, the estimates of the parameters vary smoothly with time, and the variability from day to day is not great. A potential way of increasing the amount of training data is to also use training data from the same season in previous years; this could be done using ensemble reforecasts, as proposed by Hamill et al. (2004). However, reforecasts put high demands on computational and human resources, and they were not available to us.

We close by comparing spatial BMA to other ensemble postprocessing techniques. Wilks (2002) proposed fitting mixtures of multivariate normal densities to ensemble forecasts of multivariate weather quantities. This resembles the spatial BMA technique, but does not take bias and calibration adjustments into account. Roulston and Smith (2003) proposed combining statistical and dynamical ensembles, and suggested the use of hybrid ensembles, in which the members of the dynamical ensemble are dressed with errors drawn from an archive of the best member errors. A difficulty in this approach is the identification of the best members. Wang and Bishop (2005) showed that under a wide range of scenarios the best member dressing method fails to be calibrated. They proposed a modified dressing technique, in which statistical perturbations are generated, with flexible covariance structures that are estimated from training data. This is similar to the spatial BMA technique, in that the Wang and Bishop (2005) predictive PDF is also a weighted average of multivariate normal densities, each centered at a bias-corrected member of the dynamical forecast ensemble, but the weights are all equal and do not depend on the member's skill. Fortin et al. (2006) proposed dressing kernels that depend on the rank of the member within the ensemble. This method is tailored to ensembles with exchangeable members, as opposed to the University of Washington mesoscale ensemble, or any poor person's ensemble, for which the members are not exchangeable. Raftery et al. (2005, p. 1170) discuss an adaptation of BMA to the case of exchangeable members.

Acknowledgments. We thank Jeff Baars, Fadoua Balabdaoui, F. Anthony Eckel, Yulia Gel, Eric P. Grimit, Nicholas A. Johnson, Clifford F. Mass, J. McLean Sloughter, and Patrick Tewson, for sharing code, comments, and/or data. This research was supported by the DoD Multidisciplinary University Research Initiative (MURI) program administered by the Office of Naval Research under Grant N00014-01-10745.

REFERENCES

- Anderson, J. L., 1996: A method for producing and evaluating probabilistic forecasts from ensemble model integrations. J. *Climate*, 9, 1518–1530.
- Buizza, R., 1997: Potential forecast skill of ensemble prediction and spread and skill distribution of the ECMWF ensemble prediction system. *Mon. Wea. Rev.*, **125**, 99–119.
- Byrd, R. H., P. Lu, J. Nocedal, and C. Zhu, 1995: A limited memory algorithm for bound constrained optimization. *SIAM J. Sci. Comput.*, **16**, 1190–1208.
- Chilès, J.-P., and P. Delfiner, 1999: *Geostatistics: Modeling Spatial Uncertainty.* Wiley, 695 pp.
- Cressie, N. A. C., 1993: *Statistics for Spatial Data*. rev. ed. Wiley, 900 pp.
- Dempster, A. P., N. M. Laird, and D. B. Rubin, 1977: Maximum likelihood from incomplete data via the EM algorithm. J. Roy. Stat. Soc. B, 39, 1–39.
- Du, J., S. Mullen, and F. Sanders, 2000: Removal of distortion error from an ensemble forecast. *Mon. Wea. Rev.*, **128**, 3347– 3351.
- Eckel, F. A., and C. F. Mass, 2005: Aspects of effective mesoscale, short-range ensemble forecasting. *Wea. Forecasting*, **20**, 328– 350.
- Epstein, E. S., 1969: Stochastic dynamic prediction. *Tellus*, **21**, 739–759.
- Fortin, V., A.-C. Favre, and M. Saïd, 2006: Probabilistic forecasting from ensemble prediction systems: Improving upon the best-member method by using a different weight and dressing kernel for each member. *Quart. J. Roy. Meteor. Soc.*, 132, 1349–1369.
- Gel, Y., A. E. Raftery, and T. Gneiting, 2004a: Calibrated probabilistic mesoscale weather field forecasting: The geostatistical output perturbation (GOP) method (with discussion). J. Amer. Stat. Assoc., 99, 575–588.
- —, —, , and V. J. Berrocal, 2004b: Calibrated probabilistic mesoscale weather field forecasting: The geostatistical output perturbation (GOP) method: Rejoinder. J. Amer. Stat. Assoc., 99, 588–590.
- Glahn, H. R., and D. A. Lowry, 1972: The use of model output statistics (MOS) in objective weather forecasting. J. Appl. Meteor., 11, 1203–1211.
- Gneiting, T., 1999: Correlation functions for atmospheric data analysis. Quart. J. Roy. Meteor. Soc., 125, 2449–2464.
- —, and A. E. Raftery, 2005: Weather forecasting using ensemble methods. *Science*, **310**, 248–249.
- ____, ___, A. H. Westveld, and T. Goldman, 2005: Calibrated

probabilistic forecasting using ensemble model output statistics and minimum CRPS estimation. *Mon. Wea. Rev.*, **133**, 1098–1118.

- —, H. Ševčíková, D. B. Percival, M. Schlather, and Y. Jiang, 2006: Fast and exact simulation of large Gaussian lattice systems in R²: Exploring the limits. *J. Comput. Graph. Stat.*, **15**, 483–501.
- Grell, G. A., J. Dudhia, and D. R. Stauffer, 2004: A description of the fifth-generation Penn State/NCAR Mesoscale Model (MM5). NCAR Tech. Note NCAR/TN-398+STR, 121 pp. [Available from MMM Division, NCAR, P.O. Box 3000, Boulder, CO 80307.]
- Grimit, E. P., and C. F. Mass, 2002: Initial results of a mesoscale short-range ensemble forecasting system over the Pacific Northwest. *Wea. Forecasting*, **17**, 192–205.
- Hamill, T. M., 2001: Interpretation of rank histograms for verifying ensemble forecasts. *Mon. Wea. Rev.*, **129**, 550–560.
- —, and S. J. Colucci, 1997: Verification of Eta-RSM short-range ensemble forecasts. *Mon. Wea. Rev.*, **125**, 1312–1327.
- —, J. S. Whitaker, and X. Wei, 2004: Ensemble reforecasting: Improving medium-range forecast skill using retrospective forecasts. *Mon. Wea. Rev.*, **132**, 1434–1447.
- Hersbach, H., 2000: Decomposition of the continuous ranked probability score for ensemble prediction systems. *Wea. Forecasting*, **15**, 559–570.
- Hoeting, J. A., D. M. Madigan, A. E. Raftery, and C. T. Volinsky, 1999: Bayesian model averaging: A tutorial. *Stat. Sci.*, 14, 382–401. [A corrected version is available online at www.stat. washington.edu/www/research/online/hoeting1999.pdf.]
- Houtekamer, P. L., and H. L. Mitchell, 1998: Data assimilation using an ensemble Kalman filter technique. *Mon. Wea. Rev.*, **126**, 796–811.
- —, and —, 2001: A sequential ensemble Kalman filter for atmospheric data assimilation. *Mon. Wea. Rev.*, **129**, 123–137.
- —, L. Lefaivre, J. Derome, H. Ritchie, and H. L. Mitchell, 1996: A system simulation approach to ensemble prediction. *Mon. Wea. Rev.*, **124**, 1225–1242.
- Ihaka, R., and R. Gentleman, 1996: R: A language for data analysis and graphics. J. Comput. Graph. Stat., 5, 299–314.
- Leith, C. E., 1974: Theoretical skill of Monte Carlo forecasts. Mon. Wea. Rev., **102**, 409–418.
- Mitchell, H. L., C. Charette, C. Chouinard, and B. Brasnett, 1990: Revised interpolation statistics for the Canadian data assimilation procedure: Their derivation and application. *Mon. Wea. Rev.*, **118**, 1591–1614.
- Molteni, R., R. Buizza, T. N. Palmer, and T. Petroliagis, 1996: The ECMWF ensemble system: Methodology and validation. *Quart. J. Roy. Meteor. Soc.*, **122**, 73–119.
- Paciorek, C. J., and M. J. Schervish, 2006: Spatial modelling using a new class of nonstationary covariance functions. *Environmetrics*, **17**, 483–506.
- Palmer, T. N., 2002: The economic value of ensemble forecasts as a tool for assessment: From days to decades. *Quart. J. Roy. Meteor. Soc.*, **128**, 747–774.
- Purser, R. J., W. S. Wu, D. F. Parrish, and N. M. Roberts, 2003: Numerical aspects of the application of recursive filters to variational statistical analysis. Part II: Spatially inhomogeneous and anisotropic general covariances. *Mon. Wea. Rev.*, 131, 1536–1548.
- Raftery, A. E., T. Gneiting, F. Balabdaoui, and M. Polakowski, 2005: Using Bayesian model averaging to calibrate forecast ensembles. *Mon. Wea. Rev.*, **133**, 1155–1174.
- Richardson, D. S., 2000: Skill and relative economic value of the

ECMWF ensemble prediction system. *Quart. J. Roy. Meteor. Soc.*, **126**, 649–667.

- Roulston, M. S., and L. A. Smith, 2003: Combining dynamical and statistical ensembles. *Tellus*, 55A, 16–30.
- Scherrer, S. C., C. Appenzeller, P. Eckert, and D. Cattani, 2004: Analysis of the spread-skill relations using the ECMWF ensemble prediction system over Europe. *Wea. Forecasting*, **19**, 552–565.
- Schlather, M., 2001: Simulation and analysis of random fields. *R News*, **1** (2), 18–20.
- Sloughter, J. M., A. E. Raftery, T. Gneiting, and C. Fraley, 2007: Probabilistic quantitative precipitation forecasting using Bayesian model averaging. *Mon. Wea. Rev.*, in press.
- Smith, L. A., and J. A. Hansen, 2004: Extending the limits of ensemble forecast verification with the minimum spanning tree. *Mon. Wea. Rev.*, 132, 1522–1528.
- Stensrud, D. J., H. E. Brooks, J. Du, M. S. Tracton, and E. Rogers, 1999: Using ensembles for short-range forecasting. *Mon. Wea. Rev.*, **127**, 433–446.
- Talagrand, O., R. Vautard, and B. Strauss, 1997: Evaluation of probabilistic prediction systems. Proc. ECMWF Workshop on Predictability, Reading, United Kingdom, ECMWF, 1–25.
- Toth, Z., and E. Kalnay, 1993: Ensemble forecasting at the NMC:

The generation of perturbations. Bull. Amer. Meteor. Soc., 74, 2317–2330.

- Wandishin, M. S., S. L. Mullen, D. J. Stensrud, and H. E. Brooks, 2001: Evaluation of a short-range multimodel ensemble system. *Mon. Wea. Rev.*, **129**, 729–747.
- Wang, X., and C. H. Bishop, 2005: Improvement of ensemble reliability with a new dressing kernel. *Quart. J. Roy. Meteor. Soc.*, 131, 965–986.
- Wilks, D. S., 2002: Smoothing forecast ensembles with fitted probability distributions. *Quart. J. Roy. Meteor. Soc.*, **128**, 2821– 2836.
- —, 2004: The minimum spanning tree histogram as a verification tool for multidimensional ensemble forecasts. *Mon. Wea. Rev.*, **132**, 1329–1340.
- Wilson, L. J., S. Beauregard, A. E. Raftery, and R. Verret, 2007: Calibrated surface temperature forecasts from the Canadian ensemble prediction system using Bayesian model averaging. *Mon. Wea. Rev.*, **135**, 1364–1385.
- Wood, A. T. A., and G. Chan, 1994: Simulation of stationary Gaussian processes in [0, 1]^d. J. Comput. Graph. Stat., 3, 409–432.