



# 1 On the Effect of Model Parameters on Forecast Objects

2 Caren Marzban<sup>1,2\*</sup>, Corinne Jones<sup>2</sup>, Ning Li<sup>2</sup>, Scott Sandgathe<sup>1</sup>

<sup>1</sup> Applied Physics Laboratory

<sup>2</sup> Department of Statistics

Univ. of Washington, Seattle, WA 98195 USA

## 3 ABSTRACT

4 Many physics-based numerical models produce a gridded, spatial field of forecasts, e.g., a  
5 temperature “map.” However, the field for some quantities such as precipitation generally  
6 consists of spatially coherent and disconnected “objects.” Certain features of these objects  
7 (e.g., number, size, and intensity) are generally of interest. Here, a methodology is developed  
8 for assessing the impact of model parameters on features of forecast objects. Although, in  
9 principle, the objects can be defined by any means, here they are identified via clustering  
10 algorithms. The methodology is demonstrated on precipitation forecasts from a mesoscale  
11 numerical weather prediction model.

12 The author’s copyright for this publication is transferred to University of Washington.

\*Corresponding Author: [marzban@stat.washington.edu](mailto:marzban@stat.washington.edu)



## 13 1. Introduction

14 Complex, physics-based numerical models of natural phenomena often have parameters -  
15 henceforth, model parameters - whose values are generally not *a priori* specified. In such sit-  
16 uations it is important to infer the manner in which the model parameters affect the outputs  
17 of the model (i.e., forecasts, or predictions), and often the techniques of Sensitivity Analysis  
18 (SA) are employed to assess the effects. There is a wide range of techniques from relatively  
19 simple one-at-a-time method (also known as the Morris method) where each model param-  
20 eter is varied individually (e.g., Yu et al. (2013)), to multivariate approaches motivated by  
21 statistical methods of experimental design (Montgomery 2009) where the values of the model  
22 parameters are varied according to some optimization criterion. Alternative approaches can  
23 be found in Backman et al. (2017) where algorithmic differentiation is used, and in Kalra  
24 et al. (2017) where the underlying physics equations are integrated using quadrature meth-  
25 ods. And yet another alternative is the adjoint method, commonly used in meteorological  
26 circles (Errico 1997).

27 It is difficult to classify these methods into a simple taxonomy (Bolado-Lavin and Badea  
28 2008), but the terms Local and Global have been used to denote two broad categories  
29 (Saltelli et al. 2010, 2008); generally, local methods employ some sort of derivative of the  
30 model output with respect to inputs, while global techniques rely on a decomposition of the  
31 variance of the output in terms of the variance explained by the inputs. Comparisons of  
32 the various approaches are not common-place, because each approach is usually suited for  
33 a specific application where other methods may not be practically feasible. However, an  
34 example of the comparison of one global approach and one local (adjoint) approach on the



35 Lorenz '63 model (Lorenz 1963) has been performed by Marzban (2013).

36 Another possible classification criterion is based on the purpose of the SA. Some SA  
37 work is performed for assessing how model parameters impact the model itself, not as a  
38 means to some other goal. For example, Lucas et al. (2013) uses a global SA method to  
39 explore the effect of model parameters on the probability of model crashes. By contrast,  
40 sometimes SA is performed as an intermediate step to another goal, such as the calibration  
41 of the model (Safta et al. 2015; Hacker et al. 2011; Laine et al. 2012; Ollinaho et al. 2014).  
42 All of these classification criteria are imperfect, as there exist works which fall “between”  
43 Global versus Local, or SA-only versus SA-for-calibration; some examples include Roebber  
44 (1989); Roebber and Bosart (1989); Robock et al. (2003). The work reported here falls into  
45 the Local and SA-only category; as such, although the proposed methodology can be used  
46 for calibration, no attempt is made to do so here.

47 In many SA studies, the output of the model (i.e., the response variable in the SA) is  
48 usually a single or a handful of scalar quantities. But there are situations in which the output  
49 is a gridded spatial field, e.g., temperature forecasts over a spatial region. Every grid point  
50 reflects a forecast at that location, and for a quantity like temperature the field as a whole  
51 has a smooth, continuous nature. SA is more complicated for precipitation fields, where  
52 the model output is a quantity whose spatial structure is not smooth and/or continuous.  
53 Indeed, there may be a coherent set of grid points that receive no precipitation at all, while  
54 an adjacent set of grid points will reflect a complex pattern of precipitation. In short, the  
55 spatial field of such quantities will contain “objects” within which precipitation does occur,  
56 surrounded by regions of little or no precipitation.

57 For such discrete fields, the assessment of the quality of the forecasts has given rise



58 to a wide range of specialized techniques generally referred to as spatial verification (or  
59 evaluation) (Ahijevych et al. 2009; Baldwin et al. 2001, 2002; Brown et al. 2002; Casati  
60 et al. 2004; Davis et al. 2006a,b; Du and Mullen 2000; Ebert 2008; Ebert and McBride 2000;  
61 Gilleland et al. 2009; Hoffman et al. 1995; Keil and Craig 2007; Marzban and Sandgathe  
62 2006, 2008; Marzban et al. 2008, 2009; Nachamkin 2004; Roberts and Lean 2008; Wealands  
63 et al. 2005; Wernli et al. 2008; Venugopal et al. 2005). A subset of these methods employs  
64 the notion of an object explicitly. In some applications, the object is defined subjectively  
65 - for example, by human experts. In other applications statistical methods for clustering  
66 (Everitt 1980) are used to identify/define objects within the field (Marzban and Sandgathe  
67 2006, 2008). This clustering approach, which has been re-examined by Lakshmanan and  
68 Kain (2010), and more recently by Wang et al. (2015), is the basis of the object-identification  
69 procedure used in the present work.

70 Although no spatial verification/evaluation is done here, the importance of objects within  
71 the forecast field, and the development of clustering techniques for identifying them, calls  
72 for a SA framework wherein one can assess the effect of model parameters on the objects. In  
73 meteorology certain features of the clusters/objects are of special interest; they include size,  
74 location, intensity, and shape. Also, the assessment of sensitivity is highly intertwined with  
75 that of statistical significance. As such, the methodology developed here can be viewed as  
76 a SA with a multivariate response, wherein one can assess the impact (both the magnitude  
77 and the statistical significance) of model parameters on object features.

78 The model employed to demonstrate the methodology is COAMPS<sup>®</sup> (Hodur 1997), for  
79 which some SA work has already been done. Doyle et al. (2011) and Jiang and Doyle (2009)  
80 examine the effect of model parameters on mountain waves. Motivated by the work of Holt



81 et al. (2011) who studied the effect of 11 model parameters on various characteristics of the  
82 forecasts, Marzban et al. (2014) used a global (variance-based) SA to study the effect of the  
83 same parameters and their interactions on mean (across the forecast domain) precipitation,  
84 and the center-of-gravity of precipitation.

## 85 2. Method

86 The methodology described in this paper involves two other techniques developed pre-  
87 viously by some of the authors of this paper. In one, cluster analysis is used for identifying  
88 objects (Marzban and Sandgathe 2006, 2008; Marzban et al. 2008, 2009); in the other, SA  
89 is performed to assess the effect of model parameters on non-spatial features (e.g., domain  
90 mean) of the forecast field (Marzban et al. 2014). This section describes these components,  
91 puts forth the SA model, proposes means of assessing sensitivity and statistical significance,  
92 and describes the data used to demonstrate the methodology.

### 93 a. Data

94 The inputs of the numerical model examined here are 11 model parameters, and the  
95 outputs are forecasts of precipitation at each of  $45 \times 72$  grid points, with a spacing of  $81\text{km}$ ,  
96 covering the entire continental US, including coastal regions, and portions of Canada and  
97 Mexico. The SA method developed here requires data - technically, *computer data* - which  
98 are created by generating an ensemble (or sample) of inputs values, assimilating surface  
99 observations, and then running the model forward to produce 24h forecasts of precipitation



100 amount at each grid point. As such, the SA results are contingent on the nature of this data,  
101 and consequently, care must be taken in the data-generation step of the methodology.

102 In order to include a wide range of weather phenomena, the data include 120 days from  
103 February 16 through July 2, 2009. Confirmed by visual examination of all 120 forecasts, this  
104 temporal period includes a comprehensive series of midaltitude synoptic systems traveling  
105 across the northern portion of the domain. These synoptic systems extend down into the  
106 southeastern US early in the period and are replaced by subtropical convective systems in the  
107 late spring and summer months. This subtropical activity also occurs in the southwestern  
108 portion of the domain (west coast of Mexico) during June and July in association with the  
109 southwest monsoon. The only apparent atypical weather appears to be a greater amount  
110 of convective activity off the east coast of the US associated with quasi-stationary or slow  
111 moving frontal systems during the period.

112 It is important that the data cases are as independent as possible. To that end, the 120  
113 days are sampled at 3-day intervals in order to minimize temporal dependency, leading to  
114 40 days for the analysis.

115 For each of the 40 days, 99 different values for 11 parameters are generated by Latin  
116 Hypercube Sampling (LHS). Said differently, for each day, a sample of size 99 is taken from  
117 the 11-dimensional space of the model parameters. This so-called “space-filling” sampling  
118 scheme assures that no two of the 99 points have the same value for any of the 11 parameters.  
119 It can be shown that this property leads to more precise estimates (at least, no less-precise  
120 estimates) than alternative sampling schemes (Cioppa and Lucas 2007; Montgomery 2009;  
121 Marzban 2013). LHS is appropriate when the model parameters are all continuous quantities  
122 (i.e., taking values on the Real line). For discrete or categorical inputs, Latin Square Designs



123 (LSD) or Fractional Factorial Designs (FFD) can be employed to produce optimal samples  
124 (Montgomery 2009); these cases will be demonstrated in a separate article.

125 The 11 model parameters are shown in Table 1; the choice of these parameters is ex-  
126 plained in Holt et al. (2011). As mentioned in that paper, these parameters were chosen  
127 for their anticipated sensitivity (through model tests and discussions with developers) of the  
128 parameterizations in an effort to chose parameters most likely to produce changes in the  
129 model output precipitation fields. Also, to focus on heavy precipitation, only the grid points  
130 whose convective precipitation amount exceeds the 90th percentile of precipitation across  
131 the domain are analyzed.

132 A very similar data set is used in Marzban et al. (2014) to assess the sensitivity of the  
133 average and center-of-gravity of precipitation (across the domain) with respect to the model  
134 parameters. Here, however, the precipitation fields are first subjected to cluster analysis  
135 (Sect. 3d), and then six cluster features (Sect. 3e) are employed as response variables in a  
136 multivariate SA.

### 137 *b. Statistical Model*

138 The SA methodology in (Marzban et al. 2014) is a variance-based approach which allows  
139 one to identify linear or nonlinear relationships between the forecast quantities and the model  
140 parameters, and even interactions between the model parameters. As a first approximation,  
141 however, it is sufficient to estimate only the linear (i.e., main) effects, because nonlinear and  
142 interaction effects are often much smaller than main effects; see pages 192, 230, 272, 314,  
143 329 in Montgomery (2009), and pages 33-34 in Li et al. (2006). For this reason a linear



144 regression-based model is adequate. Specifically, the effects of the model parameters are  
145 assessed via the least-squares estimate of the regression coefficients  $\beta_i$  in

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{11} x_{11} + \delta, \quad (1)$$

146 where  $x_i$  denote standardized model parameters,  $y$  is the forecast quantity (e.g., some cluster  
147 feature), and  $\delta$  represents any source of variability in  $y$  other than from the model param-  
148 eters. This linear model is further justified by the results (shown below) because when it  
149 is specialized to the case of one cluster (i.e., the entire spatial domain), it reproduces the  
150 results of the variance-based approach reported in Marzban et al. (2014).

151 There exists a realization of Eq. (1) in which the response is vector-valued; the model is  
152 called Multivariate Multiple Regression (MMR), wherein Eq. (1) is understood as a vector  
153 equation, where  $y$ ,  $\alpha$ , and  $\beta_i$  are all vectors. Ideally one could allow each component of the  
154 response vector to represent a forecast feature of a given object. However, the number of  
155 objects/clusters varies across the 99 values of the parameters and across days in the data.  
156 Methods for estimating MMR coefficients when the number of responses is a random variable  
157 (varying across cases) are not readily available. Therefore, for each of the six features, we  
158 consider three summary measures: The minimum, median, and maximum (across the clusters  
159 in the domain) of the feature. These three quantities can be thought of as a 3-point summary  
160 of the distribution of the feature, and they serve as the three responses in MMR.

161 The median across clusters is useful, because one can assess the effect of the model  
162 parameters on a “typical” cluster; minimum and maximum across clusters are useful because  
163 they allow one to assess whether a model parameter has an effect on **any** of the clusters in a  
164 field. For example, if it is found that a particular model parameter is positively (negatively)





165 associated with the minimum (maximum) size across clusters, then one can conclude that  
166 the size of at least one of the clusters in the field is affected by that parameter. This is  
167 an important consideration, because if the size of at least one of the clusters is not affected  
168 by a parameter, then that parameter can be said to have no effect on the size of clusters.  
169 Additionally, consideration of the three summary measures, together, allows one to assess  
170 the effect of the model parameters on the distribution of the features.

171 The data on the response variables  $y$  are log-transformed to assure more bell-shaped  
172 histograms; this transformation is not necessary, but is useful when the regression coefficients  
173 are subjected to statistical tests, because many such tests assume relatively bell-shaped  
174 distributions.

### 175 *c. Significance Tests*

176 Testing the coefficients in the MMR model involves performing a large number of sta-  
177 tistical tests ( $40 \times 11 \times 6 \times 3$ ): one on each of 40 days, for each of 11 parameters, for each  
178 of six cluster features, and for each of three summary measures across clusters. A large  
179 number of tests, in turn, leads to an exponential growth in the probability of making *some*  
180 Type I error - a fact known as the multiple hypothesis testing problem (Montgomery 2009).  
181 A standard procedure in statistics for taming Type I errors is to divide the task into two  
182 stages (Montgomery 2009). In the first stage, one performs a single, often-called omnibus,  
183 hypothesis test of whether any of the parameters have an effect on any of the responses.  
184 In the present application, such a test reduces the number of tests to  $40 \times 6$ . If the null  
185 hypothesis cannot be rejected, then one performs no more tests, and the conclusion of the

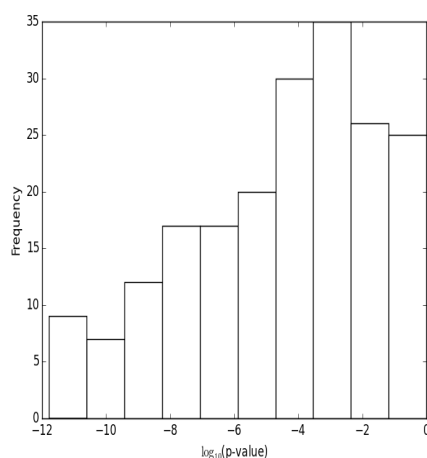


FIG. 1. Histogram of p-values from the multivariate tests across all days and response variables.

186 analysis is that there is no evidence that any of the parameters have an effect on any of the  
187 responses. If, however, the null hypothesis is rejected, then one proceeds to the second stage,  
188 i.e., testing each of the  $40 \times 6$  effects, separately.

189 For the first stage, omnibus tests are readily available within MMR models (DelSole and  
190 Yang 2011; Fox et al. 2013; Rencher and Christensen 2012). Here, these tests were performed,  
191 yielding extremely small p-values, i.e., highly significant results (see Fig. 1), necessitating  
192 the second stage analysis.

193 For the second stage, a number of methods have been developed, again for the purpose  
194 of taming Type I errors; two of the more commonly employed methods are due to Tukey  
195 and Dunnett (Montgomery 2009). But these tests are generally complex procedures which  
196 in the end still involve a simplistic comparison of a p-value with a prespecified significance  
197 level. Although sufficient for hypothesis testing, these p-values provide no information on



198 the magnitude of the effect. For this reason, instead, we adopt the more qualitative approach  
199 of examining the boxplot of the estimated regression coefficients directly.

200 The boxplots (shown in the next section) are generated and analyzed as follows. For  
201 each of the six cluster features, the response vector  $y$  is set to the minimum, median, and  
202 maximum (across clusters in the whole field) of that feature. For each of these three response  
203 variables, boxplots of the regression coefficients for the 11 model parameters are produced.  
204 The degree of overlap between each boxplot and the number zero reflects a visual (though  
205 qualitative) assessment of both the statistical significance and the magnitude of the effect of  
206 the corresponding model parameter on the response: If zero is well within the span of the  
207 boxplot, then one cannot conclude anything regarding the effect; if the boxplot is significantly  
208 above (below) zero, then one can conclude that the corresponding parameter has a positive  
209 (negative) effect on the response in question; and in such a case, the “distance” of the boxplot  
210 relative to zero provides a visual indication of the magnitude of the effect.

#### 211 *d. Cluster Analysis*

212 There exists a wide range of clustering methods, each with their respective parameters  
213 (Everitt 1980). At one extreme, there exists a class of clustering methods wherein the  
214 desired number of cluster,  $NC$ , is specified by the user. A proven example in this class is  
215 called Gaussian Mixture Model (GMM) clustering (McLachlan and Peel 2000). At the other  
216 extreme, there exist clustering routines where  $NC$  does not play a role at all. One such  
217 method is called Density-Based Spatial Clustering of Applications with Noise (DBSCAN)  
218 (Ester et al. 1996). DBSCAN has two parameters, here denoted  $\epsilon$  and  $\text{min\_samples}$ . Roughly



219 speaking,  $\epsilon$  is the maximum distance between two grid points in order for them to be in the  
220 same cluster, and `min_samples` is the minimum number of grid points necessary to form a  
221 cluster.

222 These two approaches are selected here for demonstration because they allow for two  
223 very different ways in which a user can inject *a priori* knowledge into the analysis. For  
224 example, in some applications it may be more natural to specify the number of clusters,  
225 in which case GMM is a natural choice. On the other hand, DBSCAN is more natural  
226 if the user has knowledge of the typical size and distance between clusters. For example,  
227 consider a situation wherein the grid-spacing is relatively large (as is the case in this paper,  
228 i.e.,  $81\text{km}$ ), allowing one to examine only large scale precipitation. Although time of year  
229 and location are also important, but if one were to focus only on winter months in, say, the  
230 Pacific Northwest, then it is reasonable to set  $\epsilon$  to 3 or 4. By contrast, if one is considering  
231 jet streaks, e.g., where some maximum wind speed value is reached, then  $\epsilon$  can be closer  
232 to 1. As for `min_samples`, 4 or 5 are reasonable values for both precipitation and jet streak  
233 events, at the model resolution used here.

234 In addition to the way in which the respective parameters are handled, another reason  
235 why these two clustering methods are used here is that they occupy two other extremes in the  
236 family of clustering algorithms: GMM clustering belongs to a class of model-based algorithms  
237 (Banfield and Raftery 1993; Fraley and Raftery 2002) common in statistics circles because  
238 they are conducive to performing statistical tests, while DBSCAN assumes no underlying  
239 model, and for this reason is often employed in machine learning applications. For the SA  
240 component of the methodology developed here, it is not necessary for the objects to be  
241 defined by these or any other clustering algorithm; the objects may be defined by any other



242 criterion or even by human experts.

243 *e. Cluster Features*

244 In spatial verification some of the errors that are of interest include displacement, in-  
245 tensity, size/area, and shape error. The estimation of these errors presumes the ability to  
246 compute, respectively, the location, intensity, area, and shape of a cluster. Here, the latitude  
247 and longitude of the centroid of a cluster are taken as coordinates of its location; intensity is  
248 measured by the median (across the spatial extent of the cluster) of precipitation; and area  
249 is measured by the number of grid points in a cluster. The shape of a cluster, in GMM, is an  
250 ellipse, because that is the cross-section (i.e., level-set) of a bivariate Gaussian. Then, the  
251 eccentricity and orientation of the semi-major axis of the ellipse are natural for quantifying  
252 the shape of clusters. In DBSCAN, clusters are not restricted to have any specific shape.  
253 In order to be able to compare the two clustering algorithms, here an ellipse is “fitted” to  
254 the clusters, and again the eccentricity and orientation of the semi-major axis is used to  
255 represent the shape of the cluster. (Technically, the direction of the semi-major axis is de-  
256 fined to be the direction of the first eigenvector of the covariance matrix computed from the  
257 coordinates of all the grid points in a given cluster. The length of the semi-major axis is set  
258 to the largest eigenvalue.)

259 In short, the six cluster features examined here are latitude, longitude, intensity, area,  
260 orientation, and eccentricity. Also, recall that for each of these features, three summary  
261 measures are computed - minimum, median, and maximum - and used as the multivariate  
262 response vector in MMR (see Sect. 3b).

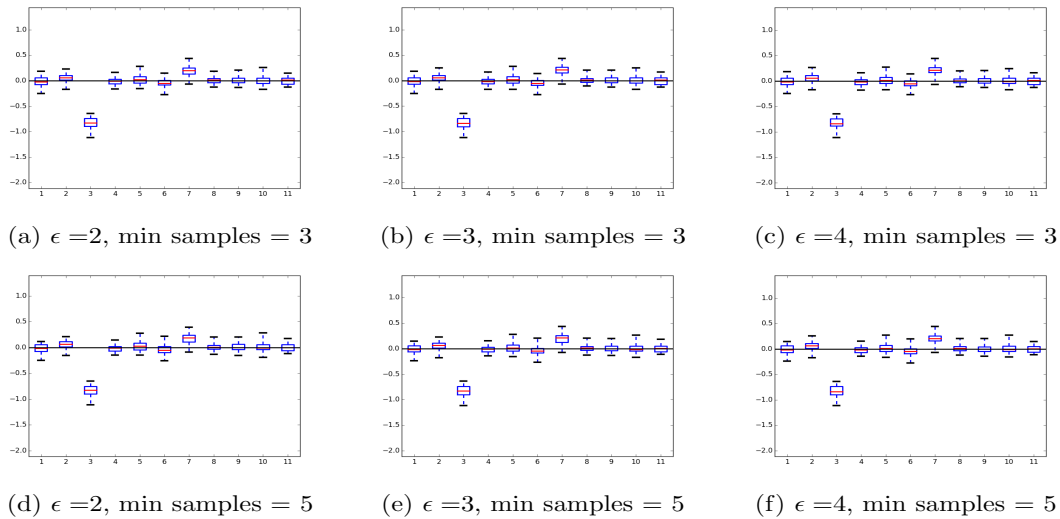


FIG. 2. Estimated regression coefficients (i.e. sensitivity of the model parameters) with median precipitation of the clusters as the response, after clustering with DBSCAN with various parameter values.

### 263 3. Results

264 As mentioned earlier, 40 forecasts are produced, each with 99 different values of 11 pa-  
265 rameters in COAMPS. Each forecast field is clustered, and three summary measures (mini-  
266 mum, median, and maximum, all across clusters) are computed, each for six cluster features  
267 (latitude, longitude, etc.). First, an omnibus test is performed to test whether any of the  
268 11 parameters have an effect on any of the three summary measures, on each day and for  
269 each cluster feature. Then, six MMR models are set up mapping the 11 parameters to three  
270 response variables. The daily variability - displayed as boxplots (e.g., Fig. 2 and Fig. 3) - for  
271 each of the regression coefficients offers a visual assessment of both the statistical significance  
272 and the magnitude of the effect of each parameter.



273 In the first stage of the analysis the response variable is a 3-dimensional vector, and  
274 an omnibus test is performed to test if any of the 11 parameters have an effect on any of  
275 the three response variables, for each day and each cluster feature. Such a test reduces the  
276 number of tests from  $(40 \times 11 \times 6 \times 3)$  to  $(40 \times 6) = 240$ . The individual p-values are  
277 not shown here, but for DBSCAN their histogram is shown in Fig. 1. Evidently, all of  
278 the comparisons yield extremely small p-values. At a significance level of 0.05, out of the  
279 240 tests, 29 p-values are not significant when using DBSCAN and 59 are not significant  
280 when using GMM. By examining the p-values, the majority of the non-significant results are  
281 associated with the tests when the response is the eccentricity of a cluster. If one applies  
282 the Bonferroni correction (Devore and Farnum 2005) to the significance level in order to  
283 account for the multiple tests, the significance level becomes  $0.05/(40 \times 6) = 2 \times 10^{-4}$ . At  
284 this significance level there are many more nonsignificant comparisons: 87 for DBSCAN and  
285 94 for GMM. Upon making this correction, in addition to eccentricity some of the other  
286 features also emerge as being unaffected by any of the 11 parameters. Further details of  
287 these results are presented below.

288 Figure 2 shows the sensitivity results when the response is the median (across clusters) of  
289 precipitation intensity, and DBSCAN is employed with different parameters. The analogous  
290 results for GMM with different values of  $NC$  are not shown here, but they are similar. Recall  
291 that the variability displayed in each boxplot is due to the 40 days examined. First, note  
292 that all of the panels are mostly similar to one another, which implies that the sensitivity  
293 results are mostly unaffected by the parameters of the clustering algorithm.

294 It can also be seen that many of the 11 parameters have a histogram/boxplot of values  
295 mostly around zero. In other words, when considered across multiple days most of the 11



296 model parameters have no effect on the median of precipitation, The most obvious exception  
297 is parameter 3, which by virtue of having mostly negative values for its regression coefficient,  
298 is negatively associated with median precipitation. Parameter 7 not only has a weaker  
299 effect (because the median of the corresponding boxplot is closer to zero), it is also not  
300 as statistically significant (because zero falls well within the span of the boxplot). This  
301 parameter is positively associated with precipitation intensity in the typical (median) cluster,  
302 i.e., increasing the parameter leads to more intense clusters; more, below. All of these findings  
303 are consistent with those found for convective precipitation in Marzban et al. (2014) where  
304 a variance-based sensitivity was performed without any clustering at all. This consistency  
305 adds justification to the local/regression-based SA adopted here, i.e., Eq. (1).

306 Figure 3 shows the effect of the model parameters on the latitude and longitude of the  
307 clusters (top two rows), amount of precipitation (middle row) in the clusters, and the area  
308 and orientation of the clusters (bottom two rows). The three columns correspond to the  
309 minimum, median, and maximum of a feature. Eccentricity has also been examined, but  
310 the results are not shown here because it is not affected by any of the 11 parameters; this  
311 conclusion is consistent with the results of the F-test performed in the first stage, mentioned  
312 above.

313 Examination of all of the panels suggests that parameters 4, 5, 8, 9, 10, 11 have little or  
314 no effect on any of the object features. By contrast, parameters 1, 2, 3, 6, and 7 appear to  
315 have varying effects depending on the object feature. Also, the orientation (in addition to  
316 eccentricity) of the clusters is unaffected by any of the parameters.

317 The strongest effects are from parameters 3 and 7 on the amount of precipitation. This  
318 relationship was already examined in Fig. 2; but now the same pattern can be seen in the



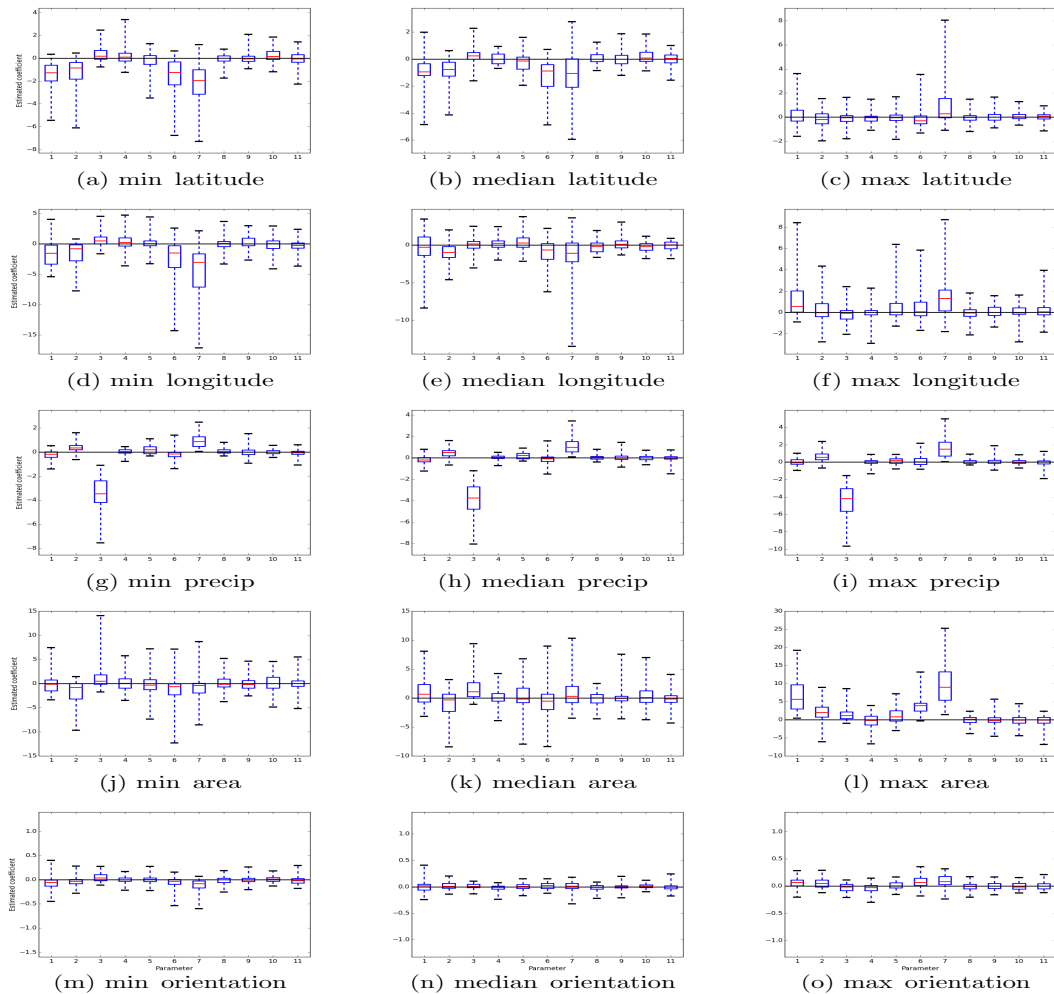


FIG. 3. Estimated regression coefficients (i.e. sensitivity of the model parameters) on three summary measures (minimum, median, maximum) of different cluster features (latitude, longitude, amount of precipitation, and area and orientation of clusters. Eccentricity is not shown (see text). The clustering is done with DBSCAN with  $\epsilon = 2\sqrt{2}$ , `min_samples = 3`.



319 minimum, median, and maximum intensity (panels g, h, i in Fig. 3), which implies that the  
320 effect of parameters 3 and 7 is to shift down and up, respectively, the whole distribution of  
321 precipitation intensity.

322 The next strongest effects are those of parameters 1 and 7 on maximum area (panel l).  
323 Given that these two parameters have no effect on the minimum and median area (panels j  
324 and k), it follows that these parameters affect only the right tail of the distribution of size. In  
325 other words, by contrast to precipitation intensity whose distribution shifts when parameter  
326 7 is varied, the distribution of size is stretched when that parameter changes. Parameter 6,  
327 too, appears to have an effect on maximum area, but to a lesser extent, both statistically  
328 and in magnitude.

329 Whereas parameter 1 tends to stretch out the distribution of area to the right, it appears  
330 to have the opposite effect on the minimum and median longitude of the clusters. The effect  
331 is weak in magnitude, but statistically significant. It does not affect the maximum longitude  
332 (panel f), and so, it stretches the distribution of longitude on the left, causing clusters to  
333 appear with smaller longitude, which given the encoding of the data used here, means to the  
334 west. Parameters 2, 6, and 7 appear to have the same effect as parameter 1.

335 The latitude appears to be weakly affected by some of the parameters. For example,  
336 parameter 7, and to a much lesser degree parameter 1, is positively associated with median  
337 and maximum latitude, but negatively associated with minimum latitude. In other words,  
338 increasing parameter 7 increases the width of the distribution of latitude values, causing  
339 them to be more spread out along the latitudes.

340 All of the above conclusions are based on clustering with DBSCAN with  $\epsilon = 2\sqrt{2}$  and  
341  $\text{min\_samples}=3$ . To test the robustness of these results the same analysis was repeated but



342 with GMM as the cluster algorithm and with  $NC = 3$ . The results (not shown here) are  
343 mostly the same. One relatively clear difference between the DBSCAN and GMM results is  
344 in the effect of parameters 1 and 7 on area; whereas with DBSCAN those parameters have  
345 an effect only on the maximum area, the results based on GMM suggest a significant effect  
346 on all three cluster features (minimum, median, and maximum area).

347 Further differences between DBSCAN and GMM sensitivity results are found when one  
348 performs a multivariate test for the effect of the model parameters across **all** days. For  
349 DBSCAN, the p-values corresponding to each of the six cluster features are all found to be  
350 nearly zero. So, some of the model parameters do have a significant effect on some of the  
351 features. The same is true for GMM, with the exception of latitude and eccentricity for which  
352 there is no evidence of an effect (p-values 0.435 and 0.290, respectively). It may appear that  
353 these results are contradictory, but they are not because the respective parameters of the  
354 two clustering algorithms have not been tuned to render them comparable. Specifically, the  
355 DBSCAN parameters are  $\epsilon = 2\sqrt{2}$  and  $\text{min\_samples}=3$ , while for GMM the parameter  $NC$   
356 is set to three. In other words, the differences are due to the way in which the two clustering  
357 algorithms handle their respective parameters. As mentioned earlier, such differences do not  
358 point to defects in the methodology; they simply reflect the choice of what the user considers  
359 to be an object.



## 360 4. Conclusion and Discussion

361 It is shown that by employing methods of cluster analysis and sensitivity analysis one  
362 can assess the magnitude and statistical significance of the effect of model parameters on  
363 features (location, intensity, size, and shape) of objects within forecast fields. The framework  
364 also allows one to assess the impact of the model parameters on the distribution of forecast  
365 features. For example, one can reveal the model parameters that affect the overall location  
366 and/or width of the distribution of object features, and those which impact the shape of the  
367 distribution, e.g., by stretching out the left and/or right tail. The approach does not point to  
368 any “optimal” values of the model parameters, for that would require optimizing the model  
369 parameters to maximize some measure of agreement between forecasts and observations. In  
370 other words, although the work here lays the foundation for tuning the model parameters  
371 for the purpose of improving forecasts in terms of metrics that arise naturally in spatial  
372 verification/evaluation methods, no such tuning is performed here.

373 Given the novelty of the proposed framework, some recommendations are in order. The  
374 choice of the clustering algorithm depends on the specific user. Indeed, there are situations  
375 in which clusters/objects in a field are identified by human experts. For these reason, no  
376 specific clustering algorithm is recommended. A similar philosophy is adopted with respect  
377 to the values of the parameters of the clustering algorithms; they may be specified by the  
378 user, or varied across a range of values, depending on the specific application. Although  
379 there exist statistical criteria that lead to unique values for the parameters, the criteria  
380 involve the optimization of some other quantity, e.g., Akaike Information Criterion (AIC) or  
381 Bayesian Information Criterion (BIC). As such, the ambiguity in the choice of the clustering



382 algorithm, or the values of their parameters, is simply replaced with the ambiguity of selecting  
383 the appropriate criterion. Therefore, again, no attempt is made to optimize the values of  
384 the parameters. It is assumed that the user has sufficient information about the underlying  
385 physics to either specify the number of physical objects (or a range thereof), or the typical  
386 size and distance between physical objects.

387 It is worth pointing out that at least in meteorology, it is not uncommon for different  
388 human experts to have different notions of an object in the forecast field. As such, the  
389 ambiguities discussed above are not specific to clustering algorithms, but are inherent to  
390 any object-based approach. In spite of this inherent ambiguity, many spatial verification  
391 techniques generally rely on some notion of an object. The main reason is that accounting  
392 for objects in a forecast field is a first step in the verification/evaluation process, and the  
393 manner in which objects are defined is of secondary importance.

394 While this paper is primarily about a methodology, it is worthwhile to provide a possible  
395 physical explanation for at least the strongest results in the COAMPS application. The  
396 strongest influence or sensitivity is from parameter 3, the fraction of available precipitation  
397 fed back to the grid from the Kain-Fritsch scheme. Increasing this fraction reduces con-  
398 vective precipitation and, based on the results in Marzban et al. (2014), increases stable  
399 precipitation, while not affecting total precipitation. It also is responsible for weakening  
400 the convective precipitation, i.e., increasing the number of weak systems. The next largest  
401 sensitivity is from parameter 7, which controls the temperature difference required to ini-  
402 tiate convective precipitation. Again, as shown in Marzban et al. (2014), this parameter  
403 also controls a trade-off between convective and stable precipitation and has little effect on  
404 total precipitation (along with parameter 1). Parameters 1 and 7 do increase the area of



405 convective precipitation in large precipitation events but not in smaller (areal) precipitation  
406 events, likely due to the trade-off between stable and convective precipitation in large events  
407 such as frontal systems and mesoscale clusters. This process may also explain the apparent  
408 increase in east-west areal coverage and the intensification of precipitation events, as found  
409 here.

410 Several generalizations of the proposed methodology are possible. In Marzban et al.  
411 (2008) it has been shown that clustering can be done not only in the 2-dimensional space of  
412 latitude and longitude of each grid point, but also in the 3-dimensional space that includes  
413 the amount of precipitation at each grid point. In fact, one may argue that the inclusion  
414 of more meteorological quantities in the clustering phase ought to lead to more meteorolog-  
415 ically relevant objects being identified. In turn, this is more likely to lead to more realistic  
416 representation of the effect of the parameters on the object features. The object features  
417 may also be extended or revised. For example, here the shape of an object is approximated  
418 by an ellipse. But it is possible to use more sophisticated methods of shape analysis (Book-  
419 stein 1991; Lack et al. 2010; Micheas et al. 2007; Lakshmanan et al. 2009) to model more  
420 complex shapes. Another possible generalization is to allow for interactions between model  
421 parameters. Although the statistical model used here does account for covariance between  
422 the model parameters, and between the response variables, no explicit interaction is intro-  
423 duced. The inclusion of such terms is straightforward, and is unlikely to lead to overfitting,  
424 at least in linear models such as MMR.



## 425 5. Code and/or data availability

426 The code and the data analyzed here occupy about 4.0G of computer space, and are avail-  
427 able upon request from the corresponding author, or from <https://doi.org/10.5281/zenodo.1043542>

## 428 6. Competing Interests

429 The authors declare that they have no conflict of interest

## 430 7. Acknowledgments

431 This work has received support from Office of Naval Research (N00014-12-G-0078 task  
432 29) and National Science Foundation (AGS-1402895). The authors are grateful to James D.  
433 Doyle and Nicholas C. Lederer for providing invaluable support.

434

## 435 REFERENCES

436 Ahijevych, D., Gilleland, D. E., Brown, B. G., and Ebert, E. E.: Application of spatial veri-  
437 fication methods to idealized and NWP-gridded precipitation forecasts, *Wea. Forecasting*,  
438 24, 1485–1497, 2009.

439 Backman, J., Wood, C., Auvinen, M., Kangas, L., Hannuniemi, H., Karppien, A., and



- 440 Kukkonen, J.: Sensitivity analysis of the meteorological pre-processor MPP-FMI 3.0 using  
441 algorithmic differentiation, *Geosci. Model Dev. Discuss.*, (in review), 2017.
- 442 Baldwin, M. E., Lakshmivarahan, S., and Kain, J. S.: Verification of mesoscale features in  
443 NWP models, in: *Amer Meteor. Soc., 9th Conf. on Mesoscale Processes*, pp. 255–258, Ft.  
444 Lauderdale, FL., 2001.
- 445 Baldwin, M. E., Lakshmivarahan, S., and Kain, J. S.: Development of an “events-oriented”  
446 approach to forecast verification, in: *15th Conf. Numerical Weather Prediction*, San An-  
447 tonio, TX, 2002.
- 448 Banfield, J. D. and Raftery, A. E.: Model-based Gaussian and non-Gaussian clustering,  
449 *Biometrics*, 49, 803–821, 1993.
- 450 Bolado-Lavin, R. and Badea, A. C.: Review of sensitivity analysis methods and experience  
451 for geological disposal of radioactive waste and spent nuclear fuel, *JRC Scientific and*  
452 *Technical Report*, Available online, 2008.
- 453 Bookstein, F. L.: *Morphometric Tools for Landmark Data: Geometry and Biology*, Cam-  
454 bridge, 1991.
- 455 Brown, B. G., Mahoney, J. L., Davis, C. A., Bullock, R., and Mueller, C.: Improved ap-  
456 proaches for measuring the quality of convective weather forecasts, in: *16th Conference*  
457 *on Probability and Statistics in the Atmospheric Sciences*, pp. 20–25, Orlando, FL, 2002.
- 458 Casati, B., Ross, G., and Stephenson, D.: A new intensity-scale approach for the verification  
459 of spatial precipitation forecasts, *Met. App.*, 11, 141–154, 2004.





- 460 Cioppa, T. and Lucas, T.: Efficient nearly orthogonal and space-filling latin hypercubes,  
461 Technometrics, 49(1), 45–55, 2007.
- 462 Davis, C., Brown, B., and Bullock, R.: Object-based verification of precipitation forecasts.  
463 Part I: Methodology and application to mesoscale rain areas, Mon. Wea. Rev., 134, 1772–  
464 1784, 2006a.
- 465 Davis, C. A., Brown, B., and Bullock, R.: Object-based verification of precipitation forecasts.  
466 Part II: Application to convective rain systems, Mon. Wea. Rev., 134, 1785–1795, 2006b.
- 467 DelSole, T. and Yang, X.: Field Significance of Regression Patterns, J. Climate, 24, 5094–  
468 5107, 2011.
- 469 Devore, J. and Farnum, N.: Applied Statistics for Engineers and Scientists, Thomson, 2005.
- 470 Doyle, J. D., Jiang, Q., Smith, R. B., and Grubii, V.: Three-dimensional characteristics of  
471 stratospheric mountain waves during T-REX, Mon. Wea. Rev., 139, 3–23, 2011.
- 472 Du, J. and Mullen, S. L.: Removal of Distortion Error from an Ensemble Forecast, Mon.  
473 Wea. Rev., 128, 3347–3351, 2000.
- 474 Ebert, E.: Fuzzy verification of high-resolution gridded forecasts: a review and proposed  
475 framework, Meteor. Appl., 15 (1), 51–64, 2008.
- 476 Ebert, E. E. and McBride, J. L.: Verification of precipitation in weather systems: determi-  
477 nation of systematic errors, Jour. Hydrology, 239, 179–202, 2000.
- 478 Errico, R. M.: What is an Adjoint Model?, Bull. Amer. Meteor. Soc, 78, 2577–2591, 1997.



- 479 Ester, M., Kriegel, H.-P., Sander, J., and Xu, X.: A density-based algorithm for discovering  
480 clusters in large spatial databases with noise, in: Proceedings of Knowledge Discovery and  
481 Data Mining (KDD-96), vol. 96(34), pp. 226–231, 1996.
- 482 Everitt, B. S.: Cluster Analysis, Heinemann Educational Books, London, 1980.
- 483 Fox, J., Friendly, M., and Weisberg, S.: Hypothesis tests for multivariate linear models using  
484 the car package, *The R Journal*, 5(1), 39–52, 2013.
- 485 Fraley, C. and Raftery, A.: Model-Based Clustering, Discriminant Analysis, and Density  
486 Estimation, *Journal of the American Statistical Association*, 97, 611–631, 2002.
- 487 Gilleland, D. E., Ahijevych, D., Brown, B. G., Casati, B., and Ebert, E. E.: Intercomparison  
488 of Spatial Forecast Verification Methods, *Wea. Forecasting*, 24, 1416–1430, 2009.
- 489 Hacker, J. P., Snyder, C., Ha, S.-Y., and Pocerlich, M.: Linear and non-linear response to  
490 parameter variations in a mesoscale model, *Tellus A*, 63, 429–444, doi:10.1111/j.1600-0870.  
491 2010.00505.x, 2011.
- 492 Hodur, R. M.: The Naval Research Laboratorys Coupled Ocean/Atmosphere Mesoscale  
493 Prediction System (COAMPS), *Mon. Wea. Rev.*, 125, 1414–1430, 1997.
- 494 Hoffman, R. N., Liu, Z., Louis, J.-F., and Grassotti, C.: Distortion representation of forecast  
495 errors, *Mon. Wea. Rev.*, 123, 2758–2770, 1995.
- 496 Holt, T. R., Cummings, J. A., Bishop, C. H., Doyle, J. D., Hong, X., Chen, S., and Jin, Y.:  
497 Development and testing of a coupled ocean-atmosphere mesoscale ensemble prediction  
498 system, *Ocean Dynamics*, 61, 1937–1954, doi:10.1007/s10236-011-0449-9, 2011.



- 499 Jiang, Q. and Doyle, J. D.: The impact of moisture on Mountain Waves, *Mon. Wea. Rev.*,  
500 137, 3888–3906, 2009.
- 501 Kalra, T. S., Aretxabaleta, A., Seshadri, P., Ganju, N. K., and Beudin, A.: Sensitivity  
502 Analysis of a Coupled Hydrodynamic-Vegetation Model Using the Effectively Subsampled  
503 Quadratures Method, *Geosci. Model Dev. Discuss.*, (in review), 2017.
- 504 Keil, C. and Craig, G. C.: A displacement-based error measure applied in a Regional En-  
505 semble Forecasting System, *Mon. Wea. Rev.*, 135(9), 3248–3259, 2007.
- 506 Lack, S. A., Limpert, G. L., and Fox, N. I.: An object-oriented multiscale verification scheme,  
507 *Wea. Forecasting*, 25(1), 79–92, 2010.
- 508 Laine, M., Solonen, A., Haario, H., and Järvinen, H.: Ensemble prediction and parameter  
509 estimation system: the method, *Q. J. R. Meteorol. Soc.*, 138, 289–297, 2012.
- 510 Lakshmanan, V. and Kain, J. S.: A Gaussian Mixture Model Approach to Forecast Verifi-  
511 cation, *Wea. Forecasting*, 25(3), 908–920, 2010.
- 512 Lakshmanan, V., Hondl, K., and Rabin, R.: An Efficient, general-purpose technique for  
513 identifying storm cells in geospatial image, *J. Atmos. Oceanic Technol.*, 26, 523–537, 2009.
- 514 Li, X., Sudarsanam, N., and Frey, D. D.: Regularities in data from factorial experiments,  
515 *Complexity*, 11(5), 32–45, 2006.
- 516 Lorenz, E. N.: Deterministic non-periodic flow, *J. Atmos. Sci.*, 20, 130–141, 1963.
- 517 Lucas, D. D., Klein, R., Tannahill, J., Ivanova, D., Brandon, S., Domyancic, D., and Zhang,



- 518 Y.: Failure analysis of parameter-induced simulation crashes in climate models, *Geosci.*  
519 *Model Dev.*, 6, 1157–1171, 2013.
- 520 Marzban, C.: Variance-based Sensitivity Analysis: An illustration on the Lorenz '63 model,  
521 *Mon. Wea. Rev.*, 141(11), 4069–4079, 2013.
- 522 Marzban, C. and Sandgathe, S.: Cluster analysis for verification of precipitation fields, *Wea.*  
523 *Forecasting*, 21(5), 824–838, 2006.
- 524 Marzban, C. and Sandgathe, S.: Cluster Analysis for Object-Oriented Verification of Fields:  
525 A Variation, *Mon. Wea. Rev.*, 136, 1013–1025, 2008.
- 526 Marzban, C., Sandgathe, S., and Lyons, H.: An Object-oriented Verification of Three NWP  
527 Model Formulations via Cluster Analysis: An objective and a subjective analysis, *Mon.*  
528 *Wea. Rev.*, 136 (9), 3392–3407, 2008.
- 529 Marzban, C., Sandgathe, S., Lyons, H., and Lederer, N.: Three Spatial Verification Tech-  
530 niques: Cluster Analysis, Variogram, and Optical Flow, *Wea. Forecasting*, 24(6), 1457–  
531 1471, 2009.
- 532 Marzban, C., Sandgathe, S., Doyle, J. D., and Lederer, N. C.: Variance-Based Sensitivity  
533 Analysis: Preliminary Results in COAMPS, *Mon. Wea. Rev.*, 142, 2028–2042, 2014.
- 534 McLachlan, G. J. and Peel, D.: *Finite Mixture Models*, John Wiley & Sons, Hoboken, NJ  
535 USA, 2000.
- 536 Micheas, A. C., Fox, N. I., Lack, S. A., and Wikle, C. K.: Cell identification and verification  
537 of QPF ensembles using shape analysis techniques, *J. Hydrology*, 343, 105–116, 2007.



- 538 Montgomery, D. C.: Design and Analysis of Experiments, Wiley & Sons, 7th edition, 2009.
- 539 Nachamkin, J. E.: Mesoscale verification using meteorological composites, *Mon. Wea. Rev.*,  
540 132, 941–955, 2004.
- 541 Ollinaho, P., ärvinen, H., Bauer, P., Laine, M., Bechtold, P., Susiluoto, J., and Haario, H.:  
542 Optimization of NWP model closure parameters using total energy norm of forecast error  
543 as a target, *Geosci. Model Dev.*, 7(5), 1889–1900, 2014.
- 544 Rencher, A. C. and Christensen, W. F.: *Methods of Multivariate Analysis*, John Wiley &  
545 Sons, Inc., Hoboken, NJ, USA, doi:10.1002/9781118391686.ch10, 2012.
- 546 Roberts, N. M. and Lean, H. W.: Scale-Selective Verification of Rainfall Accumulations from  
547 High-Resolution Forecasts of Convective Events, *Mon. Wea. Rev.*, 136, 78–97, 2008.
- 548 Robock, A., Luo, L., Wood, E. F., Wen, F., Mitchell, K. E., Houser, P., Schaake, J. C.,  
549 Lohmann, D., Cosgrove, B., Sheffield, J., Duan, Q., Higgins, R. W., Pinker, R. T., Tarpley,  
550 J. D., Basara, J. B., and Crawford, K. C.: Evaluation of the North American Land Data  
551 Assimilation System over the southern Great Plains during warm seasons, *J. Geophys.*  
552 *Res.*, 108, 8846–8867, 2003.
- 553 Roebber, P.: The role of surface heat and Moisture Fluxes Associated with large-scale ocean  
554 current meanders in maritime cyclogenesis, *Mon. Wea. Rev.*, 117, 1676–1694, 1989.
- 555 Roebber, P. and Bosart, L.: The sensitivity of precipitation to circulation details. part i: an  
556 analysis of regional analogs, *Mon. Wea. Rev.*, 126, 437–455, 1989.



- 557 Safta, C., Ricciuto, D., Sargsyan, K., Debusschere, B., Najm, H., Williams, M., and Thorn-  
558 ton, P.: Global sensitivity analysis, probabilistic calibration, and predictive assessment for  
559 the data assimilation linked ecosystem carbon model, *Geosci. Model Dev.*, 8, 1899–1918,  
560 2015.
- 561 Saltelli, A., Ratto, M., Andres, T., Campolongo, F., Cariboni, J., Saisana, M., and Tarantola,  
562 S.: *Global Sensitivity Analysis: The Primer*, Wiley Publishing, 2008.
- 563 Saltelli, A., Annoni, P., Azzini, I., Campolongo, F., Ratto, M., and Tarantola, S.: Variance  
564 based sensitivity analysis of model output: Design and estimator for the total sensitivity  
565 index, *Computer Physics Communications*, 181, 259–270, 2010.
- 566 Venugopal, V., Basu, S., and Foufoula-Georgiou, E.: A new metric for comparing precip-  
567 itation patterns with an application to ensemble forecasts, *J. Geophys. Res.*, 110, D8,  
568 D08 111 10.1029/2004JD005 395, 2005.
- 569 Wang, Y. H., Fan, C. R., Zhang, J., Niu, T., Zhang, S., and Jiang, J. R.: Forecast Verification  
570 and Visualization based on Gaussian Mixture Model Co-estimation, *Computer Graphics*  
571 *Forum*, 34, 99–110, 2015.
- 572 Wealands, S. R., Grayson, R. B., and Walker, J. P.: Quantitative comparison of spatial  
573 fields for hydrological model assessment: some promising approaches, *Advances in Water*  
574 *Resources*, 28, 15–32, 2005.
- 575 Wernli, H., Paulat, M., Hagen, M., and Frei, C.: SAL - A Novel Quality Measure for the  
576 Verification of Quantitative Precipitation Forecasts, *Mon. Wea. Rev.*, 136, 4470–4487,  
577 2008.



578 Yu, Y. Y., Finke, P. A., Wu, H. B., and Guo, Z. T.: Sensitivity analysis and calibration of a  
579 soil carbon model (SoilGen2) in two contrasting loess forest soils, Geosci. Model Dev., 6,  
580 29–44, 2013.



ID	Name (Unit)	Description	Default	Range
1	delt2KF ( $^{\circ}C$ )	Temperature increment at the LCL for KF trigger	0	-2, 2
2	cloudrad ( $m$ )	Cloud radius factor in KF	1500	500, 3000
3	prcpfrac	Fraction of available precipitation in KF, fed back to the grid scale	0.5	0, 1
4	mixlen	Linear factor that multiplies the mixing length within the PBL	1.0	0.5, 1.5
5	sfclfx	Linear factor that modifies the surface fluxes	1.0	0.5, 1.5
6	wfctKF	Linear factor for the vertical velocity (grid scale) used by KF trigger	1.0	0.5, 1.5
7	delt1KF ( $^{\circ}C$ )	Another method to perturb the temperature at the LCL in KF	0	-2, 2
8	autocon1 ( $\frac{kg}{m^3s}$ )	Autoconversion factors for the microphysics	0.001	1e-4, 1e-2
9	autocon2 ( $\frac{kg}{m^3s}$ )	Autoconversion factors for the microphysics	4e-4	4e-5, 4e-3
10	rainsi ( $\frac{1}{m}$ )	Microphysics slope intercept parameter for rain	8.0e6	8.0e5, 8.0e7
11	snowsi ( $\frac{1}{m}$ )	Microphysics slope intercept parameter for snow	2.0e7	2.0e6, 2.0e8

KF = Kain-Fritsch, PBL = Planetary Boundary Layer, LCL = Lifted Condensation Level

TABLE 1. The 11 parameters studied in this paper. Also shown are the default values, and the range over which they are varied.