

Inaccuracies with Multimodel Postprocessing Methods Involving Weighted, Regression-Corrected Forecasts

DANIEL HODYSS, ELIZABETH SATTERFIELD, AND JUSTIN MCLAY

Naval Research Laboratory, Monterey, California

THOMAS M. HAMILL

Physical Sciences Division, NOAA/Earth System Research Laboratory, Boulder, Colorado

MICHAEL SCHEUERER

Cooperative Institute for Research in Environmental Sciences, University of Colorado, Boulder, Colorado

(Manuscript received 27 May 2015, in final form 27 October 2015)

ABSTRACT

Ensemble postprocessing is frequently applied to correct biases and deficiencies in the spread of ensemble forecasts. Methods involving weighted, regression-corrected forecasts address the typical biases and underdispersion of ensembles through a regression correction of ensemble members followed by the generation of a probability density function (PDF) from the weighted sum of kernels fit around each corrected member. The weighting step accounts for the situation where the ensemble is constructed from different model forecasts or generated in some way that creates ensemble members that do not represent equally likely states. In the present work, it is shown that an overweighting of climatology in weighted, regression-corrected forecasts can occur when one first performs a regression-based correction before weighting each member. This overweighting of climatology results in an increase in the mean-squared error of the mean of the predicted PDF. The overweighting of climatology is illustrated in a simulation study and a real-data study, where the reference is generated through a direct application of Bayes's rule. The real-data example is a comparison of a particular method referred to as Bayesian model averaging (BMA) and a direct application of Bayes's rule for ocean wave heights using U.S. Navy and National Weather Service global deterministic forecasts. This direct application of Bayes's rule is shown to not overweight climatology and may be a low-cost replacement for the generally more expensive weighted, regression-correction methods.

1. Introduction

Ensemble prediction methods are now ubiquitous in weather and climate prediction. Multiple forecast simulations are generated, typically initialized with distinct model states that approximate draws from the distribution of analysis uncertainty. Effects of model imperfections on forecast uncertainty may be simulated through the use of multiple forecast models, multiple parameterization suites, a diversity of constants in the parameterizations, and/or stochastic prediction techniques. While ensemble weather prediction systems

have improved greatly, the predictions are still frequently affected by systematic errors, including biased ensemble mean forecasts and often an insufficiency of ensemble spread. Consequently, much attention has been paid in recent years to statistical postprocessing techniques, whereby the current guidance is adjusted based on relationships noted between past forecasts and observations/analyses. In many circumstances, the goal is to produce a probability density function (PDF) that is as sharp as possible while remaining reliable (Gneiting et al. 2007).

Many approaches have been proposed for statistical postprocessing of ensemble forecasts. Non-homogeneous Gaussian regression (NGR; Gneiting et al. 2005) assumes an underlying Gaussian distribution for the posterior and estimates a state-dependent mean and spread of the distribution using the ensemble mean

Corresponding author address: Dr. Daniel Hodyss, Naval Research Laboratory, Marine Meteorology Division, 7 Grace Hopper Ave., Stop 2, Monterey, CA 93943.
E-mail: daniel.hodyss@nrlmry.navy.mil

and spread as predictors. Forecast analogs (Hamill and Whitaker 2006; Delle Monache et al. 2013; Hamill et al. 2015) have been demonstrated to work well with large training sample sizes.

Other methods, however, approach the postprocessing problem, either explicitly or implicitly, as a problem in kernel density estimation (KDM). In these methods, kernels are placed at predetermined values of the variable to be postprocessed in order to enhance the ensemble size and/or smooth the resulting shape of the PDF for probabilistic forecasting. Examples of such methods include ensemble dressing (Roulston and Smith 2003; Wang and Bishop 2003; Fortin et al. 2006), which increases the ensemble size and spread by adding synthetic new ensemble members centered on existing or adjusted member forecasts; Bayesian model averaging (BMA; Raftery et al. 2005, hereafter R05; Wilson et al. 2007); and related techniques like ensemble kernel density model output statistics (EKDMOS; Glahn et al. 2009) and ensemble regression (Unger et al. 2009), which statistically adjust raw ensemble forecast guidance and produce PDFs through a weighted sum of kernels centered on the adjusted forecasts.

In BMA, the postprocessing takes place in multiple steps, following R05. If members have exchangeable errors, then this first step has been applied in several different ways: regression correction on the ensemble mean (Hamill 2007), bulk bias removal (Wilson et al. 2007), and separate regression corrections of each forecast (R05; Fraley et al. 2010). If members do not have exchangeable error statistics, then linear regression corrections are commonly applied individually to each member. The second step of BMA is the application of kernels to each weighted member. R05 made the implicit assumption that regression-corrected forecasts were similar in quality, so that the kernel standard deviation could simply be set to the same value for each forecast. With regression-corrected forecasts of very different quality, different kernel standard deviations are possible. In R05, the weights for each forecast and the kernel width were estimated with expectation-maximization methods (EM; Dempster et al. 1977 and references therein), though Vrugt et al. (2008) also showed that Markov chain Monte Carlo methods can also be used.

Since R05, there have been a number of critiques of BMA and related methods. Hamill (2007) discussed issues related to overfitting and Bishop and Shanley (2008) discussed issues with extreme forecasts. Wilks (2006) critiqued techniques like BMA that regress each member, demonstrating that for longer-lead forecasts, there is a tendency for the regressed values to coalesce toward a single number related to

the climatology. In such a situation, BMA typically fits wide kernels to make up for the loss of spread in the original ensemble.

In this manuscript, we show that when a weighting on each member is applied after a regression correction, and to nonexchangeable ensemble members, a new problem arises that has not yet been thoroughly discussed in the literature. Namely, techniques of this form will not produce the minimum-error variance estimate of the mean that could be obtained from the direct application of Bayes's rule; they overweight the climatological information, resulting in suboptimal forecast skill. To be clear, by "overweighting" we are referring to placing the postprocessed ensemble too close to the climatological mean. We will show several examples, from simple models to real data, in which the postprocessed ensemble obtained from a weighted, regression correction on each member is too close to the climatological mean such that the mean of the postprocessed ensemble does not deliver the minimum error variance estimate. This property of overweighting is not to be confused with a desirable (and correct) property of a postprocessed ensemble: as forecast skill decreases the mean of the postprocessed ensemble should converge to the climatological mean. In the following, we use the term overweighting to refer to any deviation from the correct weighting on the climatological mean.

In addition, we present an alternative approach that can be considered a direct application of Bayes's rule. This approach does not rely on an explicit regression, but rather relies on accurately fitting the dataset to the appropriate distributions required by Bayes's rule. Previous work has illustrated a direct application of Bayes's rule in statistical postprocessing. For example, Krzysztofowicz and Evans (2008) introduced the Bayesian processor of forecasts (BPF), whereby climatological data transformed to a normal distribution provides the prior. This prior was updated based on a PDF estimated through a correction of the current forecast using regression relationships estimated from transformed forecast and observational data. Other Bayesian techniques in postprocessing include Rajagopalan et al. (2002) and Luo et al. (2007). Our approach can be considered an application of the BPF of Krzysztofowicz and Evans (2008), in that this work extends the BPF method to the ensemble case and directly accounts for correlation through the chain rule of probability. In addition, we carefully detail a procedure to construct likelihoods based on a function that maps the true state to the forecast.

We also highlight the fact that weighted, KDM methods suffer from an overweighting of climatology and the direct application of Bayes's rule does not. To

test the hypothesis of KDM suboptimality when a regression correction is applied before the KDM fitting procedure, we first choose one of the most popular KDMs as a reference. Henceforth, we focus upon the behavior of BMA, but nevertheless remind the reader that the methods of Glahn et al. (2009) and Unger et al. (2009) should perform similarly. We will first construct a hypothetical scenario of two forecasts that are purposely constructed to lack the correct climatological behavior that random draws from the true distribution would have. In this situation these forecasts should benefit from statistical postprocessing that involves convolving them with climatological information. Because Bayesian methods provide the optimal combination of information, we then examine whether or not the mean of the posterior estimated from BMA is equivalent to the mean estimated from a direct application of Bayes's rule, which provides the minimum error-variance estimate (section 2). Section 3 illustrates in detail how to perform direct Bayesian estimation to be used as a control postprocessing method in the further evaluation of BMA. This method extends the BPF concepts to work with ensemble data, allowing for members to be equally or unequally likely and accounting for possible correlation between members. Section 4 discusses issues related to the size of training dataset; there are challenges with both BMA and Bayesian processing related to the "curse of dimensionality" (Bellman 2003). In section 5, we compare BMA to a Bayesian method through a real-data study, a prediction of ocean wave height based on deterministic global weather forecasts from two operational forecast centers. The Bayesian process in this case will include an additional step, the log transformation of the data; this renders the data more Gaussian before the analytic evaluation through Bayes's rule. Section 6 provides some discussion and conclusions.

2. Differences in the weighting of climatology in Bayes's rule and BMA

In this section we consider the similarity of BMA and related algorithms to a direct application of Bayes's rule. Differences are examined by considering the smallest-possible ensemble, two forecasts, under the simplest of possible assumptions, a forecast with data drawn from a distribution with zero-mean errors. We use this to illustrate the way that BMA will result in an over-weighting of climatological information.

Specifically, imagine that the true state for the physical system under consideration to be drawn from a "climatological" pdf we label, $p(x)$, with the property that x is a random variable drawn from a Gaussian

defined as $\mathcal{N}(\mu, P)$. We have available $N_f = 2$ unbiased forecasts of today's true state, $x = x_t$:

$$\mathbf{x}_f = \begin{bmatrix} x_f^1 \\ x_f^2 \end{bmatrix} = \begin{bmatrix} x_t + \varepsilon_1 \\ x_t + \varepsilon_2 \end{bmatrix}. \quad (2.1)$$

The perturbations ε_1 and $\varepsilon_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$, that is, the errors have zero mean and are Gaussian distributed with a 2×2 covariance matrix \mathbf{R} . The n th diagonal element of \mathbf{R} is denoted as r_n , and the covariance between the errors in the two forecasts will be denoted as ρ . Equation (2.1) implies that $p(\mathbf{x}_f | x)$ is also Gaussian with the statistics of the perturbations $\mathcal{N}(x_t \mathbf{1}, \mathbf{R})$, where $\mathbf{1}$ is the one vector with length N_f . This experimental design is admittedly simpler than what often occurs in actual weather prediction; here forecasts are assumed to have zero bias and no state-dependent error. The simple design here makes the underlying issue with KDMs easy to explain because we have the true posterior available as reference.

a. Bayesian solution

By Bayesian solution we mean in this section that we evaluate, using the correct components of Bayes's rule, the formula for the conditional density. We will hereafter take this expression as the correct answer and evaluate BMA with regards to its ability to recreate this expression.

A straightforward application of Bayes's rule would take the following form:

$$p(x | \mathbf{x}_f) = \frac{p(\mathbf{x}_f | x)p(x)}{p(\mathbf{x}_f)}, \quad (2.2)$$

where $p(\mathbf{x}_f)$ is $\sim \mathcal{N}(\mu \mathbf{1}, P\mathbf{I} + \mathbf{R})$. Deriving the posterior for this simple problem can be done analytically in several ways. In appendix A we show how to minimize the error variance around a weighted linear estimator [as in Eq. (2.3a)], which delivers the mean and variance of the posterior for this problem. One then obtains that $p(x | \mathbf{x}_f)$ is $\mathcal{N}(\bar{x}_{\text{Bayes}}, Q)$ with

$$\bar{x}_{\text{Bayes}} = w_1 x_f^1 + w_2 x_f^2 + w_3 \mu, \quad (2.3a)$$

$$Q = w_3 P. \quad (2.3b)$$

The weights are

$$w_1 = \frac{r_2 - \rho}{r_1 + r_2 - 2\rho + \frac{r_1 r_2 - \rho^2}{P}} \quad (2.4a)$$

$$w_2 = \frac{r_1 - \rho}{r_1 + r_2 - 2\rho + \frac{r_1 r_2 - \rho^2}{P}} \quad (2.4b)$$

$$w_3 = 1 - w_1 - w_2. \quad (2.4c)$$

These weights deliver the minimum-error variance estimate and show how the relative weight between the forecasts and climatology is determined in Eq. (2.3a). It is important to realize that while these weights do sum to one as seen in Eq. (2.4c), it is not correct to interpret these weights as probabilities. In fact, we will show below that these weights can individually be larger than one and can even be negative. Equations (2.3b) and (2.4c) show that the error variance of the posterior mean in Eq. (2.3a) is simply a fraction of the climatological variance P ; the weight applied to P is comparatively small when the sum of weights for the original forecasts is close to one, which happens when forecast-error variances r_1 and/or r_2 are $\ll P$.

b. Bayesian model averaging

Here we derive the BMA weights for this same problem following R05. Equation (2) in that article states that for this experimental design

$$p(x | x_f^1, x_f^2) = m_1 g(x | x_c^1) + m_2 g(x | x_c^2). \quad (2.5)$$

The x_c^n are statistically postprocessed forecasts that are created by separately regressing truth against each forecast to obtain coefficients a and b , such that

$$x_c^n = a_n x_f^n + b_n. \quad (2.6)$$

The kernel associated with the n th regression-corrected forecast is thus

$$g(x | x_c^n) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \frac{(x - x_c^n)^2}{\sigma^2} \right]. \quad (2.7)$$

The weights m_n and the kernel width parameter σ are commonly estimated using the EM algorithm (Dempster et al. 1977). This method can be extended to nonnormally distributed variables [as in the Gamma-distribution kernels for quantitative nonzero precipitation of Slughter et al. (2007); also see the current section 6] or to bivariate distributions [e.g., wind components in Slughter et al. (2010)]. As discussed in Fraley et al. (2010), for the case of ensemble forecasts designed to be exchangeable, the forecast regression coefficients and BMA weights should be constrained to be equal. In the present situation, we omit the refinement of σ through CRPS minimization as discussed by R05 as the value of σ is not the determining factor in the results to be presented below. We do note that if $r_1 \ll r_2$ or vice versa, the assumption of a single value of σ is a further approximation, and ideally the EM procedure should fit different values for each member. We have rerun the experiments of this section with a

modified EM procedure that obtained different σ values for unequally likely forecasts and found that this only affects the weights w_1 and w_2 slightly, and definitely not by enough to explain the difference in BMA and the Bayesian technique.

The regression coefficients in Eq. (2.6) are

$$a_n = \frac{P}{P + r_n}, \quad (2.8a)$$

$$b_n = \frac{r_n}{P + r_n} \mu, \quad (2.8b)$$

assuming a training set of infinite length. These can also be experimentally verified through regression analysis. By inspection, the regressed forecast is simply a weighted linear combination of the original forecast and the climatological mean. We will see next that setting the weight with respect to the climatological mean here in Eqs. (2.8a) and (2.8b) is significant as the relative weighting between climatology and each forecast has now been fixed.

Recall that the BMA mean is simply

$$\bar{x}_{\text{BMA}} = m_1 x_c^1 + m_2 x_c^2. \quad (2.9)$$

This implies that by using Eqs. (2.6), (2.8a), and (2.8b), we may write the BMA mean estimate as in Eq. (2.3a):

$$\bar{x}_{\text{BMA}} = w_1^b x_f^1 + w_2^b x_f^2 + w_3^b \mu, \quad (2.10)$$

where

$$w_1^b = m_1 \frac{P}{P + r_1}, \quad (2.11a)$$

$$w_2^b = m_2 \frac{P}{P + r_2}, \quad (2.11b)$$

$$w_3^b = m_1 \frac{r_1}{P + r_1} + m_2 \frac{r_2}{P + r_2}, \quad (2.11c)$$

$$m_1 + m_2 = 1. \quad (2.11d)$$

The question we wish to answer in this section is whether or not the weights in Eqs. (2.4a)–(2.4c) are equal to Eqs. (2.11a)–(2.11c). In order for Eqs. (2.4a)–(2.4c) to equal Eqs. (2.11a)–(2.11c) the EM algorithm must choose the weights m_1 and m_2 appropriately. Note, however, that this amounts to setting the weights w_1^b , w_2^b , w_3^b equal to their counterparts in Eqs. (2.4a)–(2.4c) by solving for m_1 and m_2 . This turns Eqs. (2.11a)–(2.11d) into four equations in two unknowns (i.e., the system is overdetermined). This suggests that Eqs. (2.11a)–(2.11c) cannot in all circumstances produce the Bayesian weights and therefore the BMA mean, Eq. (2.10), cannot properly produce the minimum error variance estimate.

c. A first look

This section will focus on forecasts with properties that allow one to analytically derive the basic ideas of this manuscript in their simplest form. Here, the forecasts are assumed to be equally likely and have uncorrelated errors. We fully realize that this is not a realistic model of the errors for an ensemble of weather predictions, especially short-lead predictions; in such situations the forecast errors are likely to have some correlation. We do this because this simple model allows us to extract the behavioral characteristics of the algorithms in the most insightful way. Section 2d will examine a more general situation numerically to show the more general aspects of these results.

If one could be absolutely certain that there was no bias in the collection of forecasts in Eq. (2.1) it might seem as though not performing the “bias correction” stage in BMA would be sensible. In fact, this is not true. If one were to perform the BMA algorithm of section 2b without the first-stage bias correction then the resulting mean would simply be the weighted average of the forecasts:

$$\bar{x}_{\text{BMA1}} = m_1 x_f^1 + m_2 x_f^2. \quad (2.12)$$

Note that we would also obtain Eq. (2.12) in the case where we simply performed a “bulk-bias” correction (Wilson et al. 2007) on each forecast because these forecasts have no bulk bias.

For simplicity, let us imagine that these two forecasts are of equal quality (with error variance equal to r) and uncorrelated. In this case, we know the result of the EM algorithm would be that $m_1 = m_2 = 1/2$. Hence, this version of BMA results in the arithmetic mean of the two forecasts and the entire disregard of the information provided by the climatological prior. Note that this result that Eq. (2.12) disregards the climatological prior is true even when the forecasts are correlated and unequally likely.

The error variance of Eq. (2.12) is then

$$Q_{\text{BMA1}} = \frac{r}{2} \quad (2.13)$$

Note, however, that in this same case the Bayesian result in Eq. (2.3a) would be

$$\bar{x}_{\text{Bayes}} = w_1 x_f^1 + w_2 x_f^2 + w_3 \mu, \quad (2.14)$$

where the weights are

$$w_1 = w_2 = \frac{P}{2P + r}, \quad (2.15a)$$

$$w_3 = \frac{r}{2P + r}. \quad (2.15b)$$

The error variance of the Bayesian result in Eq. (2.14) is

$$Q_{\text{Bayes}} = \frac{r}{2P + r} P. \quad (2.16)$$

An examination of Eqs. (2.13) and (2.16) shows that

$$Q_{\text{Bayes}} \leq Q_{\text{BMA1}} \quad (2.17)$$

with equality in the limit as $r \rightarrow 0$ [i.e., when it is reasonable to ignore climatology then Eq. (2.12) is a reasonable choice]. Therefore, the Bayesian result is superior to the BMA result without the first-stage bias correction.

Perhaps a more interesting comparison, however, is to compare Eq. (2.12) to the BMA result where a bias correction is applied even though we have constructed a scenario where there is no bias in the forecasts:

$$\bar{x}_{\text{BMA2}} = w_1^b x_f^1 + w_2^b x_f^2 + w_3^b \mu, \quad (2.18)$$

where

$$w_1^b = w_2^b = \frac{1}{2} \frac{P}{P + r} \quad (2.19a)$$

$$w_3^b = \frac{r}{P + r} \quad (2.19b)$$

and again we used the fact that the result of the EM algorithm will be $m_1 = m_2 = 1/2$. Note that in the limit as $r \rightarrow 0$ that Eq. (2.18) is identical to Eq. (2.12). A direct comparison of Eqs. (2.15b) and (2.19b) reveals that $w_3^b \geq w_3$, which implies an overweighting of climatology.

Using Eqs. (2.19a) and (2.19b) in Eq. (2.18) obtains the following:

$$\bar{x}_{\text{BMA2}} = \frac{P}{P + r} \bar{x}_{\text{BMA1}} + \frac{r}{P + r} \mu. \quad (2.20)$$

Hence, performing the bias correction has led to the BMA mean being a weighted average of Eq. (2.12) and the climatological mean. Note that the error variance of Eqs. (2.18) and (2.20) is

$$Q_{\text{BMA2}} = \left(\frac{P}{P + r} \right)^2 \frac{r}{2} + \left(\frac{r}{P + r} \right)^2 P. \quad (2.21)$$

In Fig. 1 we present the error-variance curves of the three state estimates examined in this section. The most important feature of this figure is that the version of BMA that incorporates the first stage bias-correction step is always better (in the sense of smaller error variance) than the version that does not, even though this set of forecasts is unbiased. The reason it is preferable to “bias correct” in the BMA method even when the

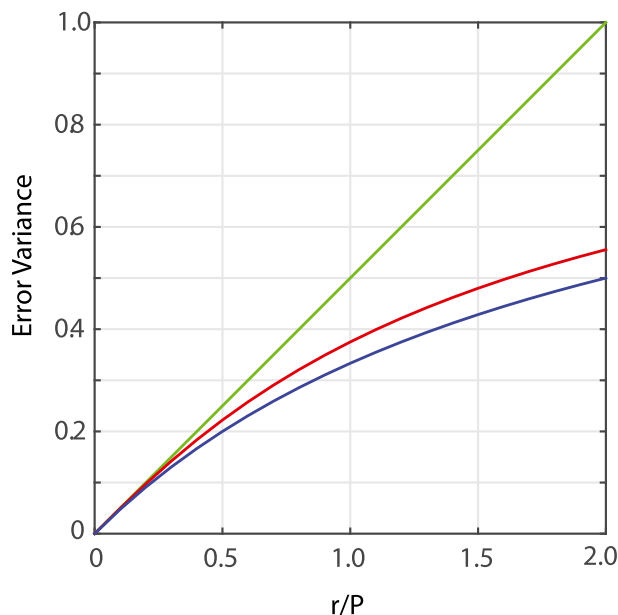


FIG. 1. The error variance for three state estimates: BMA without the first stage “bias correction” (green), BMA with the first stage “bias correction” (red), and the Bayesian weighting (blue).

forecasts are known to be unbiased is because this first stage of BMA is there to include climatological information from the prior. This first-stage correction step in BMA is crucial to getting the climatology into the mean of the estimated PDF, but nevertheless as we have shown the weight on the climatological mean is too great. This overweighting occurs because this two-stage version of BMA is attempting to approximate the multipredictor regression in Eq. (2.14) by first performing two *unipredictor* regressions as in Eqs. (2.8a) and (2.8b) and then weighting according to Eqs. (2.19a) and (2.19b). This approximation of multipredictor regression by *weighted, unipredictor regressions* is suboptimal and results in incorrect weights for the forecasts and climatology.

d. Numerical examination of Bayesian and BMA methods for two synthetic forecasts

We now demonstrate the extent to which BMA overweights climatology when only two forecasts are available. Further, we quantify how much the error in the state estimate is increased by this overweighting of climatology. The analytically derived Bayesian and the BMA EM-estimated weights are shown in Figs. 2 and 3, along with the differences in weights. In the construction of Figs. 2 and 3, a training sample of size 10^5 was used for the regression and EM used in BMA, and simulated forecasts and true values were generated consistent with Eq. (2.1) and the associated text.

Samples of forecasts with correlated errors were generated following the methodology of Houtekamer [1993, see his Eq. (13)].

Figure 2 demonstrates several interesting differences between Bayesian and BMA weights. At high correlations between forecasts (Figs. 2c,f,i), the optimal weighting for the Bayesian forecasts is sometimes negative, while BMA constrains the weights to be positive. Consequently, at high correlations, there are large differences in the weight applied to the first forecast between Bayes and BMA; the differences are smaller with lower correlations between forecasts (Fig. 2g). The corresponding weights applied to climatology and their BMA-Bayesian differences are shown in Fig. 3. For lower correlations between forecast errors, Bayesian and BMA methods provide similar weights to climatology when one or both of the forecasts is relatively accurate (Fig. 3g). For high correlation between the forecasts, there are also similar weights when the forecasts have nearly equal magnitudes of error variances (Fig. 3i). On the other hand, at zero correlation between forecasts, when both forecasts have moderate or large error variances, BMA overweights climatology by a substantial amount, in excess of 10% in many circumstances (Fig. 3g). At very high correlations between forecast errors (Fig. 3i), if one forecast is substantially lower in error than the other, climatology is overweighted.

Does this overweighting increase the errors of the postprocessed forecasts? To explore this, we reduce the training sample size to 100 to add realistic sampling errors, similar in magnitude to training samples in R05, Wilson et al. (2007), and others. We then verify them with another 100 synthetic forecast and truth samples. This training and validation procedure was then repeated 10^3 times to generate a large number of simulated Bayesian and BMA forecasts. Figure 4 provides the root-mean-square error (RMSE) and their differences. Note two particular situations where BMA is notably higher in error than direct Bayesian methods. The first is when there are independent forecasts (Fig. 4g) and moderate to large error variances. The second is when the two forecasts have highly correlated forecast errors, but these errors are significantly different in their variances.

It is interesting that the situation where BMA and the Bayesian result are different depends on ρ when it appears that the weights in Eqs. (2.11a)–(2.11c) are independent of ρ . Obviously, the impact of ρ is somehow implicit in the kernel weights, m_1 and m_2 . However, the result of the EM algorithm for equally likely forecasts (and for any value of ρ) is always $m_1 = m_2 = 1/2$. Therefore, the kernel weights, m_1 and

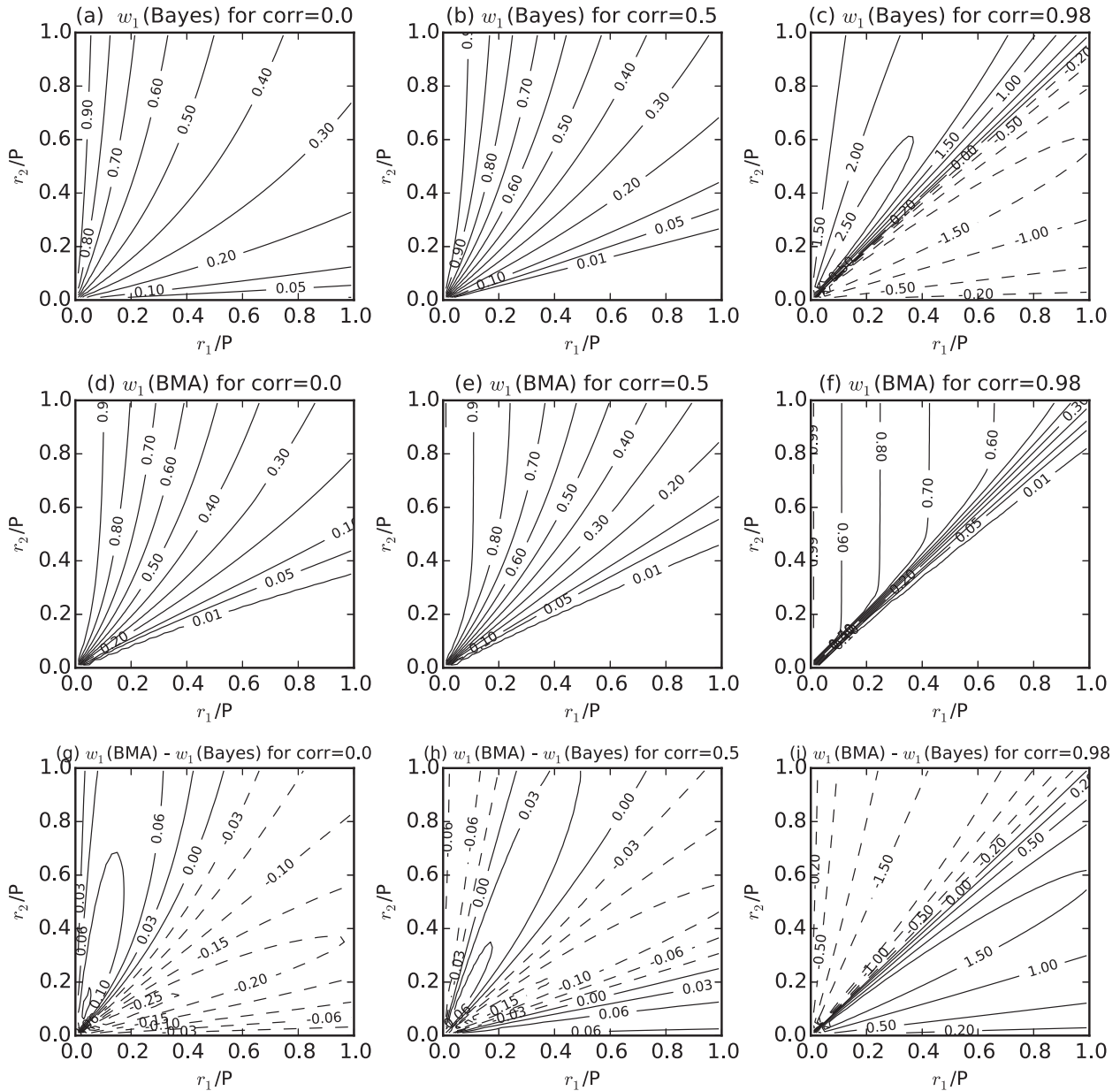


FIG. 2. Weights for the first of two members in a forecast from (a)–(c) Bayes and (d)–(f) BMA. (g)–(i) Weight differences are shown. Weights and differences are plotted as functions of the ratio of the first forecast’s error variance divided by the climatological variance x , the second forecast’s error variance divided by the climatological variance y , and the correlation between the two forecasts (the three columns).

m_2 , are in fact independent of the covariance between forecasts, ρ , in the case of equally likely forecasts. Numerical experiments (not shown) with BMA and the EM algorithm confirm that the kernel weights, m_1 and m_2 , become a function of the covariance between forecasts, ρ , only when the forecasts are not equally likely. Further understanding may be gained by noting that for the equally likely case and a correlation of one that the BMA weights in Eqs. (2.11a)–(2.11c) are

exactly in agreement with the result of regression with only *one* predictor. This is the correct answer in this case and explains why the diagonal of Figs. 2i, 3i, and 4i shows no difference between the BMA and Bayesian result. However, as the covariance ρ decreases the BMA weights in Eqs. (2.11a)–(2.11c) do not change for the equally likely case; therefore, they become more and more in error as the covariance decreases, as can be seen in Figs. 2–4. Further discussion as to how the

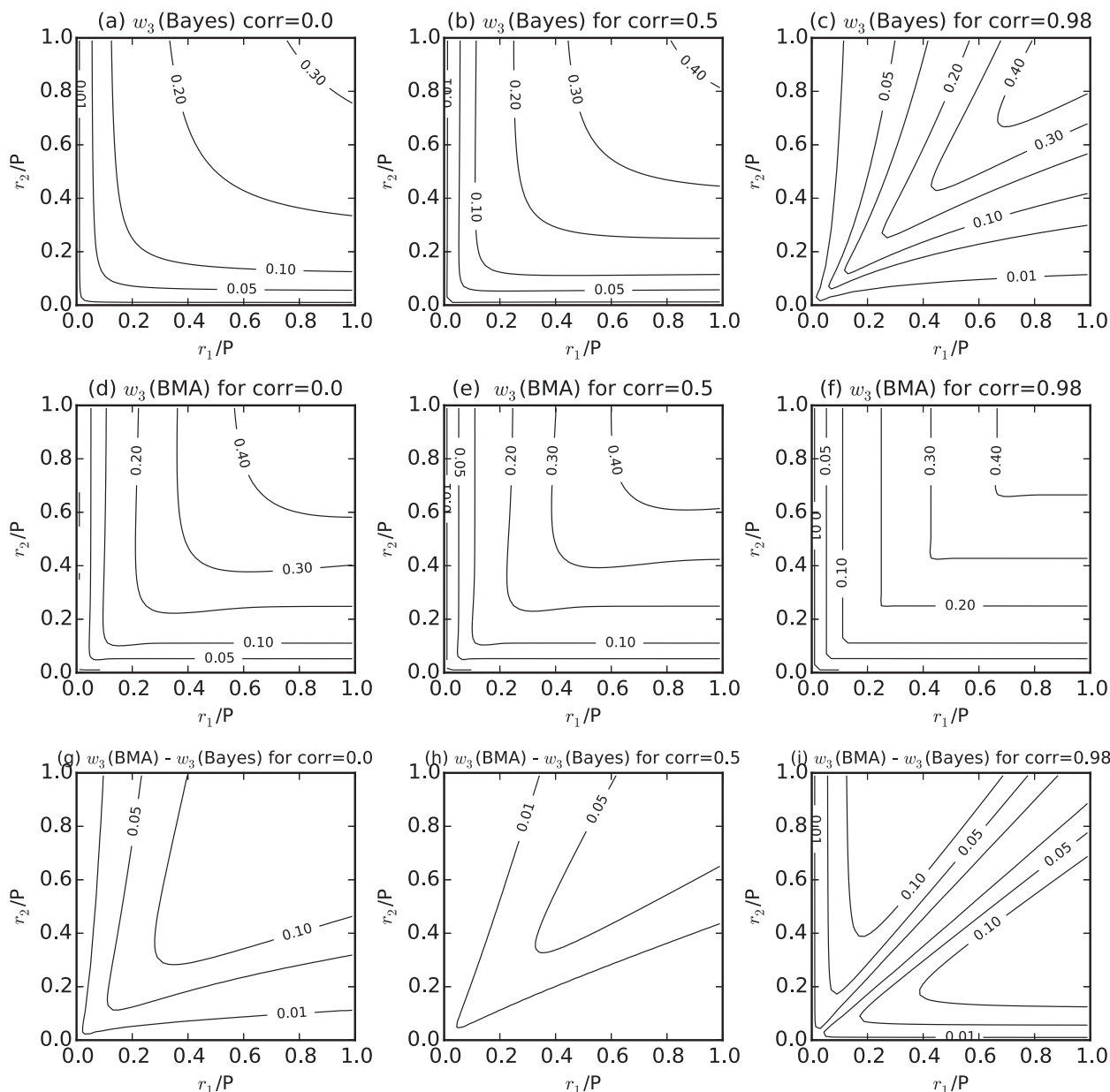


FIG. 3. As in Fig. 2, but here for the weight applied to the climatology. Note that the weight of the second forecast, not shown, is in each case one minus the weight of the first forecast and climatology.

BMA weights differ from the Bayesian weights for numbers of forecasts greater than 2 can be found in [appendix B](#).

3. Direct Bayesian estimation

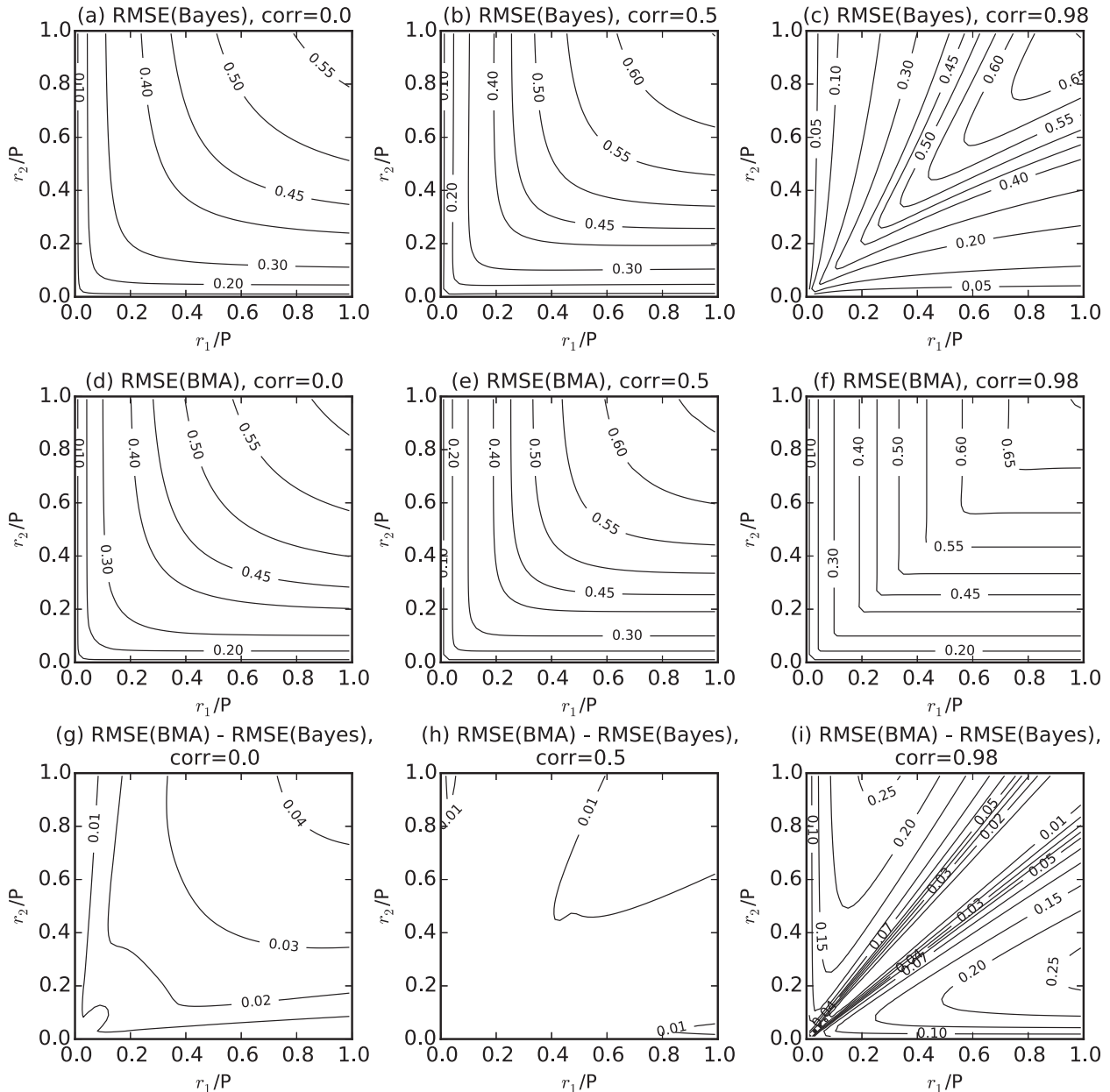
Section 2 showed that by regressing the forecasts to the truth *and then* weighting each with its own kernel was likely to lead to the wrong relative weight on each forecast as compared to the Bayesian result. This section

will illustrate a very simple method that does not begin by regressing the forecasts to the truth and therefore does not suffer from this issue.

We begin by writing Bayes's rule again as

$$p(x | \mathbf{x}_f) = \frac{p(\mathbf{x}_f | x)p(x)}{p(\mathbf{x}_f)}. \quad (3.1)$$

Note that we may use the chain rule of probability to write the forecast likelihood as



$$p(\mathbf{x}_f | x) = p(x_f^1 | x) p(x_f^2 | x, x_f^1) \cdots p(x_f^i | x, x_f^1, x_f^2, \dots, x_f^{i-1}), \quad (3.2)$$

of the characteristics of the dataset. An example of how to choose these functions will be provided in [section 6](#), and we present a method to fit generalized Gaussians in [appendix C](#).

We will assume below that the training set for the likelihoods consists of N_s pairs of forecasts and their verification, which may be the analysis from a data assimilation system or observations. Because Eq. (3.1) expresses the climatological pdf as separate from the forecast likelihoods, we may deduce the climatological

pdf from an archive distinct from our set of truth–forecast pairs. Being separate, this archive may possibly be larger, and hence we refer to the number of members in our climatological archive as N_c .

The method to create this succession of likelihoods we use here proceeds as follows. We begin by scanning the training set for the ordering of the quality of the forecasts. We then sort the forecasts from best to worst, ranked by RMSEs. We start by fitting the best forecast with a Gaussian distribution:

$$p(x_f^1 | x) = C_1 \exp \left\{ -\frac{1}{2} \frac{[x_f^1 - H_1(x)]^2}{R_1} \right\}, \quad (3.3)$$

where C_1 is simply the normalization for the pdf and $H_1(x)$ is the function that results from regression for which the truths are the predictors and the x_f^1 are the predictands. Note that $H_1(x)$ is *not* what is commonly referred to as “state-dependent bias correction”; $H_1(x)$ maps the truth to the forecast, while standard state-dependent bias correction maps the forecast to the truth. This is important, as the likelihood $p(x_f^1 | x)$ must be the distribution of x_f^1 , and if one applied a state-dependent bias correction one would obtain the wrong distribution. This distinction is critical as this aspect is what allows this procedure to create realistic distributions for the likelihoods. The variance R_1 is obtained by calculating the variance across the training set of the difference between the x_f^1 and the regressed truth {i.e. $[x_f^1 - H_1(x)]$ }. Additionally, it is important to recognize that $R_1 \neq r_1$ and that this is true for all forecasts.

This fitting procedure is then extended to the i th forecast:

$$\begin{aligned} p(x_f^i | x, x_f^1, x_f^2, \dots, x_f^{i-1}) \\ = C_i \exp \left\{ -\frac{1}{2} \frac{[x_f^i - H_i(x, x_f^1, x_f^2, \dots, x_f^{i-1})]^2}{R_i} \right\}, \end{aligned} \quad (3.4)$$

where the function $H_i(x, x_f^1, x_f^2, \dots, x_f^{i-1})$ is simply the function that results from multivariate regression for which the predictors are the $x, x_f^1, x_f^2, \dots, x_f^{i-1}$ and the predictand is x_f^i . As before, the error variance R_i is the variance of $[x_f^i - H_i(x, x_f^1, x_f^2, \dots, x_f^{i-1})]$ across the training set.

For completeness, we mention that one could simplify the succession of products in Eq. (3.2) by noting that because they are exponential the products can be rewritten as a sum within the exponential function:

$$p(\mathbf{x}_f | x) = C_1 C_2 \cdots C_{N_f} \exp \left(-\frac{1}{2} S \right), \quad (3.5)$$

$$\begin{aligned} S = & \frac{[x_f^1 - H_1(x)]^2}{R_1} + \frac{[x_f^2 - H_2(x, x_f^1)]^2}{R_2} + \cdots \\ & + \frac{[x_f^{N_f} - H_{N_f}(x, x_f^1, x_f^2, \dots, x_f^{N_f-1})]^2}{R_{N_f}}. \end{aligned} \quad (3.6)$$

After all the forecast likelihoods have been created we must now create the climatological PDF, $p(x)$. There are at least two ways to do this. The first way is to simply fit the climatological distribution to the characteristics of the climatological archive. For example, we may calculate the mean \bar{x} and variance P of the N_c members of the distribution of the truth across our archive. This information allows us to simply fit the climatological pdf as

$$p(x) = C_c \exp \left[-\frac{1}{2} \frac{(x - \bar{x})^2}{P} \right]. \quad (3.7)$$

Once Eq. (3.7) is determined all the information to evaluate Eq. (3.1) has been obtained and therefore the posterior pdf can simply be evaluated for whatever probabilistic prediction is required. Another method is to use kernel density estimation on the climatological distribution and is discussed in [appendix D](#).

The appropriateness of the choice between fitting climatology to a function or representing it with a kernel density method, such as the method described in [appendix D](#), usually comes down to the size of the training set. If the training set is very small (less than 100 samples), then one is likely to do better by simply fitting the climatological distribution to some known function with reasonably accurate characteristics. Given a relatively larger climatological training sample (hundreds or more) one could effectively employ the kernel density algorithm using either Gaussian kernels or Dirac delta kernels (i.e., particle filter). The choice as to which method performs best is likely to be dataset dependent and is left for future work.

Last, we note that the direct application of Bayes’s rule for use in postprocessing is not a new concept. The idea of leveraging prior climatological information and updating with NWP information was presented in the “Bayesian processor of forecasts” of [Krzyzstofowicz and Evans \(2008\)](#). The above procedure explicitly illustrates how to extend to the multimodel ensemble case through the use of the chain rule of probability to compute the forecast likelihood. Further, this procedure directly accounts for correlations between forecasts and is flexible enough to be applied to equally likely ensemble members, subsets of equally likely members, or unequally likely members.

4. Effects of limited training sample size

In this section we will compare BMA to the Bayesian method of [section 3](#) for different numbers of forecasts and different training lengths. The question to be answered in this section is whether or not there is some advantage for small training size to performing regression correction and then weighting each forecast (as in BMA) when compared to the method of [section 3](#).

The problem we set up will be that of [section 3](#), but we will now allow for unequally likely forecasts. As in [section 3](#), we imagine climatology to be drawn from $p(x)$, with the property that x is a random draw from $\mathcal{N}(1, 1)$. We define the forecasts as in [appendix B](#), Eq. (B.1), in which the forecasts will not be correlated here. (The same experiments that will be described below were run for correlations as high as 0.98 and the result that the Bayesian technique had smaller MSE was also found.) Two cases will be run, one with unequally likely forecasts, and the other with equally likely forecasts. In the equally likely forecast case, all forecasts will have $r = 0.5$. In the unequally likely case, the i th forecast will have an error variance of $r_i = 1/10 + (9/10)i/N_f$, where we have chosen this function to bound the forecast quality between 0.1 and 1 for any number of forecasts. The BMA algorithm will be implemented as described in [section 3](#), with a convergence criteria measured as the change in the weights being less than 10^{-5} or a maximum number of iterations of 50.

After these calculations have been performed, we calculate the mean of the predicted distribution and then calculate the RMSE with respect to the truth for each technique. We will repeat this entire calculation for 10^5 verification trials for each combination of number of forecasts and training size. For each of these trials, a different truth is drawn from climatology and truth samples are created according to Eq. (2.1). The resulting average over the 10^5 verification trials is reported as the difference between the RMSE of the above-mentioned Bayesian technique and BMA ([Fig. 5](#)). The red line in the upper-left corner of the figure denotes the region for which the number of forecasts is equal to or greater than the training size. Experiments in this region were not performed, as we do not believe that either technique can be expected to deliver sensible results when the number of forecasts (predictors) is greater than the training size. Positive values in [Fig. 5](#) indicate that BMA had larger RMSE averaged over the 10^5 trials than the Bayesian technique. By scanning both panels of [Fig. 5](#), one can see that the Bayesian method has smaller RMSE over a wide range of training sizes and numbers of forecasts. In fact, only one experiment resulted in a very weak negative result, and that is for the equally

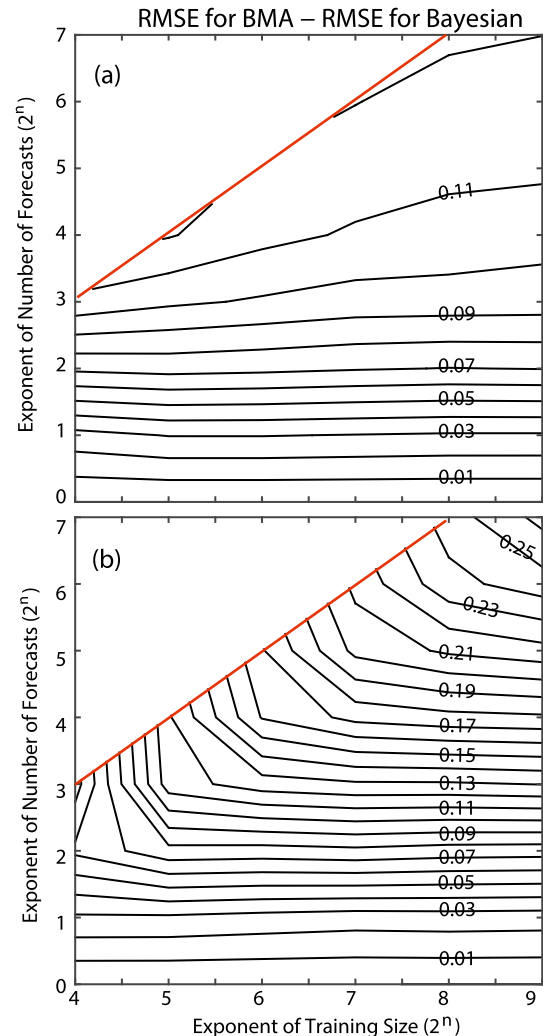


FIG. 5. MSE difference between BMA and the Bayesian particle-filter-inspired method of [section 4](#). Results from the experiment with (a) unequally likely forecasts and (b) equally likely forecasts. The diagonal red line in the upper-left corner of each plot denotes the region above which the training size is smaller than the number of predictors, and for which experiments were not performed.

likely forecast case and a training size of 16 and the number of forecasts equal to 16. Therefore, [Fig. 5](#) shows that this difference between the Bayesian result and BMA's overweighting of climatology is robust in the presence of sampling error from limited training size.

5. An application to wave forecast postprocessing

In this section we apply the previously developed ideas to multimodel forecasts of ocean significant wave heights (units will be in meters) in the North Pacific Ocean. All forecasts are obtained using the Wavewatch III ([Tolman 1997, 1999, 2001](#)) global ocean wave model

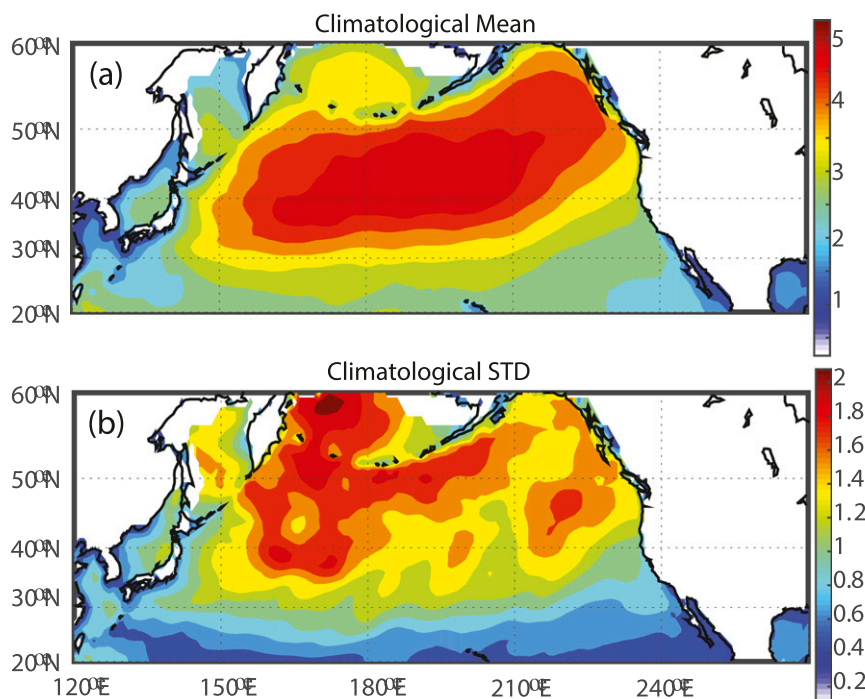


FIG. 6. The (a) climatological mean and (b) standard deviation across the training set at each grid point for the period 3 Dec 2011–31 Jan 2012.

implemented at two U.S. forecasting centers: Fleet Numerical Meteorology and Oceanography Center (FNMOC) and the National Centers for Environmental Prediction (NCEP). Wind data are obtained from the Navy Operational Global Atmospheric Prediction System (NOGAPS) (Rosmond, 1992) model and the Global Forecast System (GFS; Han and Pan 2011), respectively.

We will use two forecasts as in section 2, one from each implementation of Wavewatch III. In the following we consider the +120-h forecast lead time and obtain data at $1^\circ \times 1^\circ$ resolution. The verification dataset will be taken as the FNMOC Wavewatch III significant wave height analysis. We will postprocess each grid point of these fields that contain a wave height forecast from both models and for grid points between 20° and 60° N latitudes and between 120° and 270° E longitudes. The training set for this experiment will be a 60-day running period immediately preceding each day to be postprocessed. We will apply this training set to postprocess the forecasts for 1–28 February 2012. For example, to postprocess the forecasts for 1 February we use a training period from 3 December 2011 to 31 January 2012. The climatological mean and variance for each of the points in our postprocessing region for the 1 February training set is plotted in Fig. 6.

Wave height forecasts have the property that they are positive definite, which leads in these forecasts to highly

skewed distributions with long tails. Clearly then both BMA and the Bayesian techniques must be constructed to account for this fact. For BMA we will use a power transformation (following Yeo and Johnson 2000) such that a variable u is mapped to “log” space through the following transform:

$$u \rightarrow \log(1 + u). \quad (5.1)$$

This transformation into log space effectively pulls the tails of the distribution inward toward the center of the distribution, which results in a training dataset in log space that is more accurately fit to a Gaussian. We apply this to all data required by the algorithm, perform the BMA algorithm, and then map a variable v back to physical space using

$$v \rightarrow \exp(v) - 1. \quad (5.2)$$

The equivalent BMA kernels in Eq. (2.7) are of the following form:

$$g(x | x_c^n) = \frac{C_n}{1+x} \exp \left[-\frac{1}{2} \frac{(Lx - Lx_c^n)^2}{\sigma^2} \right], \quad (5.3)$$

where $Lx = \log(1 + x)$ and Lx_c^n is the n th bias corrected forecast in log space. We will use this representation to plot the predicted distributions from BMA.

For comparison we will also implement the direct Bayesian method of [section 3](#). Because the training set for climatology only has 60 samples, we simply fit the climatological distribution to a function. To begin, we must also account for the positive definiteness and long tails of the wave height distributions in both the forecast likelihoods and in the climatological distribution. For the forecast likelihoods we first transform the data as in [Eq. \(5.1\)](#) and then apply the equations of [section 3](#) to the transformed data. To map back we simply use the Jacobian of the transformation to obtain forecasts likelihoods of the following form:

$$p(x_f^1 | x) = \frac{C_1}{1 + x_f^1} \exp \left\{ -\frac{1}{2} \frac{[Lx_f^1 - H_1(Lx)]^2}{R_1} \right\}, \quad (5.4)$$

where $Lx_f^1 = \log(1 + x_f^1)$. It is important to note that the regression for the function H_1 is done with both the predictor and the predictand in log space. We choose here to define the function H_1 , and H_2 below, as the function that results from linear regression, but note that if the training set was larger we would prefer to use quadratic or even cubic polynomial regression. Similarly, if the training set were larger, we could use a generalized Gaussian (which would fit the first, second, and fourth moments, rather than just the first and second; please see [appendix C](#)) and would allow for a better prediction of the higher moments of the Bayesian posterior. This generalized Gaussian was applied to this dataset but no statistically distinguishable improvement was found in the RMS of the mean and we believe that this is due to the training set being too small to make use of the information about the tails of the distribution implied in the generalized Gaussian.

The second forecast's likelihood is obtained from

$$p(x_f^2 | x, x_f^1) = \frac{C_2}{1 + x_f^2} \exp \left\{ -\frac{1}{2} \frac{[Lx_f^2 - H_2(Lx, Lx_f^1)]^2}{R_2} \right\}, \quad (5.5)$$

where $Lx_f^2 = \log(1 + x_f^2)$.

Finally, we must construct the climatological distribution. We do this in the same way as the forecast likelihoods by transforming the wave data into log space, fitting a Gaussian, and transforming back using the Jacobian to obtain the following:

$$p(x) = \frac{C_c}{1 + x} \exp \left[-\frac{1}{2} \frac{(Lx - \overline{Lx})^2}{P} \right], \quad (5.6)$$

where we again emphasize that \overline{Lx} and P are the climatological mean and variance in log-transformed space.

An example of the resulting functions defined above is plotted in [Fig. 7](#) for the grid point at 50°N, 180°. At this grid point, the GFS forecasts happened to have less RMSE over the training period and were assigned as f_1 . Subsequently, the NOGAPS forecasts were assigned the variable f_2 . In [Fig. 7a](#) we plot the forecast likelihood describing the distribution of GFS forecasts given the verification [e.g., [Eq. \(5.4\)](#)]. By comparing the shape of this distribution to the one-to-one line, we can see that the GFS forecast typically overforecasts when the true wave height is small and underforecasts when the true wave height is large. Additionally, note that the width of the distribution increases as the true state being conditioned upon increases. This states that the variance in the forecast increases as the true state being conditioned upon increases and is a common property of positive definite distributions.

[Figures 7b and 7c](#) provide the forecast likelihood for the NOGAPS forecasts given the verification and the GFS forecasts [e.g., [Eq. \(5.5\)](#)]. In [Fig. 7b](#), we plot f_2 as a function of the verification and evaluated for f_1 equal to the GFS forecast on 1 February. In [Fig. 7b](#), we see that the variability in the NOGAPS wave height forecast is significantly larger than that of the GFS forecast, though it too typically overforecasts when the true wave height is small and underforecasts when the true wave height is large. We believe this larger variance in the NOGAPS forecast is due to the fact that it is being run at a lower resolution than the GFS model and, therefore, cannot develop strong-enough winds to force the wave model to large wave heights. In [Fig. 7c](#), we plot f_2 as a function of f_1 and evaluated for x equal to the verification at 0000 UTC 1 February 2012. This plot shows the implied correlation between f_2 and f_1 as seen by the training set, and shows that there is generally a weak correlation between them. In summary, these types of plots of the forecast likelihoods provide a more truly Bayesian way to perform forecast validation by showing the relationship between the forecasts and the truth and the forecasts between each other.

Evaluating Bayes's rule in [Eq. \(3.1\)](#) requires a probability density estimation of the climatological distribution. In [Fig. 7d](#) we show the PDF for the climatological distribution [e.g., [Eq. \(5.6\)](#)]. The histogram in [Fig. 7d](#) is composed of the 60 samples at this grid point and shows that the lognormal fit is a reasonable choice for such a small number of samples. In [Fig. 7e](#), we show the resulting Bayesian posterior distribution evaluated for f_2 equal to the GFS forecast on 1 February. Finally, in [Fig. 7f](#) we show the resulting BMA posterior, which we obtained after using [Eq. \(5.3\)](#) in [Eq. \(2.5\)](#). This example was chosen because it illustrates one of the ways that BMA tends to overweight climatology. By

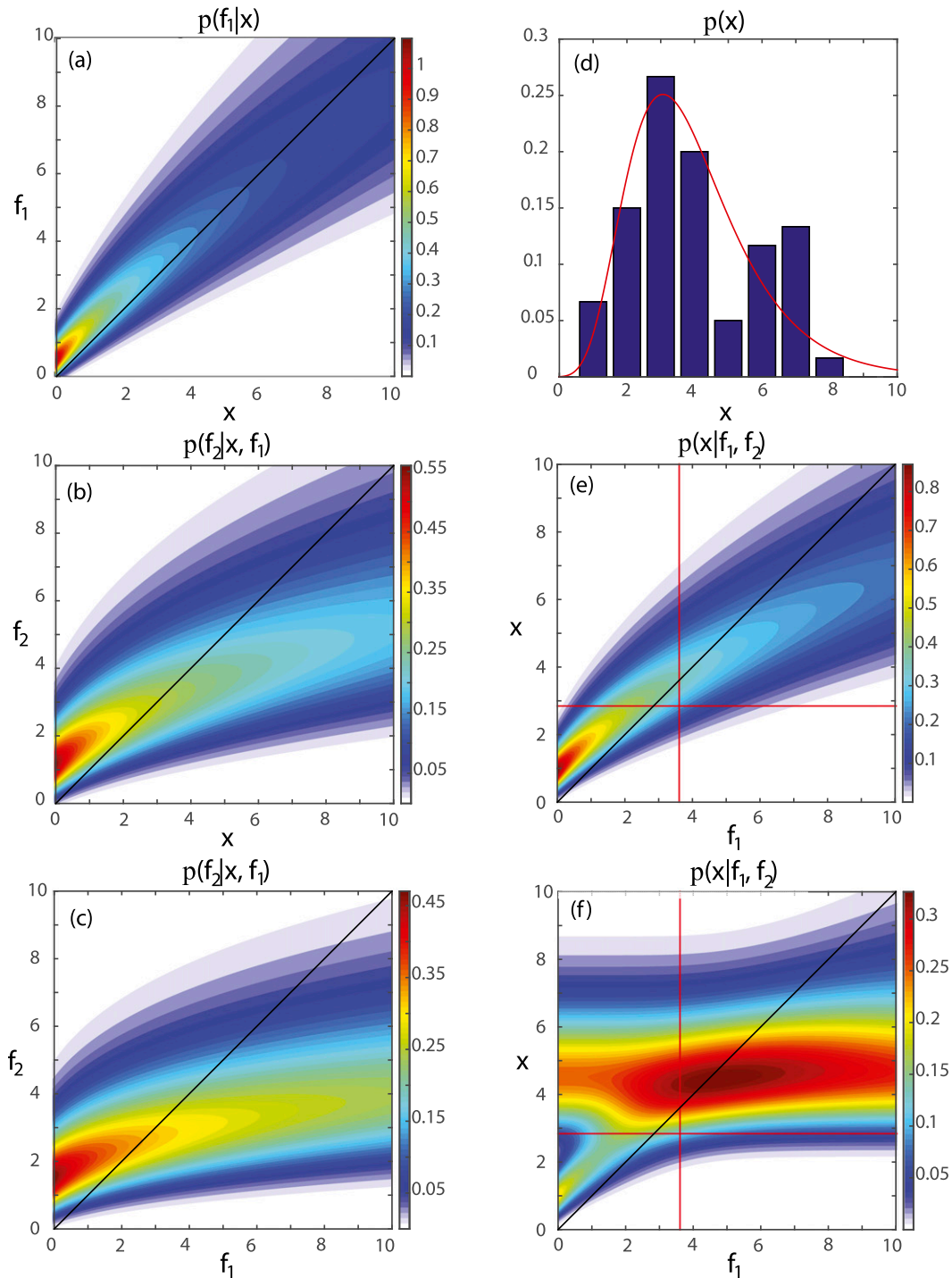


FIG. 7. The fitting of the likelihoods and climatology for a grid point in the North Pacific (50°N, 180°) and evaluated for 1 Feb 2012. (a) The forecast likelihood for the GFS model given the truth. (b),(c) The forecast likelihood for the NOGAPS model given the truth and the GFS forecast. (d) The climatological histogram as well as the fit to the lognormal distribution. The posterior distributions for the (e) Bayesian technique and (f) BMA. The black diagonal line in some panels is the one-to-one line. The red vertical (horizontal) line in (e) and (f) is the 1 Feb GFS forecast (truth).

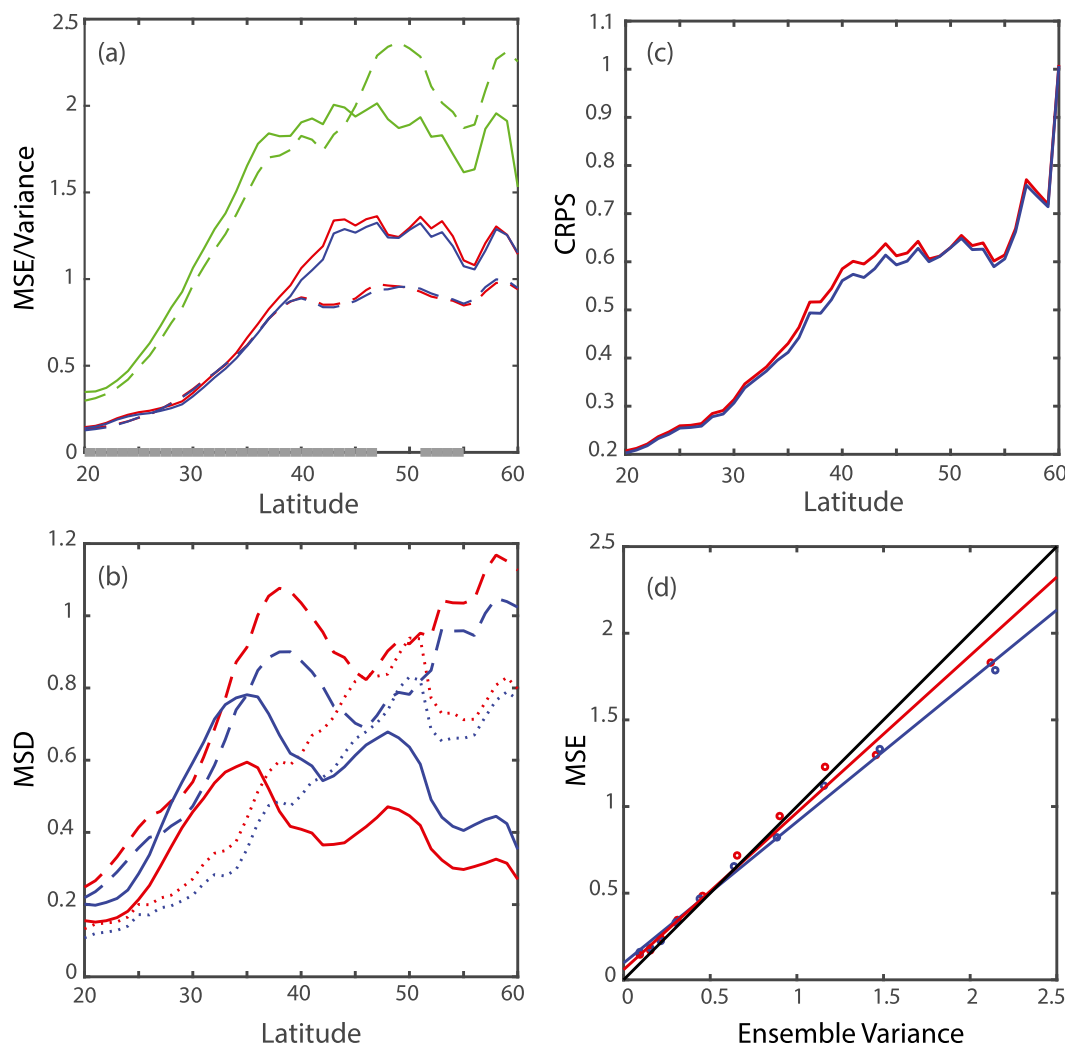


FIG. 8. Results for the postprocessing of days 1–28 Feb 2012. (a) The MSE with respect to the NOGAPS analysis as a function of latitude. Solid blue (red) is the MSE of the Bayesian technique (BMA). The green solid line is the MSE of the climatological mean. The dashed lines correspond with the same color choice for the Bayesian technique, BMA, and climatology except are the corresponding variances. The gray bar along the bottom of (a) denotes those latitudes for which the difference between the MSE of the Bayesian technique and BMA is statistically significantly different at 99% using a t test. (b) The MSD between the Bayesian technique (BMA) in blue (red). The solid line is for the posterior mean difference with respect to the climatological mean; the long dashed line is the posterior mean difference against the GFS forecast; and the short dashed line is the posterior mean difference against the NOGAPS forecast. (c) The binned spread plot for the Bayesian technique (blue) and BMA (red) using 10 equally populated bins. (d) The linear fit through all the data is also plotted as the red (BMA) and blue (Bayesian) lines.

comparing Figs. 7e and 7f, one can see that the BMA posterior is generally significantly wider and has the bulk of its probability mass fixed near the climatological mean of approximately 4 m. Additionally, one can see that the two kernels in the BMA estimation procedure become bimodal below a ~ 2 -m wave height, which we believe is unwarranted.

Further evidence that the BMA procedure is over-weighting climatology is presented in Fig. 8. In Fig. 8a, we show the mean squared error (MSE) with respect to

the verification for the 28 days of the verification period and averaged as a function of latitude. Additionally, we show the MSE for the same period for the climatological mean. Note that the direct Bayesian estimation of section 4 has a lower MSE than BMA for all latitudes except near 60°N. Note that near 60°N the climatological variance over the training period is much larger than the MSE of the climatological mean during the verification period. Hence, a procedure that erroneously over-weights climatology will appear to be accurate in the

situation where the climatological variance is larger than it should be.

To relate this example to [section 2](#), recall that a comparison of Eqs. (2.3a) and (2.10) resulted in a stronger weight on the climatological mean for BMA. In [Fig. 8b](#), we show the mean squared difference (MSD) between the posterior mean from the procedure in [section 3](#) and its three components (climatological mean, NOGAPS, and GFS forecasts). We repeat these calculations for the MSD between the posterior mean from BMA and its three components (climatological mean, NOGAPS, and GFS forecasts). The point here is that when the MSD is small this implies that the weight on that component of the posterior mean estimate must be large. Similarly, when the MSD is large for a particular component of the posterior mean estimate then the weight on that component must be small. [Figure 8b](#) shows that the posterior mean from BMA has a smaller MSD to the climatological mean than the Bayesian procedure of [section 4](#). Similarly, [Fig. 8b](#) shows that the posterior mean from BMA has a larger MSD to the two forecasts than the Bayesian procedure. We feel that this shows that the theoretical results presented in [section 2](#) can be directly seen in a real application of multimodel postprocessing.

Last, we show in [Figs. 8c and 8d](#) the binned-spread diagram and the CRPS as a function of latitude for both techniques. The binned-spread diagram is created from all the data points across the Pacific and over the 28 days of verification. We believe that these figures show that both techniques have a reasonable relationship between the width of the posterior and their squared errors. Technically there appears to be some benefit seen in the CRPS for the direct Bayesian estimation method, but we caution that this may be due solely to the better posterior mean rather than the shape of the PDF. The impact of the overweighting of climatology on the variance prediction may be too small to detect in the moments higher than the first in this example with such a small training set. Further study of the impact of the size of the training set on the quality of the higher moments will be presented in a sequel.

6. Summary and conclusions

We have shown postprocessing techniques like BMA that first apply a regression correction and then apply a weighting to ensemble members will systematically overweight climatology. This problematic treatment of climatological, or potentially other non-NWP information, results in an increase in the mean-squared error of the resulting state estimate. We demonstrated that this result held in various parameter spaces,

including unequally likely and correlated ensemble forecasts. We also showed that this result was independent of the ensemble size and the size of the training set.

We note that the issue of overweighting climatology does not arise because of the kernel-density estimation assumption, but rather because the operation of regression correcting first prematurely fixes the relative weight to climatology. We present an alternative approach based on a direct Bayesian estimation. This approach does not rely on an explicit regression-correction step, but rather relies on accurately fitting the dataset to the appropriate distributions required by Bayes's rule. Our approach can be considered an application of the BPF of [Krzysztofowicz and Evans \(2008\)](#), in that this work extends the BPF method to the ensemble case and directly accounts for correlation through the chain rule of probability. In addition, we carefully detail a procedure to construct likelihoods based on a function that maps the true state to the forecast. We also present a particle-filter-based formulation of this procedure. We demonstrate the ability of this method to extend to non-Gaussian, distributions through a log transformation, by application to multimodel ocean significant wave heights.

We feel that the most powerful application of the Bayesian technique presented here will be obtained for larger training sizes than we had available. The study of how to increase the training set size as well as what size training set is required to employ higher-order function fitting procedures will be the subject of future work.

Acknowledgments. This research is supported by the Chief of Naval Research through the NRL Base Program, PE 0601153N. We thank three anonymous reviewers for very conscientious critiques that have helped us greatly improve the manuscript.

APPENDIX A

Derivation of Eqs. (2.3a) and (2.3b)

To calculate Eqs. (2.3a) and (2.3b) we first must find the error variance as a function of the weights:

$$Q = \langle (x_t - \bar{x}_{\text{Bayes}})^2 \rangle = (1 - w_1 - w_2)^2 P + w_1^2 r_1 + 2w_1 w_2 \rho + w_2^2 r_2 + (1 - w_1 - w_2 - w_3)^2 \mu^2. \quad (\text{A.1})$$

The mean of the posterior distribution will minimize the variance in Eq. (A.1). To this end we differentiate Eq. (A.1) with respect to each weight:

$$\frac{dQ}{dw_1} = -2(1 - w_1 - w_2)P + 2w_1r_1 + 2w_2\rho - 2(1 - w_1 - w_2 - w_3)\mu^2, \quad (\text{A.2a})$$

$$\frac{dQ}{dw_2} = -2(1 - w_1 - w_2)P + 2w_2r_2 + 2w_1\rho - 2(1 - w_1 - w_2 - w_3)\mu^2, \quad (\text{A.2b})$$

$$\frac{dQ}{dw_3} = -2(1 - w_1 - w_2 - w_3)\mu^2. \quad (\text{A.2c})$$

Setting Eqs. (A.2a)–(A.2c) to zero and solving obtains the weights in Eqs. (2.4a)–(2.4c). Insertion of the weights in Eqs. (2.4a)–(2.4c) into Eq. (A.1) returns Eq. (2.3b).

APPENDIX B

The Potential to Overweight Climatology with Many Ensemble Members

Consider a situation where the ensemble forecasts were constructed a priori to have no systematic errors, but a forecaster could not be sure of this property. Generally, one would expect a more accurate resulting pdf from a KDM with a larger ensemble (with associated narrower kernels) than with a smaller ensemble (with wider kernels). Is it possible then that the issues with BMA overweighting climatology demonstrated in section 2 were due to the use of a small ensemble (of size 2)? In this section we demonstrate that the overweighting of climatology can still occur with large ensembles. This overweighting is illustrated here with a simple experimental design, one for which we can analytically derive the correct Bayesian result in the limit of very large numbers of forecasts. Here, the forecasts are assumed to be equally likely and have uncorrelated errors. Again, we fully realize that this is not a realistic model of the errors for an ensemble of weather predictions, especially short-lead predictions, where the forecast errors are likely to have some correlation. This experimental design is rather meant to illustrate the mathematical property that KDMs with a prior regression stage do not necessarily converge to the correct predicted mean as the number of forecasts approaches infinity.

As in section 2, imagine the climatology for the physical system under consideration to be drawn from $p(x)$, with the property that x is a random draw from $\mathcal{N}(\mu, P)$. We have available N_f unbiased forecasts of today's true state, $x = x_t$:

$$\mathbf{x}_f = x_t \mathbf{1} + \boldsymbol{\varepsilon}, \quad (\text{B.1})$$

where $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$. Here $\mathbf{1}$ denotes an $N_f \times 1$ vector whose entries are 1. The perturbations $\boldsymbol{\varepsilon}$ have the

property that \mathbf{R} is $N_f \times N_f$. The diagonal of \mathbf{R} will be constant and denoted as r ; off-diagonal elements of \mathbf{R} are zero. Equation (3.1) thus implies that $p(\mathbf{x}_f | x)$ is Gaussian with the statistics of the perturbations [i.e., $\mathcal{N}(\mathbf{0}, \mathbf{R})$].

a. Bayesian solution

A straightforward application of Bayes's rule would take the form of Eq. (2.2). By using the information denoted above in Eq. (B.1) we obtain that $p(x | \mathbf{x}_f)$ is $\mathcal{N}(\bar{x}_{\text{Bayes}}, Q)$ with

$$\bar{x}_{\text{Bayes}} = \mu + P\mathbf{1}^T[\mathbf{1}P\mathbf{1}^T + \mathbf{R}]^{-1}(\mathbf{x}_f - \mu\mathbf{1}), \quad (\text{B.2a})$$

$$Q = (\mathbf{1} - P\mathbf{1}^T[\mathbf{1}P\mathbf{1}^T + \mathbf{R}]^{-1}\mathbf{1})P. \quad (\text{B.2b})$$

The Sherman–Morrison–Woodbury formula (Golub and Van Loan 1989, see their section 2.1.3) states that

$$[\mathbf{1}P\mathbf{1}^T + \mathbf{R}]^{-1} = \mathbf{R}^{-1} - \frac{P}{r^2 + rPN_f} \mathbf{1}\mathbf{1}^T, \quad (\text{B.3})$$

and when this is used in Eqs. (B.2a) and (B.2b), we obtain

$$\bar{x}_{\text{Bayes}} = \frac{r}{r + PN_f} \mu + \frac{P}{r + PN_f} \mathbf{1}^T \mathbf{x}_f, \quad (\text{B.4a})$$

$$Q = \left(1 - \frac{PN_f}{r + PN_f}\right)P. \quad (\text{B.4b})$$

Note that the vector product in Eq. (B.4a) is the sum of the forecasts:

$$\mathbf{1}^T \mathbf{x}_f = \sum_{n=1}^{N_f} x_f^n, \quad (\text{B.5})$$

and that Eq. (B.1) implies, in the sense of probability, that

$$\lim_{N_f \rightarrow \infty} \frac{1}{N_f} \sum_{n=1}^{N_f} x_f^n = x_t. \quad (\text{B.6})$$

Using Eqs. (B.5) and (B.6) in Eqs. (B.4a) and (B.4b) reveals that the Bayesian solution has the property that the posterior mean approaches the truth, and the posterior variance vanishes, as $N_f \rightarrow \infty$.

b. Bayesian model averaging

Recall with BMA that the PDF is constructed with weighted kernels. Following R05, their Eq. (2):

$$p(x | \mathbf{x}_f) = \sum_{n=1}^{N_f} w_n g_n(x | x_c^n). \quad (\text{B.7})$$

Here the x_c^n are the statistically postprocessed forecasts. Note that Eqs. (2.8a) and (2.8b), when applied to the model of Eq. (B.1), implies that

$$a_n = \frac{P}{P+r}, \quad (\text{B.8a})$$

$$b_n = \frac{r}{P+r} \mu, \quad (\text{B.8b})$$

when provided with a training set of infinite length; the coefficients of Eqs. (B.8a) and (B.8b) are identical for all forecasts because the forecasts are equally likely. This is precisely the reason we have chosen equally likely forecasts, as this leads to a case in which we can predict precisely the regression coefficients and the result of the EM algorithm (e.g., Fraley et al. 2010). In this case we know that the weights are all equal to $1/N_f$.

Putting these results together allows one to show that the mean from BMA is

$$\bar{x}_{\text{BMA}} = \sum_{n=1}^{N_f} w_n x_c^n = \frac{r}{P+r} \mu + \frac{P}{P+r} \frac{1}{N_f} \sum_{n=1}^{N_f} x_f^n. \quad (\text{B.9})$$

Making use of Eq. (B.6) shows that the BMA mean in the limit of an infinite number of members in Eq. (B.7) equals

$$\bar{x}_{\text{BMA}} = \frac{r}{P+r} \mu + \frac{P}{P+r} x_t. \quad (\text{B.10})$$

Therefore, the BMA mean in the limit of an infinite number of members is a weighted average between the climatological mean and the truth, while the Bayesian solution [Eqs. (B.4a) and (B.4b)] is simply the truth in that limit.

To the extent that this experimental design resembles possible ensembles, it shows that BMA can substantially overweight climatology. Note, however, that in this experimental design, had the regression analysis been applied to \bar{x} instead of to the individual members, it can be shown that the posterior mean would be identical to that from Bayes's rule. Therefore, the example of this section is simply meant to illustrate the behavior in the asymptotic limit of a large number of kernels. Recall that the typical result in KDMs is that the estimated PDF from KDMs become more accurate as the number of kernels increases. By contrast, however, we have shown that increasing the number of kernels (forecasts) does not eliminate this issue of overweighting climatology because this issue occurs in BMA because of the initial regression correction fixing the relative weight between the climatological mean and the forecasts and, as shown in section 2, this occurs more generally than just the specific experimental design of this section.

APPENDIX C

Fitting to a Generalized Gaussian

We fit a variable u to a generalized Gaussian by first calculating its mean \bar{u} , variance θ , and fourth moment F . The generalized Gaussian is defined as

$$p(u) = N \exp \left[- \left(\frac{|u - \bar{u}|}{\sqrt{2}\sqrt{\beta}} \right)^\alpha \right] \quad (\text{C.1})$$

where N is simply the normalization and the parameters α and β are to be fit to the variance and fourth moment. We do this using an iterative procedure in which we iterate for α_i by using the previous α_{i-1} in

$$g_i = \frac{\Gamma \left(\frac{3}{\alpha_{i-1}} \right)^2}{\Gamma \left(\frac{5 - \alpha_{i-1}}{\alpha_{i-1}} \right) \Gamma \left(\frac{1}{\alpha_{i-1}} \right)} \frac{F}{\theta^2} \quad (\text{C.2})$$

to obtain the next update:

$$\alpha_i = \frac{5}{1 + g_i}. \quad (\text{C.3})$$

Typically, fewer than 20 iterations are required for convergence. Once α is known β may be calculated from

$$\beta = \frac{\Gamma \left(\frac{1}{\alpha} \right)}{2\Gamma \left(\frac{3}{\alpha} \right)} \theta. \quad (\text{C.4})$$

Note that if the distribution of u is a standard Gaussian such that $F = 3\theta^2$ then $\alpha = 2$ and $\beta = \theta$, which reduces Eq. (C.1) to a standard Gaussian.

APPENDIX D

Kernel Density Methods for the Climatology

An alternative to fitting climatology to a prespecified function is to use a kernel density estimation procedure to represent the climatological pdf. The simplest way to do this is to use the framework of particle filtering (Doucet et al. 2000). In the particle filtering framework we assume that our samples from climatology are equally likely and have kernels that have zero width (i.e., they are the Dirac delta function rather than a Gaussian, as in BMA). We emphasize, however, that there is no requirement that the kernels have zero width; they could be Gaussians as in the equations for BMA [e.g., Eq. (2.7)]. The key defining difference between this use of a

KDM and that of BMA is that BMA is applying kernel density estimation to estimate the posterior directly, while under this technique kernel density estimation is used to represent the prior (climatological) distribution. This distinction is substantial as the number of kernels in BMA is equal to N_f , but the number of kernels in this technique is N_c , and typically we have training sets for which $N_c \gg N_f$. Because KDMs are typically most accurate for large numbers of kernels we view this aspect as beneficial.

As an example, we assume kernels of zero width (Dirac delta functions) to represent the climatological distribution. This allows one to determine the probability that the j th sample from the climatological distribution, x_j , is the true state given the forecasts:

$$p_j = \frac{p(\mathbf{x}_f | x_j)}{\sum_{j=1}^{N_c} p(\mathbf{x}_f | x_j)}, \quad (\text{D.1})$$

where N_c is the number of states from climatology that we have available in our training set. Hence, Eq. (D.1) is essentially an “analog” approach in which we search the training set for states from climatology that, as measured by our forecast likelihoods, are likely to be today’s truth. Note that the probability that the j th sample from the climatological distribution is the true state given the forecasts is also the weight for each of our delta function kernels. Hence, we may calculate the mean of the posterior as

$$\bar{x}_b = \sum_{j=1}^{N_c} p_j x_j. \quad (\text{D.2})$$

Other moments of the posterior may be calculated similarly.

In some applications, an ensemble of equally likely members proves useful. To obtain this from Eq. (D.1), we sort the x_j from smallest to largest and label this new set as $x_{(j)}$. We then reorganize the p_j into the same ordering as the $x_{(j)}$ to obtain $p_{(j)}$. The $p_{(j)}$ can now be cumulatively summed to determine the cumulative distribution function (CDF). The procedure for sampling from a CDF is well known and goes as follows. Draw a uniformly distributed random number on 0–1. Next, find the element of the CDF closest to this number. This element of the CDF corresponds to a particular $x_{(j)}$ and, therefore, this value is the correct random draw from the posterior, Eq. (3.1). Repeat this procedure any number of times to obtain the ensemble of equally likely members.

REFERENCES

- Bellman, R. E., 2003: *Dynamic Programming*. Dover Publications, 384 pp.
- Bishop, C. H., and K. T. Shanley, 2008: Bayesian Model Averaging’s problematic treatment of extreme weather and a paradigm shift that fixes it. *Mon. Wea. Rev.*, **136**, 4641–4652, doi:10.1175/2008MWR2565.1.
- Delle Monache, L., F. A. Eckel, D. L. Rife, B. Nagarajan, and K. Searight, 2013: Probabilistic weather prediction with an analog ensemble. *Mon. Wea. Rev.*, **141**, 3498–3516, doi:10.1175/MWR-D-12-00281.1.
- Dempster, A., N. Laird, and D. Rubin, 1977: Maximum likelihood from incomplete data via the EM algorithm (with discussion). *J. Roy. Stat. Soc. Stat. Methodol.*, **39B**, 1–38.
- Doucet, A., S. Godsill, and C. Andrieu, 2000: On sequential Monte Carlo sampling methods for Bayesian filtering. *Stat. Comput.*, **10**, 197–208, doi:10.1023/A:1008935410038.
- 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, doi:10.1256/qj.05.167.
- Fraley, C., A. E. Raftery, and T. Gneiting, 2010: Calibrating multimodel forecast ensembles with exchangeable and missing members using Bayesian model averaging. *Mon. Wea. Rev.*, **138**, 190–202, doi:10.1175/2009MWR3046.1.
- Glahn, B., M. Peroutka, J. Wiedenfeld, J. Wagner, G. Zylstra, B. Schulknecht, and B. Jackson, 2009: MOS uncertainty estimates in an ensemble framework. *Mon. Wea. Rev.*, **137**, 246–268, doi:10.1175/2008MWR2569.1.
- Gneiting, T., A. E. Raftery, 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, doi:10.1175/MWR2904.1.
- , F. Balabdaoui, and A. E. Raftery, 2007: Probabilistic forecasts, calibration, and sharpness. *J. Roy. Stat. Soc. Stat. Methodol.*, **69B**, 243–268, doi:10.1111/j.1467-9868.2007.00587.x.
- Golub, G. H., and C. F. Van Loan, 1989: *Matrix Computations*. 2nd ed. Johns Hopkins Press, 642 pp.
- Hamill, T. M., 2007: Comments on “Calibrated surface temperature forecasts from the Canadian ensemble prediction system using Bayesian Model Averaging.” *Mon. Wea. Rev.*, **135**, 4226–4236, doi:10.1175/2007MWR1963.1.
- , and J. S. Whitaker, 2006: Probabilistic quantitative precipitation forecasts based on reforecast analogs: Theory and application. *Mon. Wea. Rev.*, **134**, 3209–3229, doi:10.1175/MWR3237.1.
- , M. Scheuerer, and G. T. Bates, 2015: Analog probabilistic precipitation forecasts using GEFS reforecasts and climatology-calibrated precipitation analyses. *Mon. Wea. Rev.*, **143**, 3300–3309, doi:10.1175/MWR-D-15-0004.1.
- Han, J., and H.-L. Pan, 2011: Revision of convection and vertical diffusion schemes in the NCEP Global Forecast System. *Wea. Forecasting*, **26**, 520–533, doi:10.1175/WAF-D-10-05038.1.
- Houtekamer, P. L., 1993: Global and local skill forecasts. *Mon. Wea. Rev.*, **121**, 1834–1846, doi:10.1175/1520-0493(1993)121<1834:GALSF>2.0.CO;2.
- Krzysztofowicz, R., and W. B. Evans, 2008: Probabilistic forecasts from the national digital forecast database. *Wea. Forecasting*, **23**, 270–289, doi:10.1175/2007WAF2007029.1.
- Luo, L., E. F. Wood, and M. Pan, 2007: Bayesian merging of multiple climate model forecasts for seasonal hydrological

- predictions. *J. Geophys. Res.*, **112**, D10102, doi:[10.1029/2006JD007655](https://doi.org/10.1029/2006JD007655).
- Raftery, A. E., T. Gneiting, F. Balabdaoui, and M. Polakowski, 2005: Using Bayesian Model Averaging to calibrate forecasts. *Mon. Wea. Rev.*, **133**, 1155–1174, doi:[10.1175/MWR2906.1](https://doi.org/10.1175/MWR2906.1).
- Rajagopalan, B., U. Lall, and S. E. Zebiak, 2002: Categorical climate forecasts through regularization and optimal combination of multiple GCM ensembles. *Mon. Wea. Rev.*, **130**, 1792–1811, doi:[10.1175/1520-0493\(2002\)130<1792:CCFTRA>2.0.CO;2](https://doi.org/10.1175/1520-0493(2002)130<1792:CCFTRA>2.0.CO;2).
- Rosmond, T. E., 1992: The design and testing of the Navy Operational Global Atmospheric Prediction System. *Wea. Forecasting*, **7**, 262–272, doi:[10.1175/1520-0434\(1992\)007<0262:TDATOT>2.0.CO;2](https://doi.org/10.1175/1520-0434(1992)007<0262:TDATOT>2.0.CO;2).
- Roulston, M. S., and L. A. Smith, 2003: Combining dynamical and statistical ensembles. *Tellus*, **55A**, 16–30, doi:[10.1034/j.1600-0870.2003.201378.x](https://doi.org/10.1034/j.1600-0870.2003.201378.x).
- Sloughter, J. M., A. E. Raftery, T. Gneiting, and C. Fraley, 2007: Probabilistic quantitative precipitation forecasting using Bayesian model averaging. *Mon. Wea. Rev.*, **135**, 3209–3220, doi:[10.1175/MWR3441.1](https://doi.org/10.1175/MWR3441.1).
- , T. Gneiting, and A. E. Raftery, 2010: Probabilistic wind speed forecasting using ensembles and Bayesian model averaging. *J. Amer. Stat. Assoc.*, **105**, 25–35, doi:[10.1198/jasa.2009.ap08615](https://doi.org/10.1198/jasa.2009.ap08615).
- Tolman, H. L., 1997: User manual and system documentation of WAVEWATCH-III version 1.15. NOAA/NWS/NCEP/OMB Tech. Note 151, 97 pp.
- , 1999: User manual and system documentation of WAVEWATCH-III version 1.18. NOAA/NWS/NCEP/OMB Tech. Note 166, 110 pp.
- , 2001: Improving propagation in ocean wave models. *Ocean Wave Measurement and Analysis*, B. L. Edge and J. M. Hemsley, Eds., ASCE, 507–516.
- Unger, D. A., H. Van den Dool, E. O’Lenic, and D. Collins, 2009: Ensemble regression. *Mon. Wea. Rev.*, **137**, 2365–2379, doi:[10.1175/2008MWR2605.1](https://doi.org/10.1175/2008MWR2605.1).
- Vrugt, J., C. J. H. Diks, and M. P. Clark, 2008: Ensemble Bayesian Model Averaging using Markov-chain Monte Carlo sampling. *Environ. Fluid Mech.*, **8**, 579–595, doi:[10.1007/s10652-008-9106-3](https://doi.org/10.1007/s10652-008-9106-3).
- Wang, X., and C. H. Bishop, 2003: A comparison of breeding and ensemble transform Kalman filter ensemble forecast schemes. *J. Atmos. Sci.*, **60**, 1140–1158, doi:[10.1175/1520-0469\(2003\)060<1140:ACOBAE>2.0.CO;2](https://doi.org/10.1175/1520-0469(2003)060<1140:ACOBAE>2.0.CO;2).
- Wilks, D. S., 2006: Comparison of ensemble-MOS methods in the Lorenz ’96 setting. *Meteor. Appl.*, **13**, 243–256, doi:[10.1017/S1350482706002192](https://doi.org/10.1017/S1350482706002192).
- 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, doi:[10.1175/MWR3347.1](https://doi.org/10.1175/MWR3347.1).
- Yeo, I.-K., and R. A. Johnson, 2000: A new family of power transformations to improve normality or symmetry. *Biometrika*, **87**, 954–959, doi:[10.1093/biomet/87.4.954](https://doi.org/10.1093/biomet/87.4.954).